

(APPENDIX) BReSK: Bootstrapped Contrastive Representation Learning for Skeleton-Based Action Understanding

Anonymous CVPR submission

Paper ID 12

001 Appendix Overview

002 The Supplementary material is organized as follows:

- 003 • Section 1: Datasets and Evaluation Metrics
- 004 • Section 2: Evaluation Protocols
- 005 • Section 3: Framework and Implementation Details
- 006 • Section 4: Additional SOTA Comparisons
 - 007 1. Action recognition (Multi-Modality) -Linear Eval-
 - 008 uation
 - 009 2. Action Retrieval - KNN Evaluation
 - 010 3. Action Prediction Evaluation
- 011 • Section 5: Additional Ablation Studies
 - 012 1. Disentanglement process
 - 013 2. Deeper Insights, Analysis, and Comparison
 - 014 3. Representation visualization using t-SNE

015 1. Datasets and Evaluation Metrics

016 **NTU-RGB+D 60 [17]** consists of 56,880 sequences of
017 high-quality 3D skeletons, captured by the Microsoft Kinect
018 v2 sensor. We only use sequences of 3D skeletons in this
019 work, and we follow the cross-subject (CS) and cross-view
020 (CV) evaluation protocol.

021 **NTU-RGB+D 120 [13]** is an extension of NTU-RGB+D
022 60, increasing the number of action classes from 60 to 120.
023 The dataset contains a total of 114,480 video samples. We
024 use 3D skeleton sequences and follow the Cross-Subject
025 (CS) and Cross-Setup (CSet) evaluation protocols.

026 **PKU-MMD [3]** : For 3D skeleton action classification on
027 PKU-MMD, action instances are segmented based on tempo-
028 ral annotations and are partitioned into training and testing
029 sets according to the cross-subject evaluation protocol. In
030 our experiments, we use PKU-MMD II, with 5,332 samples
031 for training and 1,613 samples for testing.

Posetics [28] is built upon Kinetics400 videos, the dataset
consists of 142,000 real-world video clips across 320 ac-
tion classes, along with corresponding 2D and 3D skele-
tons. We leverage the Posetics dataset to pre-train our action
representation learning framework using skeleton data and
explore transfer learning for skeleton-based action classifi-
cation. Top-1 and Top-5 accuracy are used as evaluation
metrics.

Toyota Smarthome [4] (SmartHome) is a real-world col-
lection for daily living action classification, containing
16,115 videos across 31 action classes. It includes RGB
videos, as well as 2D and 3D skeleton data [27]. Since
the 2D data is more robust for action recognition, even in
cross-view evaluations, we use 2D data for the experiments
unless otherwise specified. For evaluation, we report the
mean per-class accuracy following the cross-subject (CS)
and cross-view (CV2) protocols.

Penn Action [32] contains 2,326 video sequences of 15
different real-world sports actions. In this work, we use 2D
skeletons obtained by LCRNet++ for our experiments, and
we report Top-1 accuracy following the standard train-test
splitting.

054 2. Evaluation Protocols

We evaluate BReSK under four complementary self-
supervised learning protocols following prior works [1, 12,
24, 25, 29, 31]. (i) We conduct **linear evaluation** to measure
the intrinsic discriminability of the learned representations
using a frozen encoder. (ii) We perform **cross-dataset trans-**
fer by pretraining on a source dataset and adapting to unseen
target datasets to assess generalization. (iii) We study **semi-**
supervised learning by fine-tuning with only 5% or 10%
labeled samples to evaluate data efficiency. (iv) We evaluate
temporal reasoning through **action prediction** from partial
observations (70%, 80%, 90% of the sequence). (iv) We
examine **action detection** in untrimmed sequences using
mAP at multiple temporal IoU thresholds (0.1, 0.3, 0.5). To-
gether, these protocols assess representation discriminability,

Category	Parameter	Lab setting & Real-world
Optimizer	Optimizer type	SGD
	Momentum	0.9
DiP Settings	num-layers	1
	Prediction-dim	128
Contrastive Settings	Queue size	8192
	Momentum	0.999
	Softmax temperature	0.2
	Projection dimension	128
Training Setup	Weight decay	0.0001
	Initial learning rate	0.01
	Learning rate schedule	[350]
	Batch size	64
	λ	[100 - 1]
	Training steps	440

Table 1. Training settings for self-supervised learning on the lab setting and real-world (NTU-RGB+D and Posetics dataset).

Parameter	Action Classification		Action Retrieval	Transfer Learning		Semi-Supervised		Action Detection
	Lab	Real	& Action Prediction	Lab	Real	Lab	Real	Lab
Initial learning rate	2		2	0.2	2	0.2		0.01
Batch size	1024		1024	32	1024	32	16	32
Training steps	80		80		80	80	100	60

Table 2. Evaluation hyperparameter across different tasks and settings (lab and real-world) for various datasets.

069 data efficiency, cross-domain generalization, early action
070 understanding, and temporal localization ability.

071 3. Framework and Implementation Details

072 Following prior state-of-the-art works [5, 25], all input skele-
073 ton sequences are pre-processed to a fixed length of 64
074 frames. We then apply the same data augmentation strat-
075 egy as in [25] to generate two different views of the input,
076 serving as the query and key sequences.

077 The query and key branches share the same architecture
078 and functionality following MoCo-v2 [6]. Below, we de-
079 scribe the query branch in detail.

080 We adopt a lightweight version of CTR-GCN [2] com-
081 posed of 3 blocks as our backbone for skeleton feature ex-
082 traction, with the output channel dimension fixed at 64. The
083 extracted features are then separated into motion view (fm)
084 and joints view (fj) components after reshaping and are
085 projected into a higher semantic space.

086 These features are followed by Dual-Branch Instance
087 Predictor **DiP**, which is an implicit lightweight feature re-
088 finement module designed to enhance representation quality
089 through cross-view interaction. It adopts an asymmetric
090 dual-branch architecture to enforce cross-view consistency
091 while avoiding representation collapse.

092 The query branch is equipped with lightweight predictor
093 modules, namely **Pred-S** and **Pred-T**, to model cross-view

relationships between spatial and temporal representations. 094
Let $q_i \in \mathbb{R}^d$ denote the query representation from view 095
 $i \in \{s, t\}$, where s and t correspond to spatial and temporal 096
streams, respectively. The predictor $f_{i \rightarrow j}(\cdot)$ maps features 097
from view i to the feature space of view j . Each predictor is 098
implemented as a two-layer MLP. 099

Finally, Bootstrapped Contrastive Learning (BoCL) loss 100
is applied to enhance the representation by jointly optimiz- 101
ing two complementary objectives, i.e., a *Cross-Branch Pre-* 102
diction Alignment Loss ($\mathcal{L}_{\text{CROP}}$) that enforces cross-view se- 103
mantic consistency, and a *Mixed Contrastive Loss* ($\mathcal{L}_{\text{MiCo}}$) 104
that sharpens inter-instance discrimination via a momentum- 105
updated memory bank. 106

The MoCo parameters are $m = 0.999$ for the EMA up- 107
date parameter, and $k = 8192$ for the queue length. We 108
pretrain our model with SGD optimizer for 440 epochs with 109
a batch size of 64, and we pretrain our framework on one 110
Nvidia H100 GPU. For the evaluation and downstream task, 111
we use the trained query branch encoder under different pro- 112
tocol evaluations. All hyperparameters used during training 113
and evaluation are detailed in Tables 1 and 2, respectively. 114

115 4. Additional SOTA Comparisons

In this section, we provide additional state-of-the-art (SOTA) 116
comparisons and evaluations across both real-world and lab- 117
setting (controlled) datasets. For the lab-setting datasets, we 118

Method	NTU-RGB+D 60 CS(%)		NTU-RGB+D 60 CV(%)		NTU-RGB+D 120 CS(%)		NTU-RGB+D 120 CSett(%)	
	J	J+B+M	J	J+B+M	J	J+B+M	J	J+B+M
CrosSCLR [10] (CVPR 21)	72.9	77.8	79.9	83.4	-	67.9	-	66.7
CMD [14] (ECCV 22)	79.8	84.1	86.9	90.9	70.3	74.7	71.5	76.1
ActCLR [11] (CVPR 23)	80.9	84.3	86.7	88.8	69.0	74.3	70.5	75.7
AimCLR [20] (AAAI 22)	74.3	78.9	79.7	83.8	63.4	68.2	63.4	68.8
RVTCLR+ [33] (ICCV 23)	74.7	79.7	79.1	84.6	-	60.0	-	68.9
UmURL [18] (ACMMM 23)	82.3	84.2	89.8	90.0	73.5	75.2	74.3	76.3
PCM ³ [30] (ACMMM 23)	83.9	87.4	90.4	93.1	76.5	80.0	77.5	81.2
HiCo [5] (AAAI 23)	81.1	83.8	88.6	90.8	72.8	75.9	74.1	77.3
SCD-Net [25] (AAAI 24)	86.6	87.3	91.7	-	76.9	-	80.1	-
USDRL [24] (AAAI 25)	84.2	87.1	90.8	93.2	76.0	79.3	76.9	80.6
BReSK (Our)	87.3	88.2	92.0	93.2	79.1	81.8	80.1	82.5

Table 3. **Action Recognition (Multi-Modality)- Linear Evaluation-** for Lab-setting datasets (NTU-RGB+D 60 and NTU-RGB+D 120)

Methods	Action Retrieval	
	Posetics Top-1 (%)	PennAction Top-1 (%)
HiCo [5](AAAI 23)	9.6	70.8
PCM ³ [30] (ACMMM 23)	7.6	73.7
ViA [29] (IJCV 24)	-	76.9
USDRL [24](AAAI 25)	12.7	76.7
BReSK-2D (Our)	13.3	78.4
HiCo [5](AAAI 23)	15.9	86.8
USDRL [24](AAAI 25)	20.2	89.8
BReSK-3D (Our)	20.8	90.2

Table 4. **Action Recognition 2D and 3D skeleton** – Combined results for *Linear Evaluation*, *Action Retrieval* and *Transfer Learning* on Real-world datasets (Posetics, SmartHome, and PennAction).

Method	70%	80%	90%	100%
DeepSCN [9]	58.2	60.2	60.0	58.6
MSRNN [7]	63.9	67.4	68.9	69.2
P-TSL [23]	77.6	80.1	81.5	82.0
SCDNet [25]	79.5	82.1	85.0	86.6
USDRL (STTR) [24]	80.3	82.4	83.7	84.2
USDRL (DSTE) [24]	81.4	83.6	84.5	85.2
BReSK-3D (Ours)	82.0	84.2	86.1	87.3

Table 5. **Action Prediction- Linear Evaluation results-** for Lab-setting datasets (NTU-RGB+D60 CS).

Method	70%	80%	90%	100%
<i>Penn Action</i>				
HiCo [5]	77.8	80.6	83.4	87.6
PCM ³ [30]	53.4	56.7	62.4	85.7
USDRL [24]	85.7	86.9	87.9	89.6
BReSK-2D (Our)	86.7	89.2	90.8	91.9
<i>Posetics</i>				
HiCo [5]	18.5	19.9	20.5	21.3
PCM ³ [30]	17.1	18.1	18.7	20.0
USDRL [24]	24.2	24.8	25.4	25.9
BReSK-2D (Our)	27.6	28.5	29.2	29.8

Table 6. **Action Prediction- Linear Evaluation results-** on Real-world datasets (Penn Action and Posetics).

Method	Transfer to PKU-MMD II	
	NTU-RGB+D 120 CS(%)	
HiCo-Transformer [5]	55.4	
CrosSCLR-B [10]	52.8	
CMD [14]	57.0	
UmURL-3 [18]	58.5	
A ² MC [26]	58.9	
USDRL+ [22]	58.3	
Heterogeneous [21]	63.1	
BReSK-3D (Our)	64.2	

Table 7. **Action Recognition- Transfer Learning Evaluation (Fine-tune) results-** for Lab-setting datasets (from NTU-RGB+D 120 to PKU-MMD II).

119 present additional experiments, including multimodal linear
 120 evaluation and action retrieval. For real-world scenarios, we
 121 adopt recent SOTA methods [5, 24] on 2D and 3D skeleton
 122 sequences to ensure a fair comparison with our approach. We
 123 evaluate action retrieval after being pretrained on Posetics
 124 dataset.

125 4.1. Action recognition (Multi-Modality) 126 -Linear Evaluation:

127 Skeleton sequence can be represented in different modalities:
 128 joint-based (J), focusing on individual joint positions to cap-
 129 ture static postures; bone-based (B), encoding the relative

positions between connected joints to highlight body part
 interactions; and motion-based (M), capturing the changes
 in joint or bone positions over time to emphasize dynamic
 movement patterns. These modalities offer complementary
 insights into both static and dynamic aspects of human movement.
 In Table 3, we report the linear evaluation results compared to SoTA on lab-setting datasets (NTU-RGB+D 60 and NTU-RGB+D 120) using joint-based and multi-modality (J+B+M). For the multi-modality case, we first train our model separately for each modality and then fuse the results from all three modalities. BReSK shows a significant improvement compared to the recent SoTA [24, 25, 30], across

Method	NTU-RGB+D 60		NTU-RGB+D 120	
	CS (%)	CV (%)	CS (%)	Cset (%)
ISC [19]	62.5	82.6	50.6	53.0
CrosSCLR-B [10]	66.1	81.3	52.2	54.9
HaLP [16]	65.8	83.6	55.8	59.0
SkeAttnCLR [8]	69.4	67.8	46.7	58.0
HiCo [5]	68.3	84.8	56.6	59.1
MAMP [15]	62.0	70.0	51.8	56.1
A ² MC [26]	<u>70.8</u>	<u>85.4</u>	<u>59.1</u>	<u>62.6</u>
Heterogeneous [21]	66.3	87.1	55.7	59.8
BReSK-3D (Our)	75.2	87.2	61.8	64.8

Table 8. **Action retrieval- KNN Evaluation results-** for Lab-setting datasets (NTU-RGB+D 60 & NTU-RGB+D 120).

Activity	Gain from BReSK-2D
Usetablet	+64.67
Walk	+29.23
Cutbread	+21.82
Eat.Snack	+19.43
Drink.Fromcan	+17.05
Pour.Frombottle	+10.83
Usetaptop	+9.54
Drink.Fromcup	-1.05
Sitdown	-2.00
Drink.Frombottle	-3.95
Pour.Fromcan	-6.66
Mean Accuracy	+4.3

Table 9. Gain from BReSK compared to ViA on SmartHome-CV2 dataset. The results are Mean per class.

Activity	Gain from BReSK-2D
Drink.Fromcup	+50.91
Leave	+47.23
Readbook	+45.09
Pour.Fromkettle	+36.99
Usetaptop	+31.41
Enter	+26.28
Usetelephone	+22.20
Pour.Fromcan	-2.30
Drink.Fromglass	-2.30
Maketea.Boilwater	-2.30
Drink.Frombottle	-4.57
Makecoffee.Pourgrains	-6.05
Mean Accuracy	+7.0

Table 10. Gain from BReSK compared to ViA [29] on SmartHome-CS dataset. The results are Mean per class.

across all evaluation protocols. While prior approaches such as HiCo [5] and A²MC demonstrate strong performance under specific settings, their results are less consistent across different protocols. In contrast, BReSK maintains high accuracy across all evaluations, indicating that the learned representations are both robust and transferable. These results highlight the effectiveness of BReSK in learning discriminative feature embeddings for skeleton-based action retrieval. As well as demonstrating strong generalization capability and effective cross-dataset transfer of the learned representations.

all benchmark datasets with different types of modality.

4.2. Action Prediction - Linear Evaluation

Tables 5 and 6 report BReSK’s performance on action prediction, demonstrating its ability to anticipate actions from partial observations. Thanks to the design of BReSK, which allows for capturing robust spatio-temporal dependencies, enabling accurate early recognition even when only a fraction of the skeleton sequence is available. This highlights the model’s strength in temporal reasoning and its potential for real-time applications, where predicting ongoing actions from limited observations is critical.

4.3. Action Retrieval - KNN Evaluation

Table 4 and Table 8 report the action retrieval performance using a K-Nearest Neighbor (KNN) evaluation on (Posetics, PennAction, NTU-RGB+D 60, and NTU-RGB+D 120). We use a K-Nearest Neighbor (KNN) classifier, a non-parametric supervised learning method, which directly reflects the quality of the representation space learned by the frozen encoder. BReSK consistently outperforms recent state-of-the-art methods, including Heterogeneous [21] and A²MC [26],

5. Additional Ablation Studies

5.1. Disentanglement process (Spatial and Temporal)

The decoupling process is considered an essential step and has proven effective in many SOTA. While USDRL [24] relies on a dense spatio-temporal encoder and mask strategy applied directly to the projected embeddings from reshaped skeleton sequences, and ViA [29] leverages orthogonal basis vectors to disentangle motion in the feature space, BReSK applies lightweight projection layers directly to intermediate features after reshaping from a motion (time) and joint perspective. This design makes the features more effective and adaptable in our framework.

Specifically, given a skeleton sequence $\mathbf{SK} \in \mathbb{R}^{C \times T \times V \times M}$, with C 2D/3D, temporal length T and V skeletal joints, intermediate features $\mathbf{fsk} \in \mathbb{R}^{C_{\text{inter}} \times T \times V}$ are first extracted. These features are then transformed into two distinct views (motion and joint) by reshaping them into $\mathbf{f}_m \in \mathbb{R}^{T \times (V \times C_{\text{in}})}$ for the motion domain and $\mathbf{f}_j \in \mathbb{R}^{V \times (T \times C_{\text{in}})}$ for the joint domain. The rearranged sequences are subsequently mapped into a shared embedding space using separate MLPs (Eq. 1), producing feature representations $\mathbf{F}_m \in \mathbb{R}^{T \times C_e}$ and $\mathbf{F}_j \in \mathbb{R}^{V \times C_e}$. This decou-

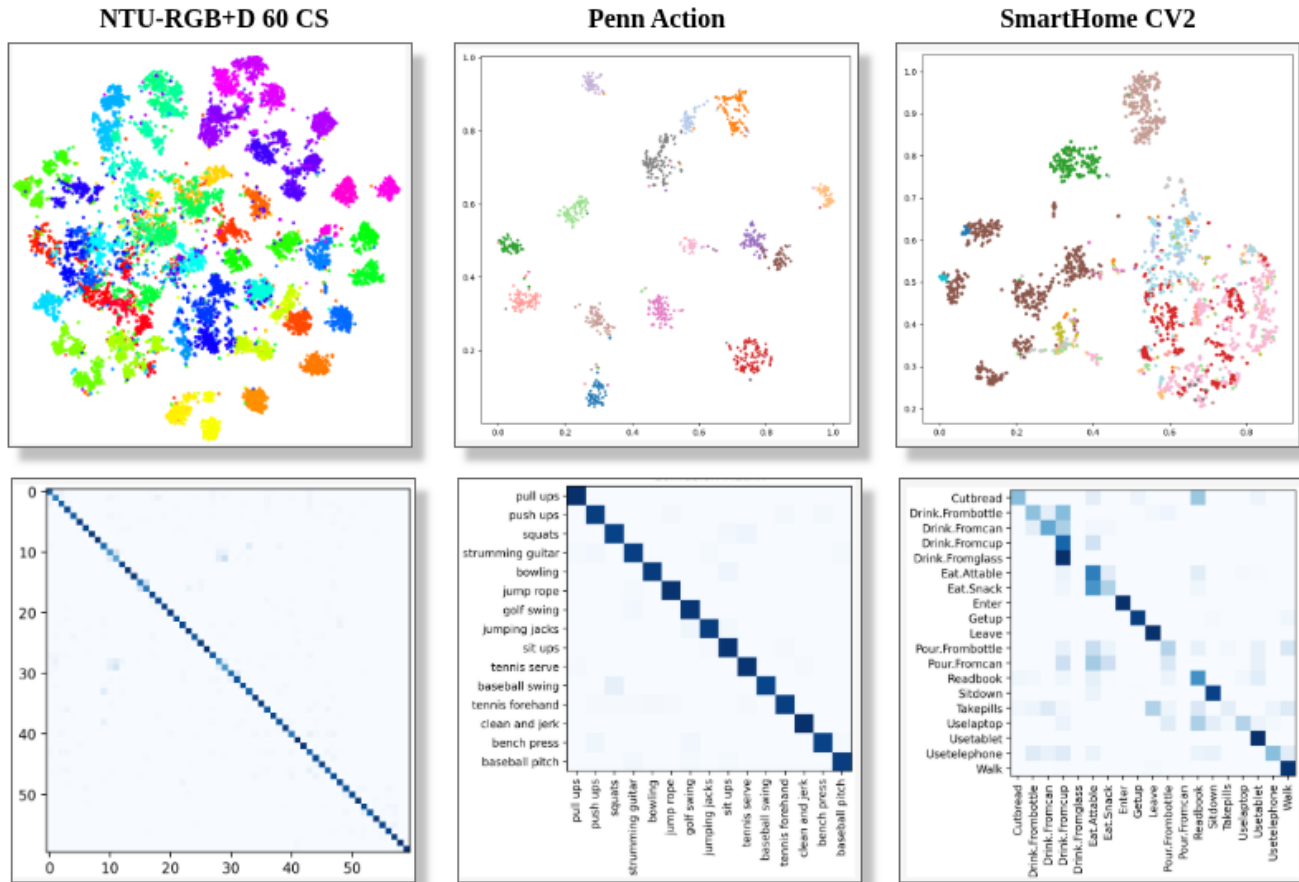


Figure 1. Confusion matrix and TSNE plot of BReSK features NTU-RGB+D 60, Penn CS datasets.

195 pling enables actions to be represented from complementary
196 perspectives.

197 5.2. Deeper Insights, Analysis and Comparison

198 The class-wise evaluation across SmartHome-CS and
199 SmartHome-CV2 provides deeper insights into the inductive
200 biases of BReSK and its ability to generalize under cross-
201 subject and cross-view conditions.

202 Table 9 shows that BReSK-2D improves over ViA with a
203 +4.3 mean accuracy gain on the **SmartHome-CV2 dataset**,
204 indicating overall better performance. The largest improve-
205 ments are observed for dynamic and subtle actions, such as
206 Usetable (+64.67), Walk (+29.23), and Cutbread (+21.82),
207 as well as Eat.Snack and Drink.Fromcan. This suggests
208 that BReSK-2D is effective at capturing fine-grained motion
209 patterns and temporal dynamics from skeleton data.

210 In contrast, slight performance drops are observed for
211 a few classes (e.g., Drink.Fromcup, Pour.Fromcan), likely
212 due to high similarity between motion patterns and limited
213 discriminative cues in skeleton representations.

214 on **SmartHome-CS** dataset, Table 10 shows that BReSK-
215 2D achieves a larger overall improvement, with a +7.0

mean accuracy gain, confirming the consistency of its
216 advantage over ViA. Similar to the previous results, the
217 strongest gains appear in subtle and dynamic actions, such
218 as Drink.Fromcup (+50.91), Leave (+47.23), Readbook
219 (+45.09), and Pour.Fromkettle (+36.99). Activities like Use-
220 laptop, Enter, and Usetelephone also show notable improve-
221 ments, highlighting the ability of BReSK-2D to capture
222 fine-grained motion dynamics from skeleton data.
223

224 However, performance drops are observed for several
225 classes (e.g., Pour.Fromcan, Drink.Fromglass, Drink.Frombottle,
226 Makecoffee.Pourgrains). These activities are object-centric
227 and involve similar interaction patterns, where distinguishing
228 the action depends heavily on object information. Since skeleton
229 data does not encode object cues, these cases remain challenging.
230

231 Taken together, these results confirm that BReSK-2D
232 consistently excels at capturing subtle and dynamic motion
233 patterns from skeleton data, leading to strong improvements
234 across diverse activities. At the same time, its limitations
235 are evident in object-centric actions, where the absence of
236 explicit object information makes it difficult to distinguish
237 between similar interaction patterns, highlighting an inherent

238 challenge of skeleton-based representations.

239 5.3. Representation visualization using t-SNE

240 To further analyze the quality of the representations learned
241 by BReSK, we conduct a qualitative and quantitative evalua-
242 tion across three benchmark datasets: NTU-RGB+D 60 CS,
243 Penn Action, and SmartHome CV2. Figure 1 illustrates the
244 t-SNE projections of the extracted features alongside their
245 corresponding confusion matrices, providing insight into the
246 discriminative capacity of our model under varying dataset
247 complexities.

248 We first assess the structure of the learned embedding
249 space using t-SNE projections (Figure 1, top row). On
250 NTU-RGB+D 60 CS, which comprises 60 action categories,
251 BReSK produces a large number of visually distinct clusters,
252 demonstrating its capacity to maintain inter-class separability
253 even under high class-count conditions. On Penn Action, the
254 15 action classes form tightly grouped and clearly separated
255 clusters, indicating that BReSK captures highly consistent
256 intra-class features while preserving strong inter-class mar-
257 gins. On the more challenging SmartHome CV2 dataset,
258 where activities are semantically fine-grained and visually
259 similar, partial cluster overlaps are observed — particularly
260 among actions sharing similar objects or motion patterns.
261 This behavior is expected and reflects the inherent difficulty
262 of fine-grained activity recognition rather than a limitation
263 of the learned representations.

264 The confusion matrices (Figure 1, bottom row) comple-
265 ment the t-SNE analysis by quantifying per-class prediction
266 accuracy. Across all three datasets, a dominant diagonal
267 structure is consistently observed, confirming that BReSK
268 correctly classifies the majority of samples for each action
269 class. On NTU-RGB+D 60 CS, the diagonal remains strong
270 across all 60 classes with minimal off-diagonal activations.
271 On Penn Action, the near-perfect diagonal corroborates the
272 clean cluster separation seen in the embedding space. On
273 SmartHome CV2, localized off-diagonal entries appear be-
274 tween semantically related classes, such as Drink.Frombottle
275 and Drink.Fromcan, which are distinguishable only by subtle
276 object-level differences, highlighting the challenging nature
277 of this benchmark.

278 The consistency between the t-SNE visualizations and the
279 confusion matrices across all three benchmarks confirms that
280 BReSK learns compact, discriminative, and generalizable
281 feature representations. The progressive increase in classifi-
282 cation complexity from Penn Action to NTU-RGB+D 60 CS
283 to SmartHome CV2 is clearly reflected in both the embed-
284 ding structure and the confusion patterns, validating that our
285 model scales gracefully with dataset difficulty. These results
286 collectively support the design choices underlying BReSK
287 and justify its effectiveness as a robust skeleton-based action
288 recognition framework.

References

- [1] 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. 1
- [2] 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 *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 13359–13368, 2021. 2
- [3] Liu Chunhui, Hu Yueyu, Li Yanghao, Song Sijie, and Liu Jiaying. Pku-mmd: A large scale benchmark for continuous multi-modal human action understanding. *arXiv:1703.07475*, 2017. 1
- [4] Srijan Das, Rui Dai, Michal Koperski, Luca Minciullo, Lorenzo Garattoni, Francois Bremond, and Gianpiero Francesca. Toyota smarthome: Real-world activities of daily living. In *ICCV*, 2019. 1
- [5] 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. 2, 3, 4
- [6] 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
- [7] Jian-Fang Hu, Wei-Shi Zheng, Lianyang Ma, Gang Wang, Jianhuang Lai, and Jianguo Zhang. Early action prediction by soft regression. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2019. 3
- [8] 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. 4
- [9] Yu Kong, Zhiqiang Tao, and Yun Fu. Deep sequential context networks for action prediction. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 3
- [10] 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, 4
- [11] 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
- [12] 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. 1
- [13] J. Liu, A. Shahroudy, M. Perez, G. Wang, L. Y. Duan, and A. C. Kot. Ntu rgb+d 120: A large-scale benchmark for 3D human activity understanding. *IEEE TPAMI*, 2020. 1
- [14] 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. 3

- 348 [15] Yunyao Mao, Jiajun Deng, Wengang Zhou, Yao Fang, Wanli
349 Ouyang, and Houqiang Li. Masked motion predictors are
350 strong 3d action representation learners. In *ICCV*, 2023. 4
- 351 [16] Anshul Shah, Aniket Roy, Ketul Shah, Shlok Kumar
352 Mishra, David Jacobs, Anoop Cherian, and Rama Chellappa.
353 Halp: Hallucinating latent positives for skeleton-based self-
354 supervised learning of actions, 2023. 4
- 355 [17] Amir Shahroudy, Jun Liu, Tian-Tsong Ng, and Gang Wang.
356 Ntu rgb+d: A large scale dataset for 3D human activity analy-
357 sis. *CVPR*, 2016. 1
- 358 [18] Shengkai Sun, Daizong Liu, Jianfeng Dong, Xiaoye Qu,
359 Junyu Gao, Xun Yang, Xun Wang, and Meng Wang. Uni-
360 fied multi-modal unsupervised representation learning for
361 skeleton-based action understanding. In *Proceedings of the*
362 *31st ACM International Conference on Multimedia*, pages
363 2973–2984, 2023. 3
- 364 [19] Fida Mohammad Thoker, Hazel Doughty, and Cees G. M.
365 Snoek. Skeleton-contrastive 3d action representation learning,
366 2021. 4
- 367 [20] Guo Tianyu, Liu Hong, Chen Zhan, Liu Mengyuan, Wang
368 Tao, and Ding Runwei. Contrastive learning from extremely
369 augmented skeleton sequences for self-supervised action
370 recognition. In *AAAI*, 2022. 3
- 371 [21] Hongsong Wang, Xiaoyan Ma, Jidong Kuang, and Jie Gui.
372 Heterogeneous skeleton-based action representation learning.
373 In (*CVPR*), 2025. 3, 4
- 374 [22] Hongsong Wang, Wanjiang Weng, Junbo Wang, Fang Zhao,
375 Guo sen Xie, Xin Geng, and Liang Wang. Foundation model
376 for skeleton-based human action understanding. *Transactions*
377 *on Pattern Analysis and Machine Intelligence*, 2025. 3
- 378 [23] Xionghui Wang, Jian-Fang Hu, Jian-Huang Lai, Jianguo
379 Zhang, and Wei-Shi Zheng. Progressive teacher-student learn-
380 ing for early action prediction. In *2019 IEEE/CVF Conference*
381 *on Computer Vision and Pattern Recognition (CVPR)*, 2019.
382 3
- 383 [24] Wanjiang Weng, Hongsong Wang, Junbo Wang, Lei He, and
384 Guosen Xie. Usdrl: Unified skeleton-based dense represen-
385 tation learning with multi-grained feature decorrelation. In
386 *Proceedings of the AAAI Conference on Artificial Intelligence*,
387 2025. 1, 3, 4
- 388 [25] Cong Wu, Xiao-Jun Wu, Josef Kittler, Tianyang Xu, Sara
389 Atito, Muhammad Awais, and Zhenhua Feng. Scd-net:
390 Spatiotemporal clues disentanglement network for self-
391 supervised skeleton-based action recognition. In *AAAI*, 2024.
392 1, 2, 3
- 393 [26] Binqian Xu, Xiangbo Shu, Jiachao Zhang, Rui Yan, and Guo-
394 Sen Xie. Attack-augmentation mixing-contrastive skeletal
395 representation learning, 2024. 3, 4
- 396 [27] Di Yang, Rui Dai, Yaohui Wang, Rupayan Mallick, Luca Min-
397 ciullo, Gianpiero Francesca, and Francois Bremond. Selective
398 spatio-temporal aggregation based pose refinement system:
399 Towards understanding human activities in real-world videos.
400 In *WACV*, 2021. 1
- 401 [28] Di Yang, Yaohui Wang, Antitza Dantcheva, Lorenzo Garat-
402 toni, Gianpiero Francesca, and Francois Bremond. Unik:
403 A unified framework for real-world skeleton-based action
404 recognition. In *BMVC*, 2021. 1
- 405 [29] Di Yang, Yaohui Wang, Antitza Dantcheva, Lorenzo Garat-
406 toni, Gianpiero Francesca, and Francois Bremond. Via: View-
invariant skeleton action representation learning via motion
retargeting. *IJCV*, 2024. 1, 3, 4
- [30] Jiahang Zhang, Lilang Lin, and Jiaying Liu. Prompted con-
trast with masked motion modeling: Towards versatile 3d
action representation learning. In *Proceedings of the ACM*
International Conference on Multimedia, 2023. 3
- [31] Jiahang Zhang, Lilang Lin, Shuai Yang, and Jiaying Liu. Self-
supervised skeleton-based action representation learning: A
benchmark and beyond. *IJCV*, 2026. 1
- [32] W. Zhang, M. Zhu, and K. G. Derpanis. From actemes to ac-
tion: A strongly-supervised representation for detailed action
understanding. In *ICCV*, 2013. 1
- [33] Yisheng Zhu, Hu Han, Zhengtao Yu, and Guangcan Liu. Mod-
eling the relative visual tempo for self-supervised skeleton-
based action recognition. In *Proceedings of the IEEE/CVF*
International Conference on Computer Vision (ICCV), 2023.
3