

# BReSK: Bootstrapped Contrastive Representation Learning for Skeleton-Based Action Understanding

Anonymous CVPR submission

Paper ID 12

## Abstract

Self-supervised learning especially contrastive learning has emerged as a powerful paradigm for skeleton-based action recognition. However, existing approaches often rely on heavy architectural refinements, such as transformer-based modules, which introduce redundancy and increase model complexity without necessarily improving representation consistency. To address this issue, we propose BReSK, a self-supervised framework that combines a bootstrap prediction with momentum contrastive learning for skeleton-based action understanding. At the core of BReSK is DiP, an asymmetric dual-branch predictor that enforces cross-view consistency through spatial and temporal predictors in the Query branch, while using exponentially moving averaged (EMA) targets in the Key branch to stabilize representation learning. In addition, we introduce BoCL, a hybrid objective that jointly optimizes a bootstrap alignment loss ( $\mathcal{L}_{Crop}$ ) and a momentum-based contrastive loss ( $\mathcal{L}_{MiCo}$ ), improving instance discrimination while reducing class confusion in the embedding space. Extensive experiments on six benchmark datasets, including NTU-RGB+D 60/120, PKU-MMD, Toyota SmartHome, Penn Action, and Posetics, show that BReSK consistently outperforms state-of-the-art methods across diverse settings while using fewer parameters.

## 1. Introduction

Human Action Understanding (HAU) [1, 3, 6, 18, 25, 26, 30, 35] is fundamental for applications such as healthcare monitoring, daily activity analysis, and sports understanding. Among various modalities, skeleton-based action recognition is particularly effective due to its robustness to appearance variations and its ability to model structured human motion [3, 6, 18, 30]. Recent advances have been largely driven by self-supervised learning (SSL), especially contrastive learning methods [16, 19, 28]. These approaches learn discriminative representations by aligning augmented views of the same skeleton while separating different samples, enabling effective learning from large-scale unlabeled data and improving generalization over supervised methods [8, 12, 16, 19, 21, 28, 32].

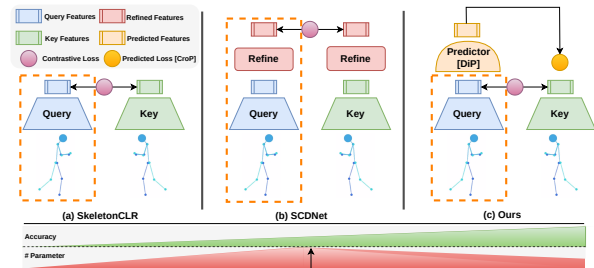


Figure 1. Illustration of the contrastive learning-based skeleton action recognition framework (top) and its accuracy-cost trade-off (bottom) compared to our proposed method **BReSK**.

Although these approaches have achieved promising performance, recent designs often rely on increasingly complex architectural components for feature refinement. For instance, SCDNet [28] introduces a Transformer-based refinement module that enhances feature representations through explicit self-correlation before projecting them into the contrastive space. While such strategies can improve representation quality, they also introduce redundant feature correlations and increase computational cost during both training and inference.

This limitation becomes more critical in real-world activity understanding scenarios. Compared with controlled laboratory datasets, real-world skeleton sequences often contain noisy joints, viewpoint variations, and irregular motion patterns, as illustrated in Fig. 1. Under these conditions, overly complex refinement modules may not necessarily lead to more robust representations, while significantly increasing computational overhead.

To address this limitation, we propose **BReSK**, a simple and efficient self-supervised framework for skeleton representation learning. Instead of explicitly refining features through complex architectural components, BReSK improves representation quality through implicit alignment in the learning objective. The framework combines a bootstrap prediction mechanism with contrastive learning, encouraging stable cross-view representation alignment while maintaining strong instance discrimination. This design strengthens positive pair consistency and reduces feature redundancy in the contrastive space.

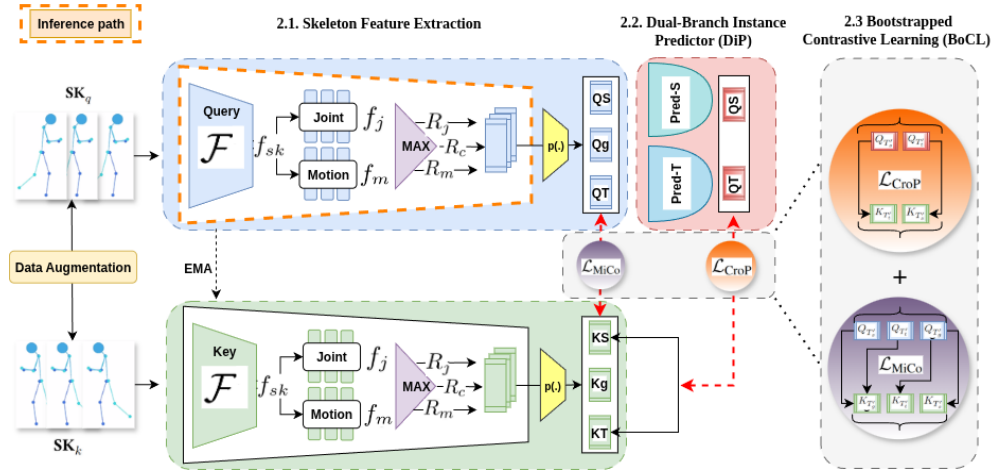


Figure 2. Overview of BReSK framework.

068 We evaluate BReSK on six benchmark datasets, i.e.,  
 069 Posetics, Toyota SmartHome PennAction, NTU-RGB+D  
 070 60/120 and PKU-MMD, covering diverse scenarios includ-  
 071 ing real-world sport and daily living activities. Across these  
 072 benchmarks, BReSK consistently outperforms recent self-  
 073 supervised methods. The learned representations remain  
 074 stable and transferable across datasets, enabling strong per-  
 075 formance in downstream tasks such as action recognition  
 076 and temporal action detection.

077 More importantly, BReSK achieves these improvements  
 078 while maintaining significantly lower computational cost.  
 079 Compared with recent methods relying on complex refine-  
 080 ment modules, our framework requires substantially fewer  
 081 parameters and lower FLOPs while preserving strong rep-  
 082 resentation capability. These results suggest that effective  
 083 skeleton representation learning can be achieved through  
 084 carefully designed learning objectives rather than increas-  
 085 ingly complex architectures.

086 In summary, our main contributions are: (i) We propose  
 087 **BReSK**, a simple and efficient self-supervised framework  
 088 that combines bootstrap representation alignment with con-  
 089 trastive learning for skeleton-based action understanding.  
 090 (ii) We introduce an asymmetric instance prediction mecha-  
 091 nism that enforces cross-view consistency using spatial and  
 092 temporal predictors while maintaining stable targets through  
 093 momentum updates. (iii) We design a hybrid learning objec-  
 094 tive that jointly optimizes bootstrap prediction alignment and  
 095 contrastive discrimination, improving representation quality  
 096 while reducing feature redundancy. (iv) Extensive experi-  
 097 ments on six benchmarks demonstrate that BReSK achieves  
 098 state-of-the-art performance while requiring significantly  
 099 lower computational cost.

## 100 2. Proposed Approach

101 In this section, we introduce **BreSK**, a self-supervised frame-  
 102 work that improves skeleton-based action understanding by  
 103 combines a bootstrap prediction with contrastive learning.

104 **Overview:** As shown in Fig. 2, the proposed framework  
 105 (BReSK) comprises three core components: (1) **Skeleton**  
 106 **Feature Extraction**; (2) **Dual-Branch Instance Predictor**;  
 107 and (3) **Bootstrapped Contrastive Learning**. Each  
 108 component is described in detail in the following sections.

### 109 2.1. Skeleton Feature Extraction:

110 The input skeleton sequence  $SK \in \mathbb{R}^{t \times J \times C}$ , where  $t$  is the  
 111 sequence length,  $J$  is the number of joints, and  $C$  denotes  
 112 joint coordinate dimensions (2D or 3D). Two augmented  
 113 views— $SK_q$  (query) and  $SK_k$  (key)—are generated using  
 114 rotation, flipping, and spatiotemporal masking [28] to im-  
 115 prove robustness. Both views are encoded using a CTR-GCN  
 116 backbone, producing  $f_{sk_q}$  and  $f_{sk_k}$ , respectively. Inspired  
 117 by MoCo-v2 [9], the query encoder is updated via backprop-  
 118 agation, while the key encoder is momentum-driven. Both  
 119 share identical structure but process different augmentations.

120 Encoder outputs contain tightly coupled spatiotemporal  
 121 features that are suboptimal for contrastive learning. To  
 122 disentangle these, we apply MLPs to separate the represen-  
 123 tations into motion view ( $f_m$ ) and joints view ( $f_j$ ) components:

$$f_m, f_j = \text{ReLU}(\text{LayerNorm}(\text{Linear}(f_{sk}))) \quad (1) \quad 124$$

125 These intermediate features are projected into the contrastive  
 126 space via a Max Pooling operation followed by a projection  
 127 layer in both the branches query and key. 128

### 129 2.2. Dual-Branch Instance Predictor (DiP)

130 **DiP** is an implicit lightweight feature refinement module de-  
 131 signed to enhance representation quality through cross-view  
 132 interaction. It adopts an asymmetric dual-branch architecture  
 133 to enforce cross-view consistency while avoiding representa-  
 134 tion collapse. 134

135 The query branch is equipped with lightweight predictor  
 136 modules, namely **Pred-S** and **Pred-T**, to model cross-view  
 137 relationships between spatial and temporal representations.  
 138 Let  $q_i \in \mathbb{R}^d$  denote the query representation from view  
 139  $i \in \{s, t\}$ , where  $s$  and  $t$  correspond to spatial and temporal 139

Methods	Linear Evaluation Posetics				Transfer Learning (Pre-trained on Posetics)						Semi-supervised							
	Top-1 (%)		Top-5 (%)		SmartHome CS		SmartHome CV2		PennAction		SmartHome CS				PennAction			
	2D	3D	2D	3D	CS (%)	CV2 (%)	CS (%)	CV2 (%)	Top-1 (%)	Top-1 (%)	5 (%)	10 (%)	5 (%)	10 (%)	5 (%)	10 (%)		
AimCLR [21] (AAAI 22)	19.2	-	39.3	-	46.6	-	48.3	-	-	-	-	-	-	-	-	-		
CMD [16] (ECCV 22)	20.4	20.4	40.5	40.7	49.0	50.8	52.5	-	89.4	87.4	-	-	-	-	-	-		
HiCLR [4] (AAAI 23)	20.1	-	39.9	-	49.1	-	52.3	-	88.7	-	-	-	-	-	-	-		
HiCo [7] (AAAI 23)	21.3	28.7	42.1	-	54.3	53.3	54.8	56.1	87.6	92.7	34.5	25.5	46.0	34.6	57.2	75.5	74.5	83.3
MAMP [17] (ICCV 23)	-	20.4	-	40.7	-	50.8	-	-	-	87.4	-	-	-	-	-	-	-	
PCM <sup>3</sup> [33] (ACMMM 23)	20.0	34.5	40.3	58.0	45.3	49.7	46.8	56.6	85.6	94.9	23.1	30.5	30.1	37.1	46.3	78.8	60.9	84.9
UmURL [20] (ACMMM 23)	-	32.7	-	56.5	-	48.1	-	56.9	-	94.8	-	22.7	-	30.6	-	78.1	-	85.9
SCD-Net [28] (AAAI 24)	27.2	-	50.6	-	54.6	-	55.5	-	90.8	-	-	-	-	-	-	-	-	-
ViA [32] (IJCV 24)	20.7	-	40.1	-	49.5	50.5	52.6	-	90.2	-	38.6	-	45.3	-	65.8	-	85.2	-
MacDiff [29] (ECCV 24)	-	23.3	-	42.1	-	51.1	-	56.2	-	91.9	-	-	-	-	-	-	-	-
USDRL [27] (AAAI 25)	25.9	34.6	48.6	58.5	55.0	49.6	53.9	56.7	89.6	95.2	37.1	34.3	43.6	40.3	58.1	73.9	70.3	86.4
<b>BReSK (Our)</b>	<b>29.8</b>	<b>36.9</b>	<b>52.7</b>	<b>61.0</b>	<b>56.5</b>	<b>53.9</b>	<b>56.9</b>	<b>57.8</b>	<b>91.9</b>	<b>95.8</b>	<b>39.7</b>	<b>37.4</b>	<b>48.3</b>	<b>46.5</b>	<b>69.5</b>	<b>80.4</b>	<b>86.1</b>	<b>89.7</b>

Table 1. **2D vs 3D Comparison** across Linear, Transfer, and Semi-supervised settings on Posetics, SmartHome, and PennAction.

Method	NTU-RGB+D 60		NTU-RGB+D 120	
	CS(%)	CV(%)	CS(%)	CS(%)
CrosSCLR [13](CVPR21)	72.9	79.9	-	-
HiCLR [4](AAAI23)	80.4	85.5	70.0	70.4
UmURL [20](ACMMM 23)	82.3	89.8	73.5	74.3
HiCo [7](AAAI 23)	81.1	88.6	72.8	74.1
ActCLR [14](CVPR 23)	80.9	86.7	69.0	70.5
ViA [32](IJCV 24)	78.1	85.8	69.2	66.9
SCD-Net [28] (AAAI 24)	<u>86.6</u>	91.7	76.9	<u>80.1</u>
ACA2Net [2](TCSVT 25)	86.0	89.6	-	-
ActCLR+ [15](TPAMI 25)	82.3	88.2	70.9	73.2
Heterogeneous [23] (CVPR 25)	80.2	88.0	70.7	73.5
USDRL [27] (AAAI 25)	84.2	90.8	76.0	76.9
USDRL+ [24] (TPAMI 25)	85.8	<u>91.8</u>	77.5	78.8
PCM <sup>3</sup> ++ [34] (IJCV 26)	84.8	91.0	76.7	-
<b>BReSK (Our)</b>	<b>87.3</b>	<b>92.0</b>	<b>79.1</b>	<b>80.1</b>

Table 2. **Action Recognition- Linear Evaluation results-** for Lab-setting datasets (NTU-RGB+D 60 and NTU-RGB+D 120)

streams, respectively. The predictor  $f_{i \rightarrow j}(\cdot)$  maps features from view  $i$  to the feature space of view  $j$ . Each predictor is implemented as a two-layer MLP:

$$f_{i \rightarrow j}(q_i) = W_2^{(i \rightarrow j)} \sigma \left( \text{LN} \left( W_1^{(i \rightarrow j)} q_i \right) \right), \quad (2)$$

where  $W_1 \in \mathbb{R}^{d_h \times d}$  and  $W_2 \in \mathbb{R}^{d \times d_h}$  are learnable parameters,  $\text{LN}(\cdot)$  denotes Layer Normalization, and  $\sigma(\cdot)$  is the ReLU activation function. In practice, we employ two cross-view predictors:

$$\hat{q}_s = f_{t \rightarrow s}(q_t), \quad \hat{q}_t = f_{s \rightarrow t}(q_s), \quad (3)$$

where  $\hat{q}_s$  and  $\hat{q}_t$  denote the predicted representations in the spatial and temporal spaces, respectively. By explicitly learning mappings across views, the proposed predictors encourage the model to capture complementary spatial-temporal semantics and improve representation alignment.

### 2.3. Bootstrapped Contrastive Learning (BoCL)

**BoCL** jointly optimizes two complementary objectives, i.e., a *Cross-Branch Prediction Alignment Loss* ( $\mathcal{L}_{\text{CroP}}$ ) that enforces cross-view semantic consistency, and a *Mixed Contrastive Loss* ( $\mathcal{L}_{\text{MiCo}}$ ) that sharpens inter-instance discrimination via a momentum-updated memory bank.

**Objective 1: Cross-Branch Prediction Alignment.** Our Dual-Branch Instance Predictor (DiP) learns to match the output of a momentum-updated teacher (key) network using stop-gradient targets and a predictor-based alignment. To enforce cross-view consistency, we apply a cosine similarity loss between predicted features and their corresponding target representations.

Method	mAP@tIoU (%)		
	0.1	0.3	0.5
AimCLR [21]	43.9	-	35.1
HiCo-Transformer [7]	32.5	31.8	28.6
HiCo-GRU [7]	50.1	48.6	44.3
SkeAttnCLR [10]	48.5	-	41.7
LAC [31]	55.2	-	58.5
SCDNet [28]	65.6	64.7	58.5
<b>BReSK (Our)</b>	<b>69.5</b>	<b>68.5</b>	<b>64.0</b>

Table 3. **Action Detection- Linear Evaluation results-** for Lab-setting datasets (PKU-MMD I).

BReSK consists of two branches: (i) a *query branch*, optimized via back-propagation and equipped with cross-view predictors, and (ii) a *key branch*, implemented as a momentum encoder updated using an exponential moving average (EMA) of the query parameters. The key branch provides stable targets and is detached from gradient computation.

Let  $k_j$  denote the target representation from view  $j$ , produced by the EMA encoder, and let  $\hat{q}_j = f_{i \rightarrow j}(q_i)$  be the predicted representation from the query branch. We normalize both vectors as:

$$\tilde{q}_j = \frac{\hat{q}_j}{\|\hat{q}_j\|_2}, \quad \tilde{k}_j = \frac{\text{sg}(k_j)}{\|k_j\|_2}, \quad (4)$$

where  $\text{sg}(\cdot)$  is the stop-gradient operator, ensuring gradients flow only through the query branch. The cross-view alignment loss is defined as:

$$\mathcal{L}_{i \rightarrow j} = \frac{1}{N} \sum_{n=1}^N \left( 1 - \tilde{q}_j^{(n)} \cdot \tilde{k}_j^{(n)} \right), \quad (5)$$

which minimizes the negative cosine similarity between predictions and targets. This encourages the query branch to align its representations with the EMA teacher across different views. Finally, the objective is symmetrized across both directions:

$$\mathcal{L}_{\text{CroP}} = \mathcal{L}_{t \rightarrow s} + \mathcal{L}_{s \rightarrow t}. \quad (6)$$

**Objective 2: Contrastive Learning.** This objective enhances discrimination in the embedding space by mixing spatial, temporal, and global views from different branches. Following [28], the MiCo loss is defined as:

$$\mathcal{L}_{\text{MiCo}} = \mathcal{L}_{\text{info}}(Q_j, K_c) + \mathcal{L}_{\text{info}}(Q_m, K_c) + \mathcal{L}_{\text{info}}(Q_c, K_j) + \mathcal{L}_{\text{info}}(Q_c, K_m), \quad (7)$$

CroP loss form	loss terms ( $Query \rightarrow Key$ )	acc
Instance	$(S \rightarrow S) + (T \rightarrow T) + (G \rightarrow G)$	86.6
Sub instance	$(S \rightarrow S) + (T \rightarrow T)$	86.4
Mixed	$(S \rightarrow G) + (T \rightarrow G)$	86.3
Cross	$(S \rightarrow T) + (T \rightarrow S)$	<b>87.3</b>

Table 4. Prediction loss forms

where  $Q$  and  $K$  represent query and key features, and  $\mathcal{L}_{\text{info}}$  is the InfoNCE loss [22].

**Final Loss.** The total self-supervised loss combines both objectives,  $\mathcal{L} = \mathcal{L}_{\text{CroP}} + \mathcal{L}_{\text{MiCo}}$ . This complementary strategy ensures that representations are internally consistent and externally discriminative, a key for robust generalization in downstream tasks.

### 3. Experiments and Analysis

We conduct extensive experiments comparing with state-of-the-art (SoTA) SSL methods to validate the effectiveness and generalization of BReSK across diverse settings. Evaluation protocols, additional implementation details, dataset descriptions, experimental details, and results are provided in the **Appendix**.

#### 3.1. Comparison with State-of-the-Art

**Action Recognition - Linear Evaluation.** As reported in Table 1 (blue), BReSK achieves SoTA performance across all real-world benchmarks, outperforming ViA [32] by **+9.1%** (Top-1) and **+12.6%** (Top-5) on **Posetics**. On controlled lab benchmarks (Table 2), BReSK consistently surpasses strong prior methods such as SCD-Net [28] on **NTU-RGB+D 60/120**, achieving a **+2.2%** improvement on NTU-120 under the Cross-Subject protocol using the joint modality.

**Action Recognition - Transfer Learning Evaluation.** We further evaluate representation transferability across datasets. In the real-world setting, we pre-train on **Posetics** and transfer to unseen datasets using **linear evaluation** (Table 1, green). BReSK outperforms prior state-of-the-art methods by **+1.5%** under the Cross-Subject protocol and **+1.4%** under Cross-View2 on **SmartHome**, demonstrating strong robustness to domain shifts and camera variations.

**Action Recognition - Semi-supervised Evaluation.** We evaluate BReSK under limited labeled data by fine-tuning with 5% and 10% of annotations on **SmartHome** and **Penn-Action** (Table 1, yellow). BReSK maintains strong performance under these low-label settings, demonstrating that slot-structured learning effectively captures transferable motion priors. Extending this evaluation to 3D skeletons (BReSK-3D), we observe consistent improvements across datasets, particularly for general actions in Posetics and sports actions in PennAction.

**Action Detection – Linear Evaluation.** We further evaluate BReSK on temporal localization using the untrimmed **PKU-MMD I** dataset. Following [28], the encoder is pre-trained on NTU-60 (Cross-Subject) and a linear classifier is trained to predict frame-level actions. As shown in Table 3, BReSK achieves significant mAP improvements over recent

Framework	Contrastive	DiP	CroP	NTU60-CS (%)	NTU120-CS (%)
Baseline	✓(NCE)	✗	✗	83.5	75.0
<b>BReSK (Ours)</b>	✓(NCE)	✓	✓	<b>85.8</b>	<b>76.7</b>
Baseline	✓(MiCo)	✗	✗	84.0	75.8
-	✓	✗	✓	86.7	76.9
-	✓	✓	✗	87.0	77.5
<b>BReSK (Ours)</b>	✓(MiCo)	✓	✓	<b>87.3</b>	<b>79.1</b>

Table 5. Ablation study of different components.

Model	Inference		Training			NTU120-CS (%)
	Params	GFlops	Params	Time (h)	Memory	
H-Transformer [5]	>100 M	118.6	—	—	—	—
GL-Transformer [11]	214 M	—	214 M	18	—	66.0
PCMB [33]	—	—	103 M	15	394.66 MB	76.5
SCDNet [28]	83.26 M	7.14	184 M	17	702.04 MB	76.9
USDRL [27]	94.86 M	—	95 M	14	362.13 MB	76.0
<b>BReSK (Ours)</b>	<b>22.03 M</b>	<b>3.75</b>	<b>83 M</b>	<b>11</b>	<b>318.74 MB</b>	<b>79.1</b>

Table 6. Computation cost on NTU-RGB+D 120 CS during inference(blue) and training process (green).

self-supervised baselines [7, 31], demonstrating that its slot-structured motion representation and unified discriminative reconstruction effectively capture temporal dynamics and enhance precise action localization in untrimmed sequences.

#### 3.2. Ablation Study

**Component-wise Analysis of BReSK:** Table 5 shows that both DiP and CroP consistently improve over the baseline across different contrastive objectives. While each component individually provides noticeable gains, their combination leads to the best performance, indicating a clear complementary effect. This is further validated under the stronger MiCo setting, where the full BReSK model achieves the highest accuracy.

**Effect of Cross-Space Feature Alignment:** Table 4 shows the importance of cross-space feature alignment in improving representation quality. While instance-level and sub-instance losses based on intra-space mappings (e.g.,  $S \rightarrow S$ ,  $T \rightarrow T$ ) yield comparable performance, they are consistently outperformed by formulations that encourage inter-space interactions. Notably, the cross loss ( $S \rightarrow T$ ,  $T \rightarrow S$ ) achieves the highest accuracy (87.3%), demonstrating that explicitly modeling complementary relationships between spatial and temporal features leads to more discriminative representations.

**Computation Cost:** Table 6 shows that BReSK achieves SOTA performance across all evaluation protocols, while simultaneously providing a significant reduction in computational cost, with up to  $\sim 2\times$  fewer inference parameters and over  $3\times$  lower GFLOPs compared to existing methods.

### 4. Conclusion

We presented BReSK, a self-supervised framework that leverages bootstrap prediction with contrastive learning for skeleton action recognition. Our design improves cross-view representation consistency while maintaining strong instance discrimination. Experiments on multiple benchmarks show that BReSK consistently outperforms existing methods across diverse settings with lower computational cost.

279  
280  
281  
282  
283  
284  
285  
286  
287  
288  
289  
290  
291  
292  
293  
294  
295  
296  
297  
298  
299  
300  
301  
302  
303  
304  
305  
306  
307  
308  
309  
310  
311  
312  
313  
314  
315  
316  
317  
318  
319  
320  
321  
322  
323  
324  
325  
326  
327  
328  
329  
330  
331  
332  
333  
334  
335  
336  
337**References**

- [1] Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lučić, and Cordelia Schmid. Vivit: A video vision transformer. *ICCV*, 2021. 1
- [2] Wenming Cao, Liangxi Qian, Yicha Zhang, Xuelong Li, and Xinpeng Yin. Asymmetric context-guided adaptive alignment network for skeleton-based action recognition. *IEEE Transactions on Circuits and Systems for Video Technology*, 2025. 3
- [3] Yuxin Chen, Ziqi Zhang, Chunfeng Yuan, Bing Li, Ying Deng, and Weiming Hu. Channel-wise topology refinement graph convolution for skeleton-based action recognition. In *ICCV*, 2021. 1
- [4] Yuxiao Chen, Long Zhao, Jianbo Yuan, Yu Tian, Zhaoyang Xia, Shijie Geng, Ligong Han, and Dimitris N. Metaxas. Hierarchically self-supervised transformer for human skeleton representation learning. In *ECCV*, 2022. 3
- [5] Yi-Bin Cheng, Xipeng Chen, Junhong Chen, Pengxu Wei, Dongyu Zhang, and Liang Lin. Hierarchical transformer: Unsupervised representation learning for skeleton-based human action recognition. In *2021 IEEE International Conference on Multimedia and Expo (ICME)*, 2021. 4
- [6] Jeonghyeok Do and Munchurl Kim. Skateformer: skeletal-temporal transformer for human action recognition. In *European Conference on Computer Vision*, 2024. 1
- [7] Jianfeng Dong, Shengkai Sun, Zhonglin Liu, Shujie Chen, Baolong Liu, and Xun Wang. Hierarchical contrast for unsupervised skeleton-based action representation learning. In *AAAI*, 2023. 3, 4
- [8] Christoph Feichtenhofer, Haoqi Fan, Bo Xiong, Ross Girshick, and Kaiming He. A large-scale study on unsupervised spatiotemporal representation learning. In *CVPR*, 2021. 1
- [9] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9729–9738, 2020. 2
- [10] Yilei Hua, Wenhan Wu, Ce Zheng, Aidong Lu, Mengyuan Liu, Chen Chen, and Shiqian Wu. Part aware contrastive learning for self-supervised action recognition. In *IJCAI*, 2023. 3
- [11] Boeun Kim, Hyung Jin Chang, Jungho Kim, and Jin Young Choi. Global-local motion transformer for unsupervised skeleton-based action learning. In *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part IV*, 2022. 4
- [12] Linguo Li, Minsi Wang, Bingbing Ni, Hang Wang, Jiancheng Yang, and Wenjun Zhang. 3d human action representation learning via cross-view consistency pursuit. In *CVPR*, 2021. 1
- [13] Linguo Li, Minsi Wang, Bingbing Ni, Hang Wang, Jiancheng Yang, and Wenjun Zhang. 3d human action representation learning via cross-view consistency pursuit. In *(CVPR)*, 2021. 3
- [14] Lilang Lin, Jiahang Zhang, and Jiaying Liu. Actionlet-dependent contrastive learning for unsupervised skeleton-based action recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 3
- [15] Lilang Lin, Jiahang Zhang, and Jiaying Liu. Self-supervised skeleton representation learning via actionlet contrast and reconstruct. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2025. 3
- [16] Yunyao Mao, Wengang Zhou, Zhenbo Lu, Jiajun Deng, and Houqiang Li. Cmd: Self-supervised 3d action representation learning with cross-modal mutual distillation. In *ECCV*, 2022. 1, 3
- [17] Yunyao Mao, Jiajun Deng, Wengang Zhou, Yao Fang, Wanli Ouyang, and Houqiang Li. Masked motion predictors are strong 3d action representation learners. In *ICCV*, 2023. 3
- [18] Lei Shi, Yifan Zhang, Jian Cheng, and Hanqing Lu. Two-stream adaptive graph convolutional networks for skeleton-based action recognition. In *CVPR*, 2019. 1
- [19] Chen Sun, Arsha Nagrani, Yonglong Tian, and Cordelia Schmid. Composable augmentation encoding for video representation learning. In *ICCV*, 2021. 1
- [20] Shengkai Sun, Daizong Liu, Jianfeng Dong, Xiaoye Qu, Junyu Gao, Xun Yang, Xun Wang, and Meng Wang. Unified multi-modal unsupervised representation learning for skeleton-based action understanding. In *Proceedings of the 31st ACM International Conference on Multimedia*, pages 2973–2984, 2023. 3
- [21] Guo Tianyu, Liu Hong, Chen Zhan, Liu Mengyuan, Wang Tao, and Ding Runwei. Contrastive learning from extremely augmented skeleton sequences for self-supervised action recognition. In *AAAI*, 2022. 1, 3
- [22] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding, 2019. 4
- [23] Hongsong Wang, Xiaoyan Ma, Jidong Kuang, and Jie Gui. Heterogeneous skeleton-based action representation learning. In *(CVPR)*, 2025. 3
- [24] Hongsong Wang, Wanjiang Weng, Junbo Wang, Fang Zhao, Guo sen Xie, Xin Geng, and Liang Wang. Foundation model for skeleton-based human action understanding. *Transactions on Pattern Analysis and Machine Intelligence*, 2025. 3
- [25] Limin Wang, Bingkun Huang, Zhiyu Zhao, Zhan Tong, Yinan He, Yi Wang, Yali Wang, and Yu Qiao. Videomae v2: Scaling video masked autoencoders with dual masking. In *CVPR*, 2023. 1
- [26] Yi Wang, Yinan He, Yizhuo Li, Kunchang Li, Jiashuo Yu, Xin Ma, Xinyuan Chen, Yaohui Wang, Ping Luo, Ziwei Liu, Yali Wang, Limin Wang, and Yu Qiao. Internvid: A large-scale video-text dataset for multimodal understanding and generation. In *ICLR*, 2024. 1
- [27] Wanjiang Weng, Hongsong Wang, Junbo Wang, Lei He, and Guosen Xie. Usdr1: Unified skeleton-based dense representation learning with multi-grained feature decorrelation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2025. 3, 4
- [28] Cong Wu, Xiao-Jun Wu, Josef Kittler, Tianyang Xu, Sara Aitio, Muhammad Awais, and Zhenhua Feng. Scd-net: Spatiotemporal clues disentanglement network for self-supervised skeleton-based action recognition. In *AAAI*, 2024. 1, 2, 3, 4
- [29] Lehong Wu, Lilang Lin, Jiahang Zhang, Yiyang Ma, and Jiaying Liu. Macdiff: Unified skeleton modeling with masked conditional diffusion. In *European Conference on Computer Vision (ECCV)*, 2024. 3

- 398 [30] Di Yang, Yaohui Wang, Antitza Dantcheva, Lorenzo Garat-  
399 toni, Gianpiero Francesca, and Francois Bremond. Unik:  
400 A unified framework for real-world skeleton-based action  
401 recognition. In *BMVC*, 2021. 1
- 402 [31] Di Yang, Yaohui Wang, Antitza Dantcheva, Quan Kong,  
403 Lorenzo Garattoni, Gianpiero Francesca, and Francois Bre-  
404 mond. Lac - latent action composition for skeleton-based  
405 action segmentation. In *ICCV*, 2023. 3, 4
- 406 [32] Di Yang, Yaohui Wang, Antitza Dantcheva, Lorenzo Garat-  
407 toni, Gianpiero Francesca, and Francois Bremond. Via: View-  
408 invariant skeleton action representation learning via motion  
409 retargeting. *IJCV*, 2024. 1, 3, 4
- 410 [33] Jiahang Zhang, Lilang Lin, and Jiaying Liu. Prompted con-  
411 trast with masked motion modeling: Towards versatile 3d  
412 action representation learning. In *Proceedings of the ACM*  
413 *International Conference on Multimedia*, 2023. 3, 4
- 414 [34] Jiahang Zhang, Lilang Lin, Shuai Yang, and Jiaying Liu. Self-  
415 supervised skeleton-based action representation learning: A  
416 benchmark and beyond. *IJCV*, 2026. 3
- 417 [35] Wentao Zhu, Xiaoxuan Ma, Zhaoyang Liu, Libin Liu, Wayne  
418 Wu, and Yizhou Wang. Motionbert: A unified perspective on  
419 learning human motion representations. In *ICCV*, 2023. 1