

SkeletonMAE: Spatial-Temporal Masked Autoencoders for Self-supervised Skeleton Action Recognition

Wenhan Wu¹, Yilei Hua², Ce Zheng³, Shiqian Wu², Chen Chen³, Aidong Lu¹

¹Department of Computer Science, University of North Carolina at Charlotte, USA

²School of Information Science and Engineering, Wuhan University of Science and Technology, China

³Center for Research in Computer Vision, University of Central Florida, USA

{wwu25, alu1}@uncc.edu; hy19797510@gmail.com

cezhen@knights.ucf.edu; shiqian.wu@wust.edu.cn; chen.chen@crcv.ucf.edu

Abstract

Fully supervised skeleton-based action recognition has achieved great progress with the blooming of deep learning techniques. However, these methods require sufficient labeled data which is not easy to obtain. In contrast, self-supervised skeleton-based action recognition has attracted more attention. With utilizing the unlabeled data, more generalizable features can be learned to alleviate the overfitting problem and reduce the demand of massive labeled training data. Inspired by the MAE [15], we propose a spatial-temporal masked autoencoder framework for self-supervised 3D skeleton-based action recognition (SkeletonMAE). Following MAE’s masking and reconstruction pipeline, we utilize a skeleton based encoder-decoder transformer architecture to reconstruct the masked skeleton sequences. A novel masking strategy, named Spatial-Temporal Masking, is introduced in terms of both joint-level and frame-level for the skeleton sequence. This pre-training strategy makes the encoder output generalizable skeleton features with spatial and temporal dependencies. Given the unmasked skeleton sequence, the encoder is fine-tuned for the action recognition task. Extensive experiments show that our SkeletonMAE achieves remarkable performance and outperforms the state-of-the-art methods on both NTU RGB+D and NTU RGB+D 120 datasets.

1. Introduction

Human Action Recognition is a fundamental research topic in computer vision, which aims to understand human behaviors and distinguish the actions [38]. With the booming development of deep learning and human pose estimation methods [2, 39, 37], human skeleton data can be efficiently extracted as a high-level but light-weighted representation, which draws great attention for human behavior

and action analysis. Thus, 3D skeleton-based action recognition has become an important research field in human action recognition.

Most recent methods focus on full-supervised learning algorithms to build their frameworks: methods based on Convolutional Neural Network(CNN) [11, 48], methods based on Recurrent Neural Networks (RNN) [43, 34], methods based on Graph Convolution Networks (GCN) [46, 33, 49, 6] and methods based on Transformer [29, 30] are widely applied in skeleton action recognition and lead to very good results. However, fully supervised action recognition is liable to overfitting. Also, it requires massive labeled training data, which is expensive and time-consuming. To alleviate these issues, *self-supervised learning* methods, which utilize unlabeled data to learn data representations, have been increasingly prevalent in skeleton action recognition. Some self-supervised approaches consider pretext tasks for skeleton representation learning using unlabeled skeleton data, such as motion reconstruction [7] and jigsaw puzzle [22]. However, such pretext-based methods focus on local features such as joint correlation and skeleton scale in the same frame, and have not fully explored the temporal information. Recently, several works [20, 14] train the contrastive-based model based on contrastive learning framework through constructing the skeleton sequences in different views by data augmentation and positive-negative pairs. Although these contrastive learning based methods emphasize high-level context information, they heavily rely on the number of the contrastive pairs in the joints for extracting skeleton features, and ignore the joint correlation information among different frames.

Recently, a new self-supervised learning approach named masked autoencoders (MAE) [15] demonstrates a strong generalization capability with remarkable performance in computer vision tasks. MAE masks a large proportion of the input image, and then forces the model

to learn a generalizable representation by using only the unmasked proportion to reconstruct the original image. However, MAE [15] can not be directly utilized for self-supervised skeleton action recognition due to the following reasons:

- The Vision Transformer (ViT) [10] architecture is used in MAE [15] to process the image input. Different from the image that does not contain temporal information, human skeleton sequences are extracted from videos with high information density, which contains fruitful semantic information: at the spatial level, joint features contain the relationships among different joints in the same frame; in temporal level, frame features represent the movements of the same joint from different frames.
- The masking strategy in MAE only focuses on the spatial domain. When processing the human skeleton sequences data, a spatial-temporal masking strategy is needed.

To address these issues, we introduce a novel skeleton-based masked autoencoder named **SkeletonMAE** for self-supervised skeleton spatial-temporal representation learning: 1) the masked input sequences are generated from the original skeleton sequences, which contain joints coordinates (spatial) information and frames (temporal) information; 2) with spatial-temporal masking strategy and encoding-decoding rule, SkeletonMAE gains reconstruction sequences from masked sequences, where the spatial and temporal information is well processed by transformer-based encoder and decoder (transformers have great potential for spatial-temporal representation learning with long-term sequence data).

The framework of SkeletonMAE is presented in Fig. 1. Specifically, the whole SkeletonMAE pipeline is designed with the following principals. During the pre-training stage, a spatial-temporal masking strategy (with pre-set frame-masking and joint masking ratios) is employed to mask out part of the input skeleton sequence in both frame-level and the joint-level (Sec. 3.2). In order to find the best trade-off point for spatial-temporal representation learning, we discuss the roles of joint-masking and frame-masking ratios and find the best ratio combination. The encoder is applied to learn the generalizable feature representation while the decoder is designed to reconstruct the missing skeletons. Since we are dealing with the skeleton sequences, we utilize Spatio-Temporal Tuples Transformer (STTFormer) [30], which is developed for processing skeleton sequences, as our network backbone instead of ViT [10]. During the fine-tuning stage, we only use the encoder with a simple output layer to predict the actions. The action recognition results show that our approach outperforms the state-of-the-art self-supervised learning methods without extra data. To summarize, we make the following contributions:

1. We propose a simple and efficient skeleton-based masked autoencoder architecture, which aims to learn comprehensive and generalizable skeleton feature representations.
2. To have a better understanding of the skeleton masking methods, we explore different masking methods and develop a novel spatial-temporal masking for skeleton data in both joint-level and frame-level. At the same time, we validate the proper combination of joint-masking ratio and frame-masking ratio.
3. We evaluate our model on NTU-RGB+D 60 and NTU-RGB+D 120 datasets, and extensive experimental results show that SkeletonMAE achieves state-of-the-art performance under self-supervised settings.

2. Related work

2.1. Supervised skeleton-based action recognition

In the pre-deep learning period, hand-craft techniques are used for extracting spatial-temporal features in skeleton-based action recognition works [44, 40, 41]. In recent years, deep learning has been widely used in skeleton action recognition fields due to its powerful ability of feature extraction and representation learning. And most of them are fully supervised. RNN-based methods (*e.g.*, LSTMs) [12, 50, 24] were widely utilized to process skeleton data. Meanwhile, CNN-based methods [35, 18, 45] were also introduced to skeleton action recognition. Nevertheless, the data representations extracted by RNNs or CNNs were too simple to present the comprehensive spatial-temporal features of skeleton data. Thus, GCN-based methods [46, 33, 49, 6] were naturally introduced to model the topological graph features from skeleton data. Recently, with the success of vision transformer (ViT) [10], transformer-based model becomes powerful architecture for sequential skeleton data analysis [29, 30, 52, 47, 5, 51, 27] due to the ability of learning global representations. Therefore, we adopt the skeleton-based transformer (STTFormer [30]) as the backbone network in our research for a better skeleton sequences processing.

2.2. Self-supervised skeleton-based action recognition

Self-supervised learning aims to extract feature representations without using labeled data, and achieves promising performance in image-based and video-based representation learning [36, 26, 9, 13]. More self-supervised representation learning approaches adopt the so-called contrastive learning manner [16, 4, 21, 3] to boost their performance. Inspired by contrastive learning architectures, recent skeleton representation learning works have achieved

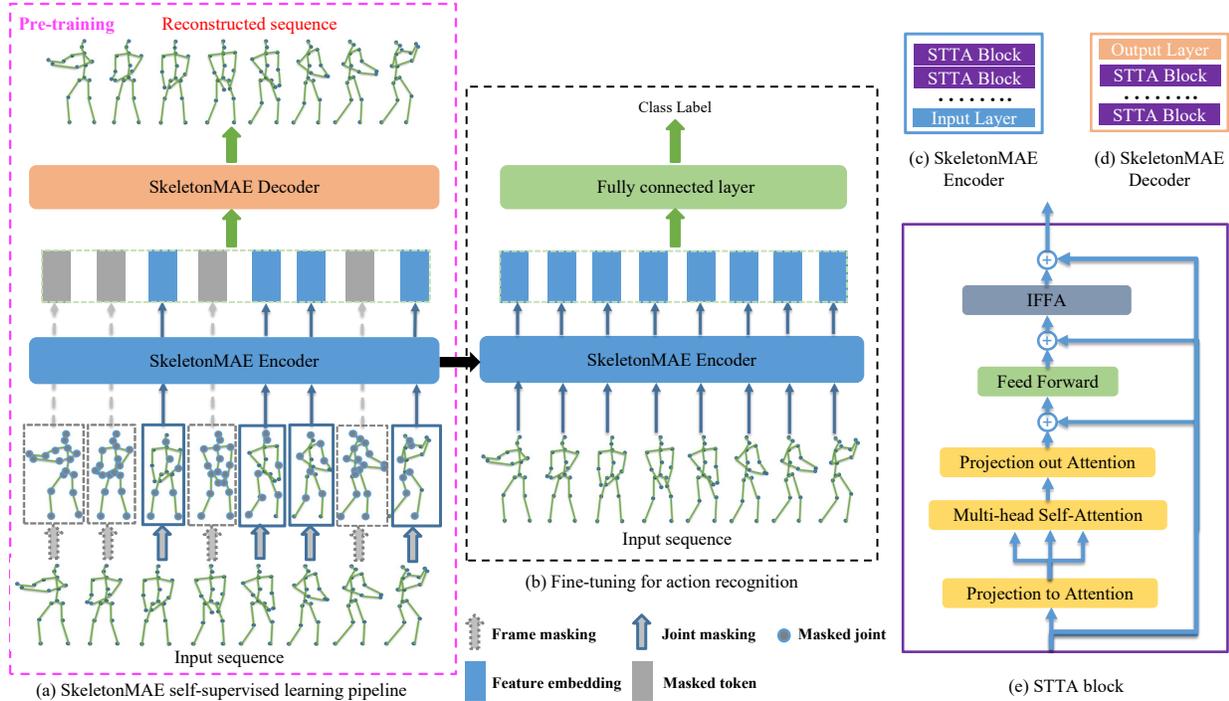


Figure 1: (a) The overall pipeline of the SkeletonMAE. During the pre-training, we utilize STTFormer to build our encoder and decoder, which consists several STTA blocks respectively. Then we only use SkeletonMAE encoder during the fine-tuning. (b) The end-to-end fine-tuning procedure for skeleton action recognition. (c) The STTFormer-based encoder structure, which is constructed by several Spatio-Temporal Tuples Attention (STTA) blocks and an input layer. (d) The STTFormer-based decoder structure, which is built by a series of STTA blocks with an output layer. (e) The structure of STTA block.

some inspiring progress in self-supervised skeleton action recognition. MS²L [22] introduced a multi-task self-supervised learning framework for extracting joints representations by using motion prediction and jigsaw puzzle recognition. CrosSCLR [20] developed a contrastive learning-based framework to learn both single-view and across-view representations from skeleton data. Following CrosSCLR, AimCLR [14] exploited an extreme data augmentation strategy to add extra hard contrastive pairs, which aims to learn more general representations from skeleton data.

2.3. Masked autoencoding

Masked autoencoding [42] is a well-structured self-supervised learning model for general representation learning, and successfully applied in BERT [8], one of the most famous self-supervised frameworks in natural language processing (NLP). The BERT model is simple and straightforward – remove part of the sequence data with the masked tokens, and predict the removed parts and calculate the loss between prediction and ground-truth data. As a result, the reconstruction sequence works well for training of the generalizable models [25, 31, 1]. Inspired by masked autoencoders and BERT, He et al. [15]. design a scal-

able self-supervised masked autoencoder(MAE) for computer vision task. With the same core concept as BERT, MAE masks parts of the image patches and rebuilds them for pre-training. Comparing with the original MAE, there are two main spotlights in our proposed SkeletonMAE: 1) a skeleton-based transformer encoder-decoder framework, the encoder processes the unmasked tokens and decoder reconstructs the original skeleton sequence; 2) a spatial-temporal masking strategy for both joint and frame level features. Following the main idea of MAE, we propose SkeletonMAE for self-supervised skeleton action recognition.

3. Methodology

In this section, we first introduce the preliminaries of SkeletonMAE in Sec. 3.1. Then, we design a spatial-temporal masking strategy for skeleton data in Sec. 3.2. Next, we analyze our SkeletonMAE for action recognition in Sec. 3.3. Finally, we present our fine-tuning procedure in Sec. 3.4.

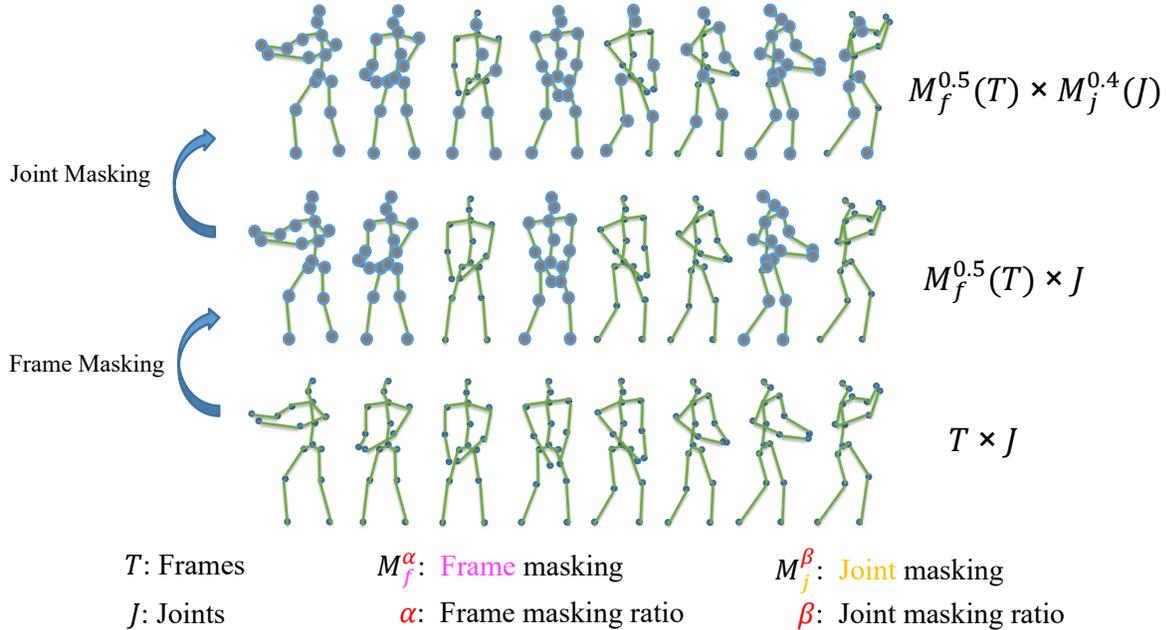


Figure 2: Illustration of the spatial-temporal masking pipeline. Based on pre-set frame-masking ratio (α) and joint-masking ratio (β), we first adopt frame masking (*i.e.* removing an entire skeleton frame) in skeleton sequence (*e.g.*, $\alpha = 0.5$), and randomly mask the joints in joint-level (*e.g.*, $\beta = 0.4$).

3.1. Preliminaries

MAE [15]. MAE is formed by an encoder and a decoder in a asymmetric way. It should be noticed that the structure of decoder is different from the encoder, which means that we can adapt some customized decoders to construct the efficient pre-training model. Specifically, the encoder in MAE is based on ViT yet only processing unmasked images: following ViT, the image patches are encoded by linear projection and added with positional embedding to be image tokens, then the tokens are processed by several transformer blocks. Only a small size of the unmasked tokens (75% patches are masked and the rest are set as the input) are loaded by encoder. As for the MAE decoder, it decodes the masked tokens with the position information based on the original image patches for reconstruction. Then the mean squared error (MSE) is calculated between masked and reconstructed tokens in pixel space. After pre-training, the pre-trained encoder with a simple classification head is applied for image classification task.

STTFormer [30]. Different from MAE which applies ViT in encoder and transformer blocks in decoder for image reconstruction, we take STTFormer to build encoder and decoder due to its skeleton-based transformer structure. Comparing with ViT which is based on image patches without temporal information, STTFormer is a skeleton data-driven transformer and shows great po-

tential in spatial-temporal data processing. Specifically, STTFormer divides skeleton data into several tuples (non-overlapping parts), and provides the self-attention module named Spatio-Temporal Tuples Attention (STTA) to extract multi-joint representations among adjacent frames. Then a feature aggregation module named Inter-Frame Feature Aggregation (IFFA) is proposed for inter-frame action integration after STTA block, improving the learning ability for similar action recognition. The structure of STTFormer is shown in Fig. 1.

3.2. Spatial-temporal masking strategy

We propose a spatial-temporal masking method to a portion of the the skeleton sequence input, the pipeline of our masking strategy is illustrated in Fig. 2.

Temporal-masking method. Fig. 2 shows our masking method at the frame-level. Based on the pre-set frame-masking ratio, a portion of the frames are randomly removed and their indices are stored, the remaining frames are then processed by spatial-masking method at the joint-level.

Spatial-masking method. As shown in Fig. 2, after implementing temporal masking method in all the input frames, the rest frames are then processed via spatial masking strategy. And based on the pre-set joint-masking ratio, we randomly mask part of the joints in every unmasked frame. It is worth noting that the indices

of the masked joints are not fixed in this randomly spatial-masking method, which means that the same joints in different frames may be masked or not. This simple approach is illustrated in Fig. 3(b). Besides this masking method, we also introduce a joint masking strategy with fixed indices, which is shown in Fig. 3(c). The joints with the same indices in different frames are all masked or not based on the joint-masking ratio. We conduct experiments to compare these two masking strategies in Sec. 4.3.

3.3. SkeletonMAE architecture

We describe the main components in SkeletonMAE, *e.g.*, encoder, decoder, reconstruction sequence, loss function and fine-tuning pipeline for skeleton action recognition. The pipeline and SkeletonMAE structure are illustrated in Fig. 1.

SkeletonMAE encoder. Our encoder is based on STTFormer and only processes the visible skeleton tokens. Given a skeleton sequence as input, we apply the frame-masking and joint-masking methods respectively. This spatially and temporally unmasked token is fed to the SkeletonMAE encoder, which maps the input to the spatial-temporal embedding features.

SkeletonMAE decoder. Our decoder also adopts STTFormer structure. Same as the decoder in MAE, the spatial-temporal embedding features is processed in SkeletonMAE decoder to reconstruct the original sequence. At the same time, in order to reserve the position information for reconstruction, positional embeddings are also introduced. The output of the decoder is the reconstructed sequence, which should be the same as the original sequence without masking.

Reconstruction. We use the mean squared error (MSE) loss to measure the consequence of reconstruction. In this case, we compute the MSE loss between original skeleton sequences and the reconstructed sequences as follows:

$$MSE = \frac{1}{N} \sum_{i=1}^N |S_i - S_i^*|^2, \quad (1)$$

where i is the index of frame, N is the number of samples, S is the input sequence, and S^* is the reconstructed sequence.

3.4. Fine-tuning for skeleton action recognition

In order to evaluate SkeletonMAE’s ability of learning skeleton representations, we load the learned parameter weights obtained from pre-training to fine-tune the model with all the training data, then the label for each action is predicted with the recognition accuracy. The procedure of fine-tuning is shown in Fig. 1 (b). Different from the latest contrastive-based self-supervised skeleton action recognition methods[20, 14], which verify the model via linear evaluation protocol, we only focus on the end-to-end fine-tuning for the skeleton action recognition tasks.

4. Experiments

4.1. Datasets

We evaluate our experiments on the following two most-used datasets: NTU-RGB+D 60 dataset [32] and NTU-RGB+D 120 dataset [23], and follow the according evaluation protocols for the experimental evaluation.

NTU-RGB+D 60 (NTU-60). NTU-60 is a large scale skeleton dataset for human skeleton-based action recognition, which contains 56,578 videos with 60 action categories and each human body contains 25 joints. There are two evaluation protocols for NTU-60: Cross-Subject (X-Sub) and Cross-View (X-View) protocols. X-Sub protocol means training data and validation data are split by different subjects, and half of the subjects are set as training sets and the rest are the test sets. X-View protocol means training data and validation data are collected from different camera views (camera 1,2 and 3). In X-View, the samples captured by camera 2 and 3 are set for training and the samples of camera 1 are set as testing set.

NTU-RGB+D 120 (NTU-120). NTU-120 is an expansion dataset of NTU-60 with 113,945 sequences with 120 action labels. There are also two evaluation protocols: Cross-Subject (X-Sub) and Corrss-Set (X-Set). In X-Sub, there are 53 subjects for training and 53 subjects for testing. In X-Set, half of the setups are split for training (even setup IDs) and the rest (odd setup IDs) are used for testing.

4.2. Experimental settings

Our experiments are performed on $8 \times$ A6000 GPUs with Pytorch [28] framework implementation. Both our pre-training and fine-tuning models are trained by Adam optimizer [19] with base learning rate 0.005 and weight decay 0.0001. The batch size is 64. The pre-training and fine-tuning epoch number are all set to 200. We also use a multi-step learning rate schedule for learning rate adjustment with gamma 0.1 and milestones are 60 epoch, 90 epoch and 110 epoch. For fair comparisons among different methods, we limit the length of the skeleton sequence to 20 frames for all experiments.

STTFormer. As mentioned in Sec 3.1, in order to learn better spatial-temporal representations, we utilize the standard STTFormer as our backbone in pre-training, which consists of a stack of STTA Blocks. Same as MAE, our method also adds sine-cosine positional embedding to both encoder and decoder inputs. For fine-tuning, we use the STTFormer-based encoder as our feature extractor.

Sequence division and patch embedding. In our research, we follow the patch embedding method in STTFormer. We first divide the original skeleton sequence into tuples. Then, since the skeleton data does not contain a large number of pixels and various noises like image data, we directly use a 1×1 Conv for patch embedding processing.

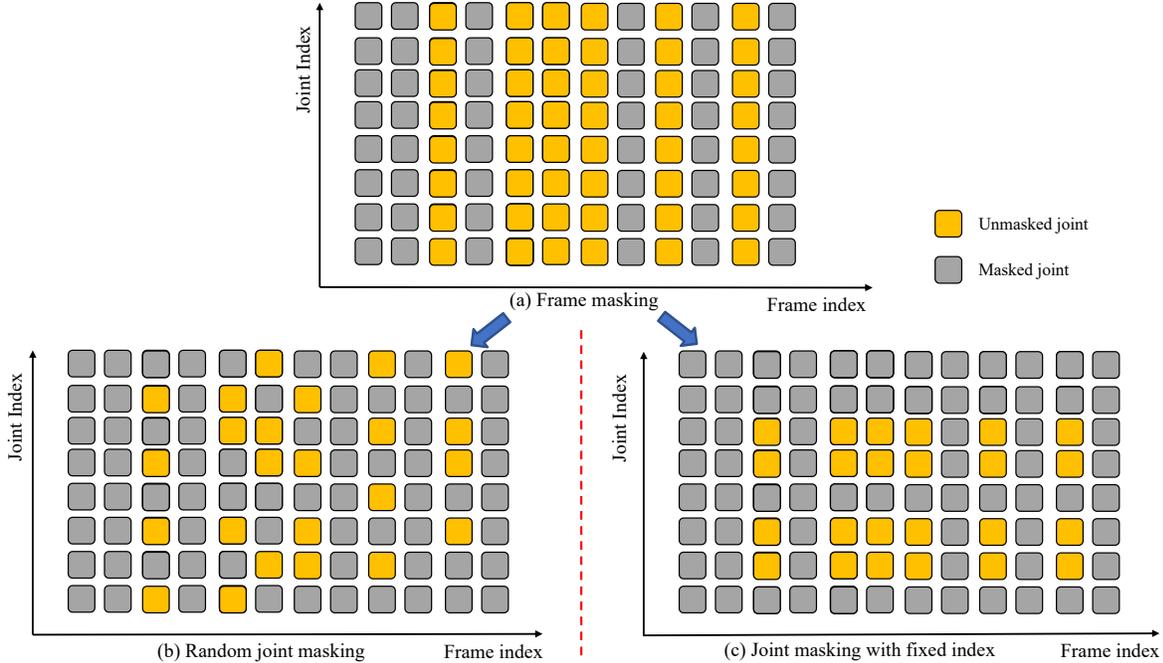


Figure 3: Illustration of two masking strategies. (a) The frame-masking is first implemented and then: (b) randomly mask the joints in spatial level; (c) mask the joint with the fixed index.

Masking settings. We implement our masking strategies before sequence division. As we discussed in Sec. 3.2, we first mask out a random subset of frames by the pre-set frame masking ratio and then mask out a random index of joints by the pre-set joint masking ratio. During experiments, we test several trials of frame-masking ratio and joint-masking ratio, finding the best trade-off combination.

Pre-training. We choose MSE loss as pre-training loss and save the best model by the minimized validation loss.

Fine-tuning. As we discussed in Sec. 3.4, we use end-to-end fine-tuning for the end task. Moreover, we choose cross entropy loss with label smoothing [17] as fine-tuning loss with smoothing rate 0.1 and save the best model by the maximized validation accuracy.

4.3. Ablation study

Different masking strategies. After performing the same degree of random frame masking, we compare the masking strategies of masking the joints randomly with the method keeping the same masked joint index over the entire sequence. The experimental results show that the pure random joint masking for visible frames is more helpful for the final fine-tuning result (in Table 6, we get our best fine-tuned recognition accuracy of 86.6% on X-sub using random masking strategy, which is 1.2% better than the best result of the masking method by fixing the joints indices). The overall results indicate that the random masking method outperforms the masking method with fixed joints

indices, which means the model learns better features with a randomly generated input than the pre-defined input. Notably, MAE experiment also shows that using a more random masking strategy is more beneficial to the final fine-tuning result.

Frame-masking ratio and joint-masking ratio. In spatial and temporal domains, we test several combinations of different frame-masking ratio and joint-masking ratio on SkeletonMAE. Following both joint index fixed and random masking strategies, we set the frame masking ratio 0.4, 0.5 and 0.6 respectively, for every decided frame masking ratio, we test different joint masking ratio (0.4, 0.5 and 0.6 respectively). As shown in Table 6, the final results on NTU-60 with X-Sub show that a frame-masking ratio of 0.4 and a joint-masking ratio of 0.5 work best in the masking method with fixed joints indices (85.4% accuracy). Using the random masking method, we achieve the best result (86.6% accuracy) in two combinations (0.5 joint-masking ratio with 0.5 or 0.4 frame-masking ratio).

Embedding dimension. Table 7 shows the ablation study on the embedding dimension of the decoder. We change the different embedding dimensions in SkeletonMAE decoder and find that the default setting with 256 dimension works better (86.6% accuracy) than the larger size (86.0% accuracy) and the small size (85.2% accuracy). We also observe that with the increasing size of embedding dimension, the number of model parameters increase as well, when we set the dimension as 512, the parameters are 11

method	frame-masking ratio	joint-masking ratio	NTU-60 X-Sub
fixed index	0.6	0.4	85.2
	0.6	0.5	84.9
	0.6	0.6	85.3
	0.5	0.4	85.3
	0.5	0.5	85.0
	0.5	0.6	84.8
	0.4	0.4	84.8
	0.4	0.5	85.4
random	0.4	0.6	85.2
	0.6	0.4	86.5
	0.6	0.5	86.0
	0.6	0.6	86.3
	0.5	0.4	86.3
	0.5	0.5	86.6
	0.5	0.6	85.7
	0.4	0.4	85.6
	0.4	0.5	86.6
	0.4	0.6	85.4

Table 1: Masking strategies with joint-masking ratio and frame-masking ratio. Specifically, there are two joint masking methods tested: fixed indices masking and randomly masking.

times larger than the parameters with dimension 128, which costs more time for training. So we choose 256 as the default embedding dimension for the following ablation studies.

embedding dimension	NTU-60 X-Sub	parameters(M)
128	85.2	3
256	86.6	11
512	86.0	33

Table 2: Ablation study on embedding dimension.

Decoder depth. Decoder depth represents the number of the STTFormer blocks. According to the last ablation experiment, we set the embedding dimension (the width of the decoder) as the default size 256, and vary the decoder depth (11, 9, 7 and 5 blocks). As the results shown in Table 8, SkeletonMAE achieves the best result (86.6% accuracy) when the decoder depth is 9. The deep depth (11 blocks with 86.5% accuracy) and shallow depth (7 blocks with 86.2% accuracy and 5 blocks with 85.7% accuracy) perform worse. According to the results from embedding dimension and decoder depth experiments, we finalize our default decoder configurations for the following experiments (256 embedding dimension and 9 blocks).

decoder depth	NTU 60 X-Sub
11	86.5
9	86.6
7	86.2
5	85.7

Table 3: Ablation study on decoder depth.

Pre-training schedule. Normally, a longer pre-training schedule will give an improvement, thus in this ablation

study, we increase the pre-training epoch from 50 epoch to 200 epoch, and test the best fine-tuned results at every 50 epoch. As it shown in Fig. 4, the best accuracy is 86.6%, so we select 200 epoch as the default pre-training epoch for the following experiments. It is worth noting that there is an impressive improvement (5.0%) between 50 epoch to 100 epoch, but a slight improvement (0.2%) between 150 epoch to 200 epoch, which means it is not cost-effective to keep increasing the pre-training epoch.

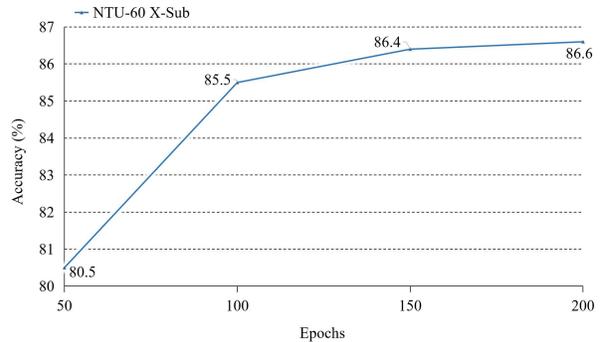


Figure 4: Ablation study on pre-training schedule.

4.4. Comparison with state-of-the-art

Self-supervised training. Notably, as we can see from Table 10, our SkeletonMAE outperforms two latest self-supervised skeleton action recognition methods: CrosSCLR [20] and AimCLR [14]. For a fair comparison, we replace their backbone networks (both of them use ST-GCN as the backbone) with STTFormer under the same settings. The results show that on NTU-60 dataset, our SkeletonMAE leads CrosSCLR 2.0% and AimCLR 2.7% on X-Sub, and also leads CrosSCLR 2.4% and AimCLR 2.5% under X-View protocol. As for the results on NTU-120 dataset, SkeletonMAE outperforms CrosSCLR by 1.8% and 1.2% on X-Sub and X-Set, and also outperforms AimCLR by 2.2% and 1.9% on X-Sub and X-Set respectively. The results indicate that our SkeletonMAE not only achieves outperforming results on the small-size dataset but also the large-size dataset.

method	backbone	NTU-60		NTU-120	
		X-Sub	X-View	X-Sub	X-Set
CrosSCLR[20]	ST-GCN	82.2	88.9	73.6	75.3
AimCLR[14]	ST-GCN	83.0	89.2	76.4	76.7
CrosSCLR[20]	STTFormer	84.6	90.5	75.0	77.9
AimCLR[14]	STTFormer	83.9	90.4	74.6	77.2
SkeletonMAE	STTFormer	86.6	92.9	76.8	79.1

Table 4: Fine-tuned results on NTU-60 and NTU-120 datasets.

Fewer labeled data training. In order to figure out the ability of spatial-temporal feature learning in fewer-

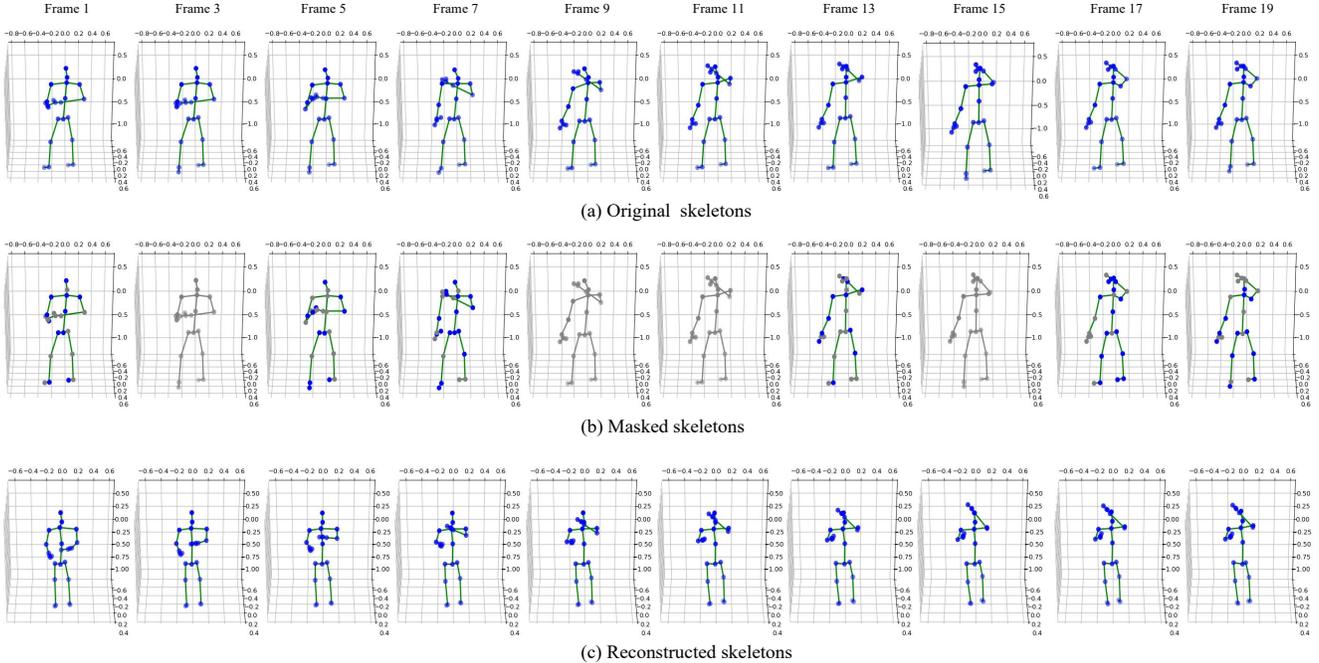


Figure 5: Visualization results on NTU-60 dataset, *drink water* action. We select the odd frames from the first 20 frames, for each frame, we visualize: (a) the input skeleton data; (b) the masked skeleton data (0.5 joint-masking ratio); (c) the reconstructed skeleton frames .

data situation, we fine-tune our pre-trained SkeletonMAE model with only 5% and 10% labeled data on both NTU-60 and NTU-120 datasets. According to Table 11, our SkeletonMAE achieves 64.4% and 68.8% on NTU-60 X-Sub and X-View with only 5% fine-tuning data, and surpasses CrossSCLR and AimCLR. Moreover, our SkeletonMAE also performs better than CrossSCLR and AimCLR with 10% labeled data (73.0% and 76.9% on NTU-60 X-Sub and X-View respectively). Meanwhile, our SkeletonMAE achieves outperformed results on NTU-120 data with 5% (50.4% on X-Sub and 52.0% on X-Set) and 10% (61.8% on X-Sub and 62.5% on X-Set) labeled data, which demonstrates a better capability of generalizability learning of our approach under the extreme fine-tuning situation.

method	backbone	label fraction	NTU-60		NTU-120	
			X-Sub	X-View	X-Sub	X-Set
CrosSCLR[20]	STTFormer	5%	63.5	66.9	50.2	50.4
AimCLR[14]	STTFormer	5%	63.9	67.5	49.0	51.8
SkeletonMAE	STTFormer	5%	64.4	68.8	50.4	52.0
CrosSCLR[20]	STTFormer	10%	71.0	75.1	58.5	60.6
AimCLR[14]	STTFormer	10%	70.2	76.2	58.6	60.5
SkeletonMAE	STTFormer	10%	73.0	76.9	61.8	62.5

Table 5: Fine-tuned results with fewer labeled data on NTU-60 and NTU-120 datasets.

4.5. Visualization

In Fig. 5, we show the visualization results of SkeletonMAE pre-training on NTU-60 dataset using randomly joint masking strategy with 0.5 joint masking ratio (a frame-

masking is applied first). We select the odd frames from the first 20 frames from the *drink water* action. As it shown: Fig. 5 (a) visualizes the input skeleton, Fig. 5 (b) shows the corresponding masking results in both frame-level (frame 3, 9, 11, 15) and joint-level (frame 1, 5, 7, 13, 17, 19), Fig. 5 (c) shows the reconstructed skeletons of the pre-training. Spatially (the visualization results in the same frame), we observe that there exist a few detailed differences between the original skeleton sequence and the reconstructed skeleton sequence, but the frameworks of the human body (e.g., the positions of arms and legs) are kept without distortion. The detail difference visualization shows the good ability of SkeletonMAE for spatial-feature learning. Temporally (the consecutive skeleton sequences), although we also observe a few variations between the original sequence and the reconstructed ones, there is no pronounced deformation in the time space (the joint motion in different frames is reserved, e.g., rising hands), which indicates that our SkeletonMAE learns temporal representations well. The overall results demonstrate that SkeletonMAE learns generalized skeleton sequences containing semantic action information, resulting a good performance in action recognition task.

5. Conclusion

We conduct a novel skeleton-based masked autoencoder named SkeletonMAE for self-supervised skeleton action recognition. In order to get a better skeleton representation learning, we apply a novel spatial-temporal masking strat-

egy in pre-training for skeleton reconstruction. The roles of different frame-ratio and joint-ratio are also discussed and implemented. With comprehensive experiments on NTU-60 and NTU-120 datasets, we show outperformed results of SkeletonMAE for skeleton action recognition.

References

- [1] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [2] Zhe Cao, Tomas Simon, Shih-En Wei, and Yaser Sheikh. Realtime multi-person 2d pose estimation using part affinity fields. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7291–7299, 2017.
- [3] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9650–9660, 2021.
- [4] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020.
- [5] 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. 2022.
- [6] Ke Cheng, Yifan Zhang, Xiangyu He, Weihang Chen, Jian Cheng, and Hanqing Lu. Skeleton-based action recognition with shift graph convolutional network. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 183–192, 2020.
- [7] Yi-Bin Cheng, Xipeng Chen, Dongyu Zhang, and Liang Lin. Motion-transformer: self-supervised pre-training for skeleton-based action recognition. In *Proceedings of the 2nd ACM International Conference on Multimedia in Asia*, pages 1–6, 2021.
- [8] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [9] Carl Doersch, Abhinav Gupta, and Alexei A Efros. Unsupervised visual representation learning by context prediction. In *Proceedings of the IEEE international conference on computer vision*, pages 1422–1430, 2015.
- [10] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- [11] Yong Du, Yun Fu, and Liang Wang. Skeleton based action recognition with convolutional neural network. In *2015 3rd IAPR Asian conference on pattern recognition (ACPR)*, pages 579–583. IEEE, 2015.
- [12] Yong Du, Wei Wang, and Liang Wang. Hierarchical recurrent neural network for skeleton based action recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1110–1118, 2015.
- [13] Spyros Gidaris, Praveer Singh, and Nikos Komodakis. Unsupervised representation learning by predicting image rotations. *arXiv preprint arXiv:1803.07728*, 2018.
- [14] Tianyu Guo, Hong Liu, Zhan Chen, Mengyuan Liu, Tao Wang, and Runwei Ding. Contrastive learning from extremely augmented skeleton sequences for self-supervised action recognition. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 762–770, 2022.
- [15] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16000–16009, 2022.
- [16] 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.
- [17] Tong He, Zhi Zhang, Hang Zhang, Zhongyue Zhang, Junyuan Xie, and Mu Li. Bag of tricks for image classification with convolutional neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 558–567, 2019.
- [18] Yonghong Hou, Zhaoyang Li, Pichao Wang, and Wanqing Li. Skeleton optical spectra-based action recognition using convolutional neural networks. *IEEE Transactions on Circuits and Systems for Video Technology*, 28(3):807–811, 2016.
- [19] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [20] Linguo Li, Minsi Wang, Bingbing Ni, Hang Wang, Jiancheng Yang, and Wenjun Zhang. 3d human action representation learning via cross-view consistency pursuit. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4741–4750, 2021.
- [21] Yunfan Li, Peng Hu, Zitao Liu, Dezhong Peng, Joey Tianyi Zhou, and Xi Peng. Contrastive clustering. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 8547–8555, 2021.
- [22] Lilang Lin, Sijie Song, Wenhan Yang, and Jiaying Liu. Ms2l: Multi-task self-supervised learning for skeleton based action recognition. In *Proceedings of the 28th ACM International Conference on Multimedia*, pages 2490–2498, 2020.
- [23] Jun Liu, Amir Shahroudy, Mauricio Perez, Gang Wang, Ling-Yu Duan, and Alex C Kot. Ntu rgb+d 120: A large-scale benchmark for 3d human activity understanding. *IEEE transactions on pattern analysis and machine intelligence*, 42(10):2684–2701, 2019.
- [24] Jun Liu, Gang Wang, Ling-Yu Duan, Kamila Abdiyeva, and Alex C Kot. Skeleton-based human action recognition with global context-aware attention lstm networks. *IEEE Transactions on Image Processing*, 27(4):1586–1599, 2017.

- [25] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.
- [26] Zelun Luo, Boya Peng, De-An Huang, Alexandre Alahi, and Li Fei-Fei. Unsupervised learning of long-term motion dynamics for videos. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2203–2212, 2017.
- [27] Yunsheng Pang, Qihong Ke, Hossein Rahmani, James Bailey, and Jun Liu. Igformer: Interaction graph transformer for skeleton-based human interaction recognition. 2022.
- [28] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019.
- [29] Chiara Plizzari, Marco Cannici, and Matteo Matteucci. Spatial temporal transformer network for skeleton-based action recognition. In *International Conference on Pattern Recognition*, pages 694–701. Springer, 2021.
- [30] Helei Qiu, Biao Hou, Bo Ren, and Xiaohua Zhang. Spatio-temporal tuples transformer for skeleton-based action recognition. *arXiv preprint arXiv:2201.02849*, 2022.
- [31] Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. *arXiv preprint arXiv:1908.10084*, 2019.
- [32] Amir Shahroudy, Jun Liu, Tian-Tsong Ng, and Gang Wang. Ntu rgb+ d: A large scale dataset for 3d human activity analysis. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1010–1019, 2016.
- [33] Lei Shi, Yifan Zhang, Jian Cheng, and Hanqing Lu. Two-stream adaptive graph convolutional networks for skeleton-based action recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12026–12035, 2019.
- [34] Sijie Song, Cuiling Lan, Junliang Xing, Wenjun Zeng, and Jiaying Liu. An end-to-end spatio-temporal attention model for human action recognition from skeleton data. In *Proceedings of the AAAI conference on artificial intelligence*, volume 31, 2017.
- [35] Tae Soo Kim and Austin Reiter. Interpretable 3d human action analysis with temporal convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 20–28, 2017.
- [36] Nitish Srivastava, Elman Mansimov, and Ruslan Salakhudinov. Unsupervised learning of video representations using lstms. In *International conference on machine learning*, pages 843–852. PMLR, 2015.
- [37] Ke Sun, Bin Xiao, Dong Liu, and Jingdong Wang. Deep high-resolution representation learning for human pose estimation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5693–5703, 2019.
- [38] Zehua Sun, Qihong Ke, Hossein Rahmani, Mohammed Bennamoun, Gang Wang, and Jun Liu. Human action recognition from various data modalities: A review. *IEEE transactions on pattern analysis and machine intelligence*, 2022.
- [39] Alexander Toshev and Christian Szegedy. DeepPose: Human pose estimation via deep neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1653–1660, 2014.
- [40] Raviteja Vemulapalli, Felipe Arrate, and Rama Chellappa. Human action recognition by representing 3d skeletons as points in a lie group. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 588–595, 2014.
- [41] Raviteja Vemulapalli and Rama Chellappa. Rolling rotations for recognizing human actions from 3d skeletal data. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4471–4479, 2016.
- [42] Pascal Vincent, Hugo Larochelle, Yoshua Bengio, and Pierre-Antoine Manzagol. Extracting and composing robust features with denoising autoencoders. In *Proceedings of the 25th international conference on Machine learning*, pages 1096–1103, 2008.
- [43] Hongsong Wang and Liang Wang. Modeling temporal dynamics and spatial configurations of actions using two-stream recurrent neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 499–508, 2017.
- [44] Jiang Wang, Zicheng Liu, Ying Wu, and Junsong Yuan. Learning actionlet ensemble for 3d human action recognition. *IEEE transactions on pattern analysis and machine intelligence*, 36(5):914–927, 2013.
- [45] Pichao Wang, Zhaoyang Li, Yonghong Hou, and Wanqing Li. Action recognition based on joint trajectory maps using convolutional neural networks. In *Proceedings of the 24th ACM international conference on Multimedia*, pages 102–106, 2016.
- [46] Sijie Yan, Yuanjun Xiong, and Dahua Lin. Spatial temporal graph convolutional networks for skeleton-based action recognition. In *Thirty-second AAAI conference on artificial intelligence*, 2018.
- [47] Jinlu Zhang, Zhigang Tu, Jianyu Yang, Yujin Chen, and Junsong Yuan. Mixste: Seq2seq mixed spatio-temporal encoder for 3d human pose estimation in video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13232–13242, 2022.
- [48] Pengfei Zhang, Cuiling Lan, Junliang Xing, Wenjun Zeng, Jianru Xue, and Nanning Zheng. View adaptive neural networks for high performance skeleton-based human action recognition. *IEEE transactions on pattern analysis and machine intelligence*, 41(8):1963–1978, 2019.
- [49] Pengfei Zhang, Cuiling Lan, Wenjun Zeng, Junliang Xing, Jianru Xue, and Nanning Zheng. Semantics-guided neural networks for efficient skeleton-based human action recognition. In *proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1112–1121, 2020.
- [50] Songyang Zhang, Xiaoming Liu, and Jun Xiao. On geometric features for skeleton-based action recognition using multilayer lstm networks. In *2017 IEEE Winter Conference on*

Applications of Computer Vision (WACV), pages 148–157. IEEE, 2017.

- [51] Ce Zheng, Matias Mendieta, Pu Wang, Aidong Lu, and Chen Chen. A lightweight graph transformer network for human mesh reconstruction from 2d human pose. *ACM Multimedia*, 2022.
- [52] Ce Zheng, Sijie Zhu, Matias Mendieta, Taojiannan Yang, Chen Chen, and Zhengming Ding. 3d human pose estimation with spatial and temporal transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11656–11665, 2021.

Supplementary material

In this supplementary material, we provide the following items for better understanding the paper:

- Detailed architectures of SkeletonMAE encoder and decoder.
- More visualization results.
- Qualitative analysis on masking strategy.

A. Detailed architecture

As it shown in Table 6, we give the detailed architecture of SkeletonMAE, including the dimensions of input and output layers, the size of each STTA block (including the input dimension D_{in} and output dimension D_{out}) in the encoder and decoder.

Given the original 3D skeleton data:

$$X_{ori} \in R^{T \times J \times D}, \quad (2)$$

where T is the number of frames of the input sequence, J is the number of joints in each frame, D is the input dimension ($J = 25$ and $D = 3$ in NTU-60 and NTU-120 datasets). Then based on the spatial-temporal masking strategy we proposed, after applying a masking method M to X_{ori} , the input skeleton data is expressed as:

$$X_{in} = M(X_{ori}), \quad (3)$$

$$X_{in} \in R^{T' \times J' \times D_{in}}, \quad (4)$$

where T' and J' represent the masked frames and masked joints following the pre-set masking approach. In SkeletonMAE encoder, according to the structure of STTA block and the data processing from STTFormer, X_{in} is processed by a series of STTA blocks for data embedding:

$$X_{out} = STTA_i(x_{in}), i \in [1, \dots, N], \quad (5)$$

$$X_{out} \in R^{T' \times J' \times D_{out}} \quad (6)$$

where N is the number of STTA blocks and D_{out} is the output dimension ($D_{out} = 3$). As for the decoder, it has an inverted structure of the decoder (the output layer is at the end of the decoder).

Moreover, we provide the sizes of query \mathbf{Q} , key \mathbf{K} and value \mathbf{V} from the STTA blocks. Specifically, the embedding dimension for Table 6 is 256 and the number of blocks in both encoder and decoder is 9 (the default setting for ablation studies). Finally, we also give more settings of the ablation studies on embedding dimension and decoder depth (Table 7, 8, 9, 10, 11).

B. More visualization

We provide more visualization results of SkeletonMAE pre-training on NTU-60 dataset with random joint-masking method (joint-masking ratio is 0.5) in Fig. 7: *type on a keyboard*; Fig. 8: *taking a selfie*; Fig. 9: *back pain*; Fig. 10: *fan self*. As we discussed before, the SkeletonMAE learns generalizable features from skeleton data without pronounced deformation. However, the detailed differences between the original skeleton and the reconstructed skeleton in the same frame are still observed (*e.g.*, the coordinates of the forearms from the reconstructed skeletons are different from the original skeletons).

C. Qualitative analysis on masking strategy

In Fig. 6, we show the qualitative analysis results of two masking strategies (index fixed and random masking methods) on NTU-60 dataset X-Sub, with different combinations of frame-masking (α) and joint-masking (β) ratios. We set the coordinates (α, β) as the values of x-axis, with a descending order of the α . We observe that the overall results of the random masking are better than the fixed index masking. Specifically, the random masking method surpasses the fixed index masking method with the largest gap (1.6%) at (0.5,0.5), and there exists a smallest gap (0.2%) at (0.4,0.6).

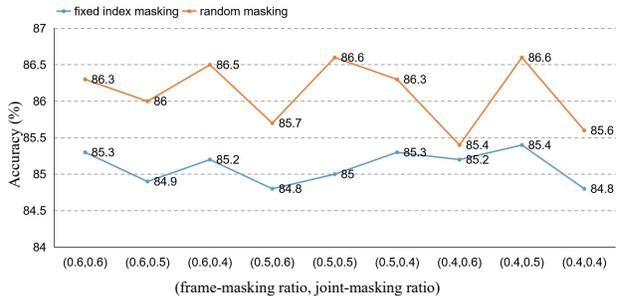


Figure 6: Qualitative analysis on different masking strategies (NTU-60 dataset X-Sub).

SkeletonMAE	layer name	input dim (D_{in})	output dim (D_{out})	QKV dim
encoder	input layer	3	64	
	Block1	64	64	16
	Block2	64	64	16
	Block3	64	128	32
	Block4	128	128	32
	Block5	128	256	64
	Block6	256	256	64
	Block7	256	256	64
	Block8	256	256	64
decoder	Block1	256	256	64
	Block2	256	256	64
	Block3	256	256	64
	Block4	256	128	64
	Block5	128	128	32
	Block6	128	64	32
	Block7	64	64	16
	Block8	64	64	16
	output layer	64	3	16

Table 6: The detailed structures of SkeletonMAE encoder and decoder, where the embedding dimension is 256 and the depth of decoder is 9.

SkeletonMAE	layer name	input dim (D_{in})	output dim (D_{out})	QKV dim
encoder	input layer	3	64	
	Block1	64	64	16
	Block2	64	128	32
	Block3	128	128	32
	Block4	128	256	64
	Block5	256	256	64
	Block6	256	512	128
	Block7	512	512	128
	Block8	512	512	128
decoder	Block1	512	512	128
	Block2	512	512	128
	Block3	512	256	64
	Block4	256	256	64
	Block5	256	128	32
	Block6	128	128	32
	Block7	128	64	16
	Block8	64	64	16
	output layer	64	3	16

Table 7: The detailed structures of SkeletonMAE encoder and decoder, where the embedding dimension is 512 and the depth of decoder is 9.

SkeletonMAE	layer name	input dim (D_{in})	output dim (D_{out})	QKV dim
encoder	input layer	3	32	
	Block1	32	32	8
	Block2	32	32	8
	Block3	32	64	16
	Block4	64	64	16
	Block5	64	128	32
	Block6	128	128	32
	Block7	128	128	32
	Block8	128	128	32
decoder	Block1	128	128	32
	Block2	128	128	32
	Block3	128	128	32
	Block4	128	64	32
	Block5	64	64	16
	Block6	64	32	16
	Block7	32	32	8
	Block8	32	32	8
	output layer	32	3	8

Table 8: The detailed structures of SkeletonMAE encoder and decoder, where the embedding dimension is 128 and the depth of decoder is 9.

SkeletonMAE	layer name	input dim (D_{in})	output dim (D_{out})	QKV dim
encoder	input layer	3	64	
	Block1	64	64	16
	Block2	64	64	16
	Block3	64	128	32
	Block4	128	128	32
	Block5	128	256	64
	Block6	256	256	64
	Block7	256	256	64
	Block8	256	256	64
decoder	Block1	256	128	64
	Block2	128	128	32
	Block3	128	64	32
	Block4	64	64	16
	output layer	64	3	16

Table 9: The detailed structures of SkeletonMAE encoder and decoder, where the embedding dimension is 256 and the depth of decoder is 5.

SkeletonMAE	layer name	input dim (D_{in})	output dim (D_{out})	QKV dim
encoder	input layer	3	64	
	Block1	64	64	16
	Block2	64	64	16
	Block3	64	128	32
	Block4	128	128	32
	Block5	128	256	64
	Block6	256	256	64
	Block7	256	256	64
	Block8	256	256	64
decoder	Block1	256	256	64
	Block2	256	256	64
	Block3	256	128	64
	Block4	128	128	32
	Block5	128	64	32
	output layer	64	3	16

Table 10: The detailed structures of SkeletonMAE encoder and decoder, where the embedding dimension is 256 and the depth of decoder is 7.

SkeletonMAE	layer name	input dim (D_{in})	output dim (D_{out})	QKV dim
encoder	input layer	3	64	
	Block1	64	64	16
	Block2	64	64	16
	Block3	64	128	32
	Block4	128	128	32
	Block5	128	256	64
	Block6	256	256	64
	Block7	256	256	64
	Block8	256	256	64
decoder	Block1	256	256	64
	Block2	256	256	64
	Block3	256	256	64
	Block4	256	256	64
	Block5	256	128	64
	Block6	128	128	32
	Block7	128	128	32
	Block8	128	64	32
	Block9	64	64	16
	Block10	64	64	16
	output layer	64	3	16

Table 11: The detailed structures of SkeletonMAE encoder and decoder, where the embedding dimension is 256 and the depth of decoder is 11.

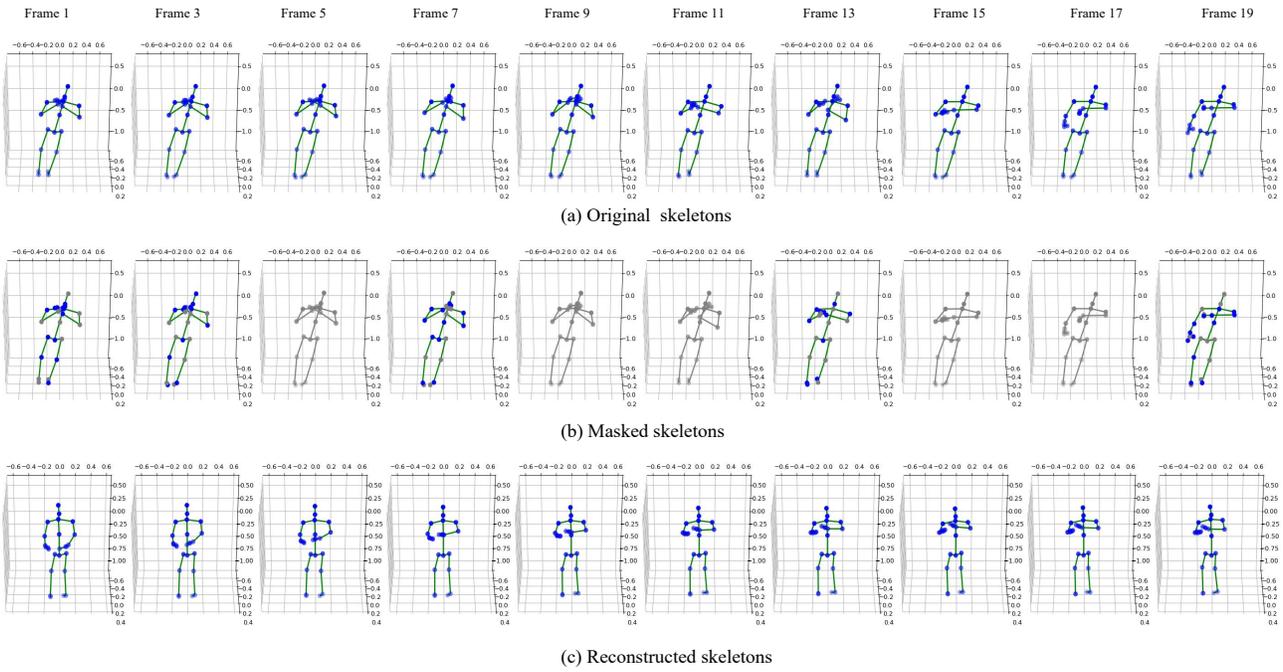


Figure 7: Visualization results on NTU-60 dataset, *type on a keyboard* action.

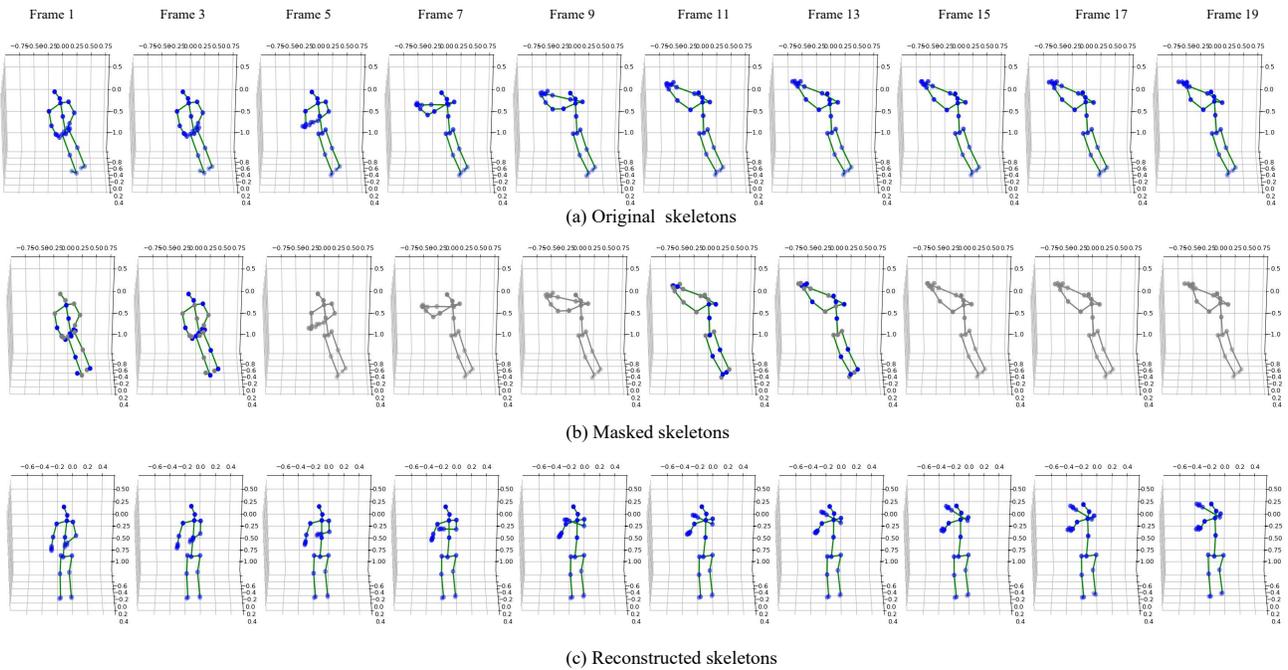


Figure 8: Visualization results on NTU-60 dataset, *taking a selfie* action.

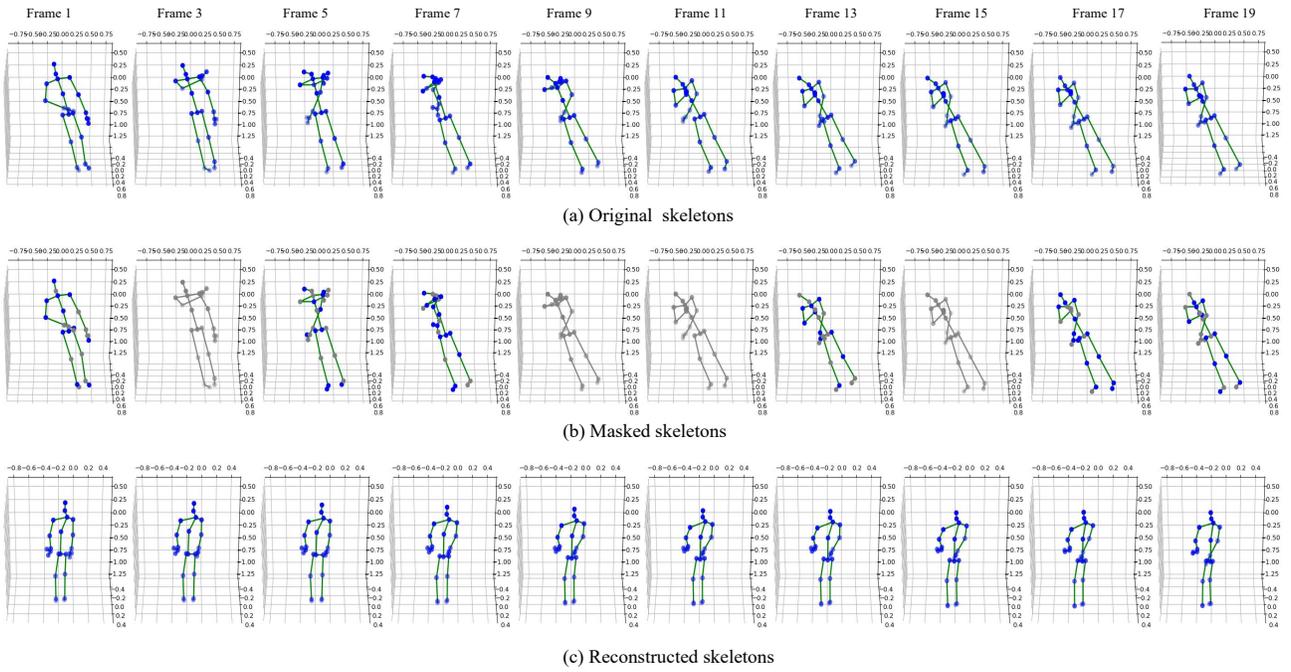


Figure 9: Visualization results on NTU-60 dataset, *back pain* action.

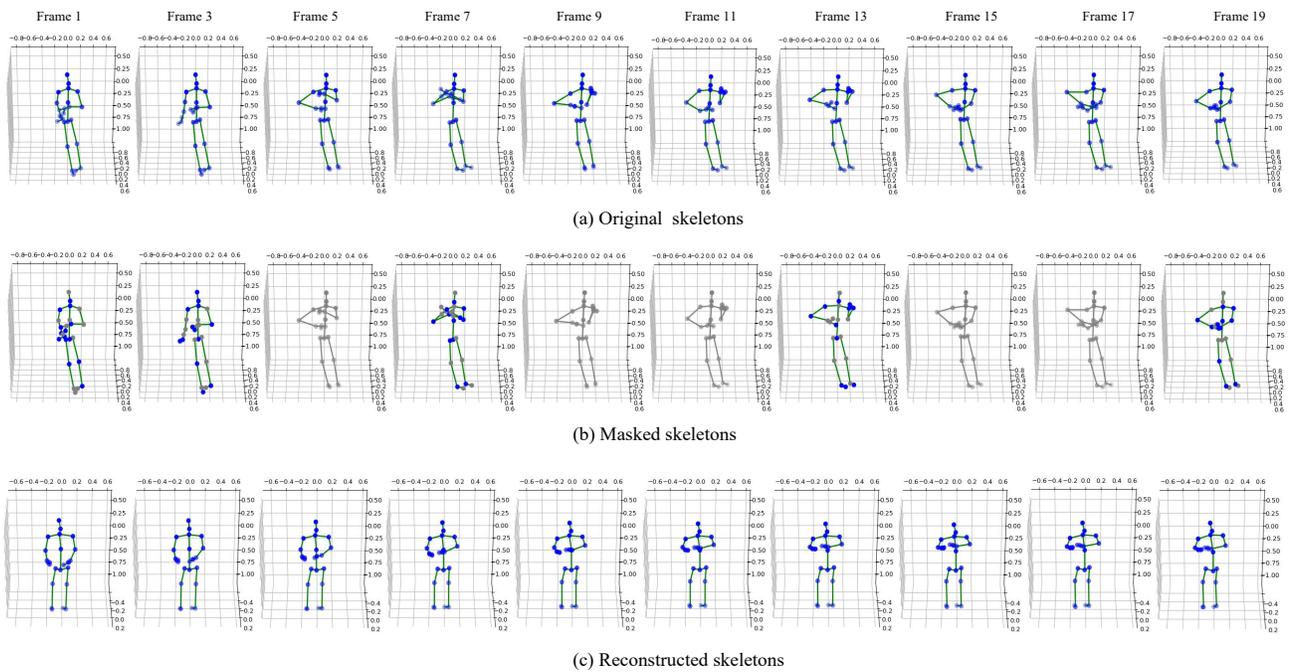


Figure 10: Visualization results on NTU-60 dataset, *fan self* action.