# Generative Data Augmentation for Skeleton Action Recognition

*Abstract*— **Skeleton-based human action recognition is a powerful approach for understanding human behaviour from pose data, but collecting large-scale, diverse, and well-annotated 3D skeleton datasets is both expensive and labor-intensive. To address this challenge, we propose a conditional generative pipeline for data augmentation in skeleton action recognition. Our method learns the distribution of real skeleton sequences under the constraint of action labels, enabling the synthesis of diverse and high-fidelity data. Even with limited training samples, it can effectively generate skeleton sequences and achieve competitive recognition performance in low-data scenarios, demonstrating strong generalisation in downstream tasks. Specifically, we introduce a Transformer-based encoder–decoder architecture, combined with a generative refinement module and a dropout mechanism, to balance fidelity and diversity during sampling. Experiments on HumanAct12 and the refined NTU-RGBD (NTU-VIBE) dataset show that our approach consistently improves the accuracy of multiple skeleton-based action recognition models, validating its effectiveness in both few-shot and full-data settings. The code will be released upon acceptance.**

## I. INTRODUCTION

Human action recognition is a key task in computer vision with applications in human-computer interaction, video surveillance, healthcare, and virtual reality. Among various modalities, 3D skeleton-based action recognition has emerged as a lightweight, privacy-preserving solution. It encodes only the positions of key joints, making it robust to appearance, lighting, and background variations, while being efficient in storage and computation.

However, acquiring large-scale, high-quality skeleton datasets remains challenging. High-precision optical motion capture systems require expensive specialised equipment, with costs often exceeding $10,000 [27]. Therefore, other datasets compromise by relying on depth sensors (e.g., Kinect V2) [33] or multi-view camera setups [48], [16], which still demand controlled environments and active subject participation. Moreover, even in carefully controlled settings, the captured data can still be highly cumbersome, noisy, especially from depth sensors, and resource-intensive. Deep learning based pose estimation methods can extract 3D poses from RGB inputs [19], [3], but the results often suffer from noise and inconsistencies, especially in unconstrained scenes.

To address the high collection cost, limited diversity, and noise in existing datasets, many approaches have explored data augmentation [41], [39], [24], [15], [6]. These methods fall into two categories: transformation-based, which apply spatial and temporal perturbations (e.g., rotation, scaling, noise), and generation-based, which synthesise new sequences using frameworks like VAEs, GANs, or diffusion

models [34], [30]. While the former often require careful tuning of hyperparameters, the latter, while capable of learning the underlying data distribution to generate realistic samples, frequently suffers from limited diversity and a strong dependence on large-scale data. Building on MDM [37], we introduce a conditional semantic encoder and the fidelity–diversity control module, and replace classifier-free guidance with classifier guidance during sampling to prioritise class alignment for recognition explicitly.

Specifically, as shown in 1, this work proposes a conditional diffusion-based data augmentation method for 3D skeleton-based action recognition. Our method can efficiently generate high-fidelity, diverse, discriminative, and label-consistent skeleton data, providing both realistic variations and strong supervision signals for downstream recognition models. In the training phase, our method employs a Transformer encoder to extract the semantic information of the original skeleton data, while incorporating action labels as supervision signals. The Transformer decoder takes the noise tokens together with the conditional representation, which integrates semantic features, temporal information, and action labels. Guided by both reconstruction and classification objectives, it progressively denoises the tokens into label-consistent skeleton sequences, while jointly capturing structural priors and maintaining label consistency. In the inference phase, the diffusion model generates realistic and label-consistent skeleton sequences. To further improve generation quality, we design a Generative Refinement Module (GRM) and introduce a sampling-time dropout mechanism to balance fidelity and diversity, encouraging the model to produce discriminative and label-consistent variations. Our method is highly efficient, requiring only a single training phase. Once trained, the model is capable of generating large-scale skeleton data during inference, while allowing explicit control over the trade-off between diversity and fidelity in the generated samples.

We conduct extensive experiments on HumanAct12 [11] and Refined NTU-RGBD (NTU-VIBE) [33], [11], evaluating generation quality and downstream recognition performance. Our method demonstrates strong generalisation, particularly in low-data scenarios, where adding synthetic samples significantly improves accuracy, reaching levels comparable to those achieved with full data training. We conducted comprehensive experiments to evaluate the effectiveness of our method. By assessing the generation results and the performance on downstream skeleton-based action recognition tasks, we demonstrated the superior performance of our approach. In scenarios with limited real data, adding synthetic samples significantly improves accuracy, reaching
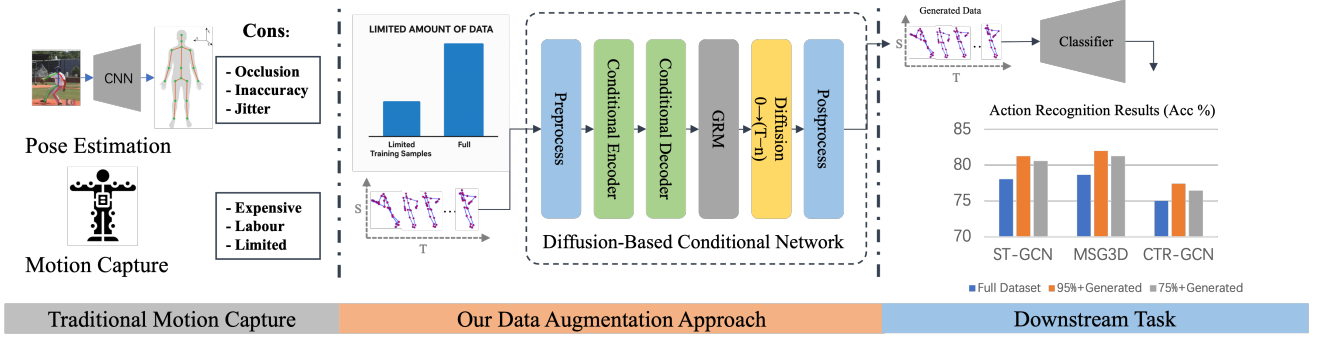
Fig. 1. Overview of our approach. With only a small set of labelled skeleton sequences, the model generates diverse and high-fidelity samples. When combined with a reduced amount of real data for training, these synthetic samples enable our skeleton action recognisers to achieve performance close to the state of the art on HumanAct12 and Refined NTU RGB+D.

levels comparable to using more real data that would be expensive and challenging to obtain. Furthermore, we conducted experiments to optimise the augmentation process by balancing diversity and fidelity in the synthetic data. Moreover, our data augmentation method demonstrates strong generality, as it can be adapted to various skeleton data formats and is compatible with a wide range of skeleton-based action recognition datasets and methods.

**Contributions:**

- We propose a conditional skeleton generation method based on diffusion models, conditioned on action labels, to generate diverse and realistic motion sequences. By generating large amounts of high-quality data from limited training samples, our approach reduces the need for costly large-scale data collection.
- We introduce a transformer encoder that extracts semantic representations from skeleton inputs and incorporates action labels as conditional signals to guide the diffusion-based generation process.
- We introduce a Generative Refinement Module (GRM) and sampling-time dropout to control fidelity and diversity in the synthetic data jointly.
- We validate our method across two datasets and multiple skeleton action recognition backbones, showing improvements in both few-shot and full-data training scenarios. Additionally, we conduct ablation studies to evaluate the contribution of each module and assess the quality of generated skeletons using standard metrics.

## II. RELATED WORK

### A. Diffusion Models.

Diffusion models [35], [36] are generative models that produce data by learning to reverse a progressive noising process. Denoising Diffusion Probabilistic Models (DDPM) [12], [36] and Denoising Diffusion Implicit Models (DDIM) have demonstrated state-of-the-art results in image generation. Conditional diffusion techniques, such as classifier guidance [8] and classifier-free guidance [13], enable fine-grained control during sampling. Beyond images, diffusion models have shown strong potential in motion generation tasks. Human motion is typically represented as sequences

of joint data in 2D, 3D, or SMPL [23], [45]. Recent works [37], [4], [29], [7], [18], [21] have shown strong success in synthesising realistic, diverse, and controllable motion sequences. While diffusion models have been explored for motion generation, no prior work has applied conditional diffusion for label-guided skeleton augmentation in recognition pipelines.

### B. Synthetic Data for Augmentation.

Data scarcity often leads to overfitting and poor generalisation in neural networks, especially under low-data regimes. Traditional augmentation methods [20] introduce simple transformations (e.g., flips, noise, crops) but are limited in diversity. Generative approaches overcome this by learning data distributions to produce new samples. Early work like DAGAN [1] and BigGAN [2] explored this idea to generate diverse image data for improving classification tasks. More recent efforts leverage text-to-image diffusion models. The study by Jahanian et al. [17] explored the feasibility of learning general-purpose visual representations from generative models instead of relying solely on original data. With the rapid development of diffusion models in recent years, this technology has become a new trend in generating training data, benefiting from its stationary training objective, high diversity, and conditional generation capabilities. [38] proposed DA-Fusion that utilised a large pre-trained text-to-image diffusion model to address the weaknesses of standard data augmentation while retaining the strengths. For skeleton data, augmentation is less explored. [26] analyses synthetic data on the fall-down detection task. [6] proposed a skeleton data augmentation method derived from observations of inaccuracies in human pose estimation. The works apply geometric perturbations (e.g., rotation, translation) or simulate occlusion. However, most do not model the complex distribution of temporal joint sequences. To our knowledge, this is the first work to apply conditional diffusion models for class-aware skeleton data augmentation, enabling label consistent generation at scale.

### C. Skeleton Action Recognition.

Early skeleton recognition methods relied on handcrafted features and classical classifiers [14], [40], but they re-
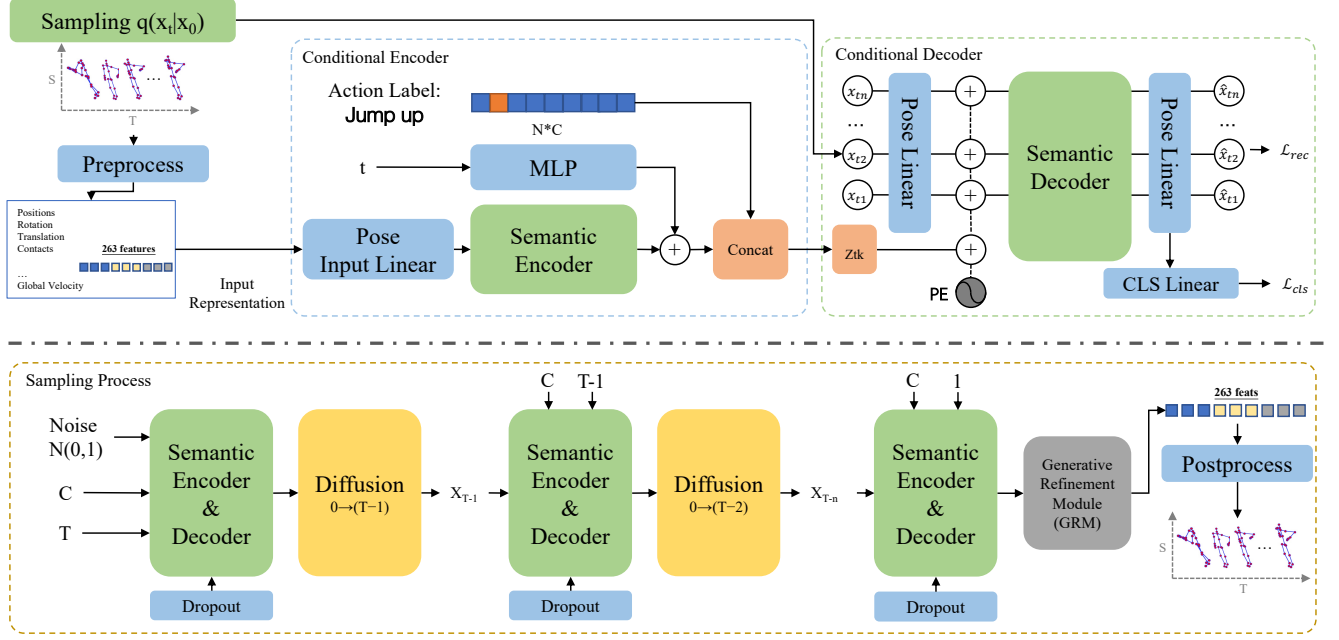
Fig. 2. Overview of our proposed network. (**Top**) **Conditional Skeleton Diffusion Module.** The encoder processes a skeleton feature sequence together with the noise step $t$ and the corresponding action label, producing a conditional representation. A Transformer-based decoder then reconstructs the clean skeleton sequence from the noise-corrupted input, guided by this representation. In addition, a lightweight classification network is introduced to encourage label-consistent generation. (**Bottom**) **Sampling Process.** The sampling input consists of the conditional representation and random noise, where the noise incorporates both label information and semantics from the original data. The decoder progressively denoises the sequence from step $T$ to 1, generating a clean skeleton motion. A Generative Refinement Module and Dropout further enhance the balance between semantic fidelity to the action and diversity of the generated motions.

quire manual feature design, are sensitive to noise/viewpoint changes, and poorly capture long-range dynamics. With the development of deep learning, recognition has shifted from handcrafted pipelines to end-to-end RNN/GCN/Transformer architectures that learn robust spatiotemporal representations from raw skeletons. Graph Convolutional Networks (GCNs) became the standard due to the ability to model the spatial and temporal relationships of skeleton data effectively. ST-GCN [42] introduced spatial-temporal graphs but incurred a high computational cost. MSG3D [22] captured multi-scale patterns; CTR-GCN [5] used channel-wise topology refinement to learn adaptive topologies and aggregates joint features for dynamic structure learning; BlockGCN [44] simplified the graph via blockwise partitioning, performing independent modelling within each block but with limited temporal modelling. We use these models as baselines to evaluate the benefit of our synthetic data.

## III. METHODOLOGY

### A. Diffusion Models Preliminary

Diffusion models [12], [25], [31] are generative frameworks that learn data distributions by simulating a forward process that gradually adds Gaussian noise, and a reverse process that removes it. In the forward process, the posterior distribution is implemented as a Markov chain that recursively adds noise to the sample through the conditional probability. This process can be denoted as:

$$q(\mathbf{x}_{1:T}|\mathbf{x}_0) = \prod_{t=1}^{T} q(\mathbf{x}_t|\mathbf{x}_{t-1}) \quad (1)$$

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1-\beta_t}\mathbf{x}_{t-1}, \beta_t \mathbf{I}) \quad (2)$$

Here $\beta_t$ is a variance schedule. Given a timestamp $t$, the $q(x_t)$ can be approximated as

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1-\bar{\alpha}_t}\varepsilon, \quad (3)$$

where $\bar{\alpha}_t = \prod_{s=1}^{t} \alpha_s$. Rather than predicting the noise $\varepsilon_t$, we follow recent work [37], [28] to directly predict the original sample $x_0$ itself from the noisy input. The training objective is:

$$L = \mathbb{E}_{\mathbf{x}_0 \sim q(\mathbf{x}_0|\mathbf{c}), t \sim [1,T]} \left[ \|\mathbf{x}_0 - G(\mathbf{x}_t, t, \mathbf{c})\|_2^2 \right], \quad (4)$$

where $t$ and $c$ denote the timestamps and condition, respectively.

### B. Conditional Diffusion Model

An overview of our pipeline is illustrated in Figure 2. The architecture is a Transformer model with action labels as conditioning signals. The model includes

- **Conditional Encoder:** A Transformer encoder to extract latent feature representations from the skeleton input data along with **timestep and label embeddings.**
- **Conditional Decoder:** A Transformer decoder reconstructs the original skeleton data, taking the encoder

TABLE I
THE 263-DIMENSIONAL FEATURE VECTOR EXPLANATIONS.

| Component | Dimensions | Description |
|---|---|---|
| Joint Positions | $22 \times 3 = 66$ | 3D coordinates $(x, y, z)$ for 22 joints |
| Joint Velocities | $22 \times 3 = 66$ | Velocity vectors for each joint |
| Joint Rotations | $22 \times 6 = 132$ | 6D rotation representations (more stable than quaternions or Euler angles) |
| Global Translation | 3 | Overall body translation in 3D space |
| Global Velocity | 3 | Global movement velocity of the body |
| **Total** | **263** | Combined total of all feature components |

features, concatenated with the skeleton data corrupted by noise through the diffusion process.

During sampling, we use a Generative Refinement Module (GRM) to discard low-fidelity generations and apply dropout to promote diversity further, ensuring that the final output is both discriminative and robust for downstream tasks.

### C. Input Representation.

Skeleton data is compact but semantically rich. In our setting, the original HumanAct12 dataset provides 3D coordinates for 22 skeletal joints. Following the HumanML3D representation [10], we convert each frame into a 263-dimensional feature vector, where the 22 joints are re-encoded to jointly capture 3D positions, local orientations, and dynamic attributes such as velocities. This extended representation offers a more comprehensive description of human motion, preserving both spatial configurations and temporal dynamics, while remaining computationally efficient compared to raw mesh or video data. Detailed construction of the 263-dimensional features is provided in I.

### D. Conditional Encoder

The input 263-dimensional feature sequence is first processed with temporal positional embeddings to preserve frame-wise order information. The action label $c$ is represented as a one-hot vector and embedded through an MLP, while the diffusion timestep $t$ is similarly mapped into the latent space. These conditional embeddings are concatenated and projected as a prefix token $z_{tk}$, which is then prepended to the feature sequence and fed into the encoder. The conditional encoder allows the model to incorporate both semantic (action label) and temporal (timestep) guidance during representation learning.

### E. Conditional Decoder

The decoder takes the noisy feature sequence together with the conditional prefix token $z_{tk}$ and performs token-level self-attention to reconstruct the underlying motion dynamics. It outputs a denoised 263-dimensional feature sequence, which is then passed through a 2-layer MLP classifier to predict the action label. This auxiliary classification objective provides label supervision, ensuring that the generated motion not only reduces diffusion noise but also remains consistent with the intended action semantics.

### F. Sampling Process

Our sampling involves predicting the clean sample $\hat{x}_0$ at each time step $t$, and then adding noise to regress it back to $x_{t-1}$. This iterative process continues from $t = T$ until $t = 0$, producing the final sample $x_0$. Unlike previous work [32], [37], which uses classifier-free guidance (occasionally masking conditions), we condition explicitly on labels throughout training and sampling, as fidelity to specific actions is essential for data augmentation. To encourage sample diversity and prevent the model from overfitting to label-conditioned patterns, we apply dropout within the denoising network during the sampling process. The stochasticity introduced in token activations allows our model to take a single action label as input and generate multiple diverse motion sequences with subtle variations not only in joint dynamics but also in higher-level semantics such as speed, thereby enriching data diversity without the need for extensive skeleton data collection.

### G. Generative Refinement Module (GRM)

The GRM evaluates generated samples $\hat{x}_0$ using a deviation measure $d(\hat{x}_0, x_0)$. Samples exceeding the threshold $\tau$ are discarded, and the retained set is defined as

$$\mathscr{S} = \{\hat{x}_0 \mid d(\hat{x}_0, x_0) \leq \tau\}, \tag{5}$$

where $x_0$ denotes the reference ground-truth sample (or its conditional embedding), $d(\cdot)$ is the deviation metric (e.g., $\ell_2$ distance in the 263-dimensional feature space), and $\tau$ is the deviation threshold. This filtering ensures that retained samples remain close to the real distribution (fidelity), while the combination with sampling-time dropout introduces diverse yet label-consistent variations.

### H. Loss function

Our total loss combines a **Reconstruction loss** and **Classification loss**. The Reconstruction loss $\mathscr{L}_{rec}$ enforces the generated samples to match the target data in the integrated 263-dimensional feature space. where $G(x_t, t, c)$ is the generated skeleton and $x_0$ is the ground truth. The Classification loss $\mathscr{L}_{cls}$ is a cross-entropy loss applied to the predicted action class of the generated data.

$$\mathscr{L}_{rec} = \mathbb{E}_{x_0, t}\left[\|x_0 - G(x_t, t, c)\|_2^2\right] \tag{6}$$

$$\mathscr{L}_{cls} = -\frac{1}{N}\sum_{i=1}^{N}\log\sigma_{y_i}(\mathbf{f}_i) \tag{7}$$

Where $\mathbf{f}i$ denotes the predicted logits for the $i$-th sample, and $\sigma y_i(\mathbf{f}_i)$ represents the predicted probability for the ground-truth class label $y_i$, obtained via the softmax function applied to $\mathbf{f}_i$. The total loss $\mathscr{L}$ adopts a weighted combination of the reconstruction loss and the classification loss, where $\lambda$ is a weighting hyperparameter used to balance.

TABLE II

COMPARISON ON **HUMANACT12** USING SKELETON-BASED ACTION RECOGNITION MODELS. RESULTS ARE REPORTED AS *mean ± std* OVER 5 INDEPENDENT RUNS; METHODS MARKED WITH * DENOTE MODELS TRAINED ON AUGMENTED DATA (REAL + SYNTHETIC). IMPROVEMENTS BROUGHT BY OUR AUGMENTED DATA ARE HIGHLIGHTED IN GREEN.

| Method | Real Data Usage | | | |
| --- | --- | --- | --- | --- |
| | 100% | 95% | 90% | 75% |
| STGCN++ [9] | 78.47 ±2.09 | 77.78 ±2.55 | 75.83 ±1.24 | 73.89 ±0.38 |
| STGCN++* | 83.19 ±2.73 (↑4.72) | 81.63 ±2.05 (↑3.85) | 81.50 ±1.47 (↑5.66) | 81.11 ±0.80 (↑7.22) |
| MSG3D [22] | 80.42 ±1.99 | 77.64 ±1.50 | 76.94 ±2.43 | 74.86 ±1.80 |
| MSG3D* | 83.11 ±3.46 (↑2.69) | 83.24 ±1.23 (↑5.60) | 81.77 ±1.18 (↑4.83) | 80.50 ±0.68 (↑5.64) |
| CTRGCN [5] | 77.78 ±1.97 | 76.94 ±2.10 | 75.56 ±1.42 | 73.61 ±2.41 |
| CTRGCN* | 79.42 ±2.02 (↑1.64) | 79.59 ±1.83 (↑2.65) | 80.16 ±2.20 (↑4.60) | 78.25 ±1.72 (↑4.64) |
| BlockGCN [44] | 77.78 ±1.30 | 75.67 ±1.30 | 75.56 ±0.76 | 75.56 ±0.90 |
| BlockGCN* | 78.91 ±0.41 (↑1.13) | 78.67 ±1.63 (↑3.00) | 78.19 ±0.38 (↑2.63) | 77.17 ±0.72 (↑1.61) |

TABLE III

COMPARISON ON THE **REFINED NTU RGB+D** DATASET USING SKELETON-BASED ACTION RECOGNITION MODELS. RESULTS ARE REPORTED AS *mean ± std* OVER 5 INDEPENDENT RUNS; METHODS MARKED WITH * DENOTE MODELS TRAINED ON AUGMENTED DATA (REAL + SYNTHETIC). IMPROVEMENTS FROM OUR AUGMENTED DATA ARE HIGHLIGHTED IN GREEN.

| Method | Real Data Usage | | | |
| --- | --- | --- | --- | --- |
| | 25% | 20% | 15% | 10% |
| STGCN++ [9] | 91.55 ±0.62 | 90.95 ±1.04 | 89.94 ±1.06 | 83.01 ±2.15 |
| STGCN++* | 92.36 ±0.33 (↑0.81) | 92.14 ±0.87 (↑1.18) | 92.07 ±0.76 (↑2.13) | 85.38 ±1.13 (↑2.37) |
| MSG3D [22] | 90.97 ±1.08 | 89.74 ±2.33 | 87.41 ±1.30 | 79.48 ±1.87 |
| MSG3D* | 92.30 ±0.39 (↑1.33) | 90.36 ±0.68 (↑0.62) | 89.90 ±1.59 (↑2.49) | 83.17 ±1.13 (↑3.69) |
| CTRGCN [5] | 90.81 ±1.07 | 90.78 ±0.20 | 87.57 ±2.69 | 79.28 ±1.46 |
| CTRGCN* | 91.13 ±1.34 (↑0.32) | 90.97 ±0.49 (↑0.19) | 89.45 ±0.35 (↑1.88) | 83.17 ±1.34 (↑3.89) |
| BlockGCN [44] | 90.03 ±0.72 | 88.51 ±1.11 | 86.70 ±1.46 | 75.05 ±1.43 |
| BlockGCN* | 90.91 ±0.54 (↑0.88) | 89.13 ±1.16 (↑0.62) | 86.05 ±1.42 (↓0.65) | 84.43 ±0.72 (↑9.38) |

$$\mathcal{L} = \mathcal{L}_{\text{rec}} + \lambda_{\text{cls}}\mathcal{L}_{\text{cls}} \qquad (8)$$

## IV. EXPERIMENTS AND RESULTS

### A. Datasets.

We evaluated our method on two benchmark datasets: HumanAct12 [11] and the Refined NTU-RGBD (NTU-VIBE) [33], [11].

- **HumanAct12** was a high-quality motion dataset derived from PHSPD [47], [46]. It contained 1,191 motion clips and over 90,000 frames across 34 fine-grained action categories. Actions included detailed labels such as *lift dumbbell with right hand* and *drink bottle left hand*, enabling conditional generation with strong label guidance.
- **Refined NTU-RGBD** was an improved version of the original NTU-RGBD [33] dataset, with 3D joint annotations recomputed using the VIBE [19] method for better consistency and realism. It included 3,902 motion

clips across 13 action categories, such as *squat down*, *sitting down*, and *throw*.

Although our experiments focus on HumanAct12 and NTU-VIBE, our method is architecture- and dataset-agnostic, and could be extended to other skeleton-based datasets with minimal modification.

### B. Usage Protocols.

For HumanAct12, we conducted experiments under **different data availability settings** by randomly sampling 75%, 90%, 95%, and 100% of the original training data to train the diffusion model. For downstream evaluation, we then augmented the selected real data with 5× synthetic samples generated by the diffusion model and tested the recognition accuracy on the validation set.

For Refined NTU-RGBD, we also considered **different data availability settings** in the few-shot regime, using only 10%, 15%, 20%, and 25% of the original training data to train the diffusion model. Each subset was similarly aug-

mented with 5× synthetic samples generated by our method during downstream evaluation. This demonstrates that even with limited real data, supplementing with label-consistent synthetic sequences can significantly improve recognition accuracy and approach the performance achieved with substantially larger real datasets.

We conclude that using 5× synthetic data represents a practical trade-off: generating substantially more data leads to redundancy and slows down the generation process, while too few synthetic samples provide only limited performance gains.

### C. Implementation Details.

Each skeleton sequence was loaded and filtered through preprocessing to ensure a minimum length. Motion data was normalised using the dataset-specific mean and standard deviation, and randomly cropped to a fixed-length window of $T = 48$ frames during training. Action labels were converted to 13-way or 34-way one-hot vectors based on predefined NTU and Humanact12 classes. We trained our method with Adam using a learning rate of $1 \times 10^{-4}$, which was decreased by 0.1 at each step. Training was conducted for 600 epochs on a single NVIDIA RTX 3090 with a batch size of 256. The Transformer encoder and decoder each consisted of 4 layers and 4 attention heads. For downstream action recognition, we used STGCN++[9], MSG3D[22], CTRGCN [5], and BlockGCN [44], applying their default configurations.

### D. Data Augmentation Evaluation

We evaluated performance by comparing the classification accuracy of the four state-of-the-art skeleton action recognition models, trained on real data only *and* reduced amounts of real data supplemented with our synthetic data. The results showed consistent gains, particularly in low-data settings. Tables II and III (HumanAct12 and NTU-VIBE, respectively) highlight accuracy gains achieved by augmenting the real data with our generated data. Performance is most pronounced when less real data is used, validating the value of our approach in few-shot contexts. On the HumanAct12 dataset, we observed that while using 95% or 90% of real data yields reasonable performance, training with only 75% real data augmented by our synthetic samples surpasses them, and even outperforms the model trained on 100% real data, demonstrating superior data efficiency and augmentation quality. On the NTU-VIBE dataset, BlockGCN with 15% data performs slightly worse, which may be due to overfitting from synthetic data. We observe consistent improvements across backbones with varying capacity — from lightweight STGCN++ to deeper BlockGCN — highlighting the generality of our approach.

### E. Data Distribution Evaluation

To examine whether our conditional diffusion augments the dataset in a label-consistent manner, we compare the distribution of *real* vs. *synthetic* samples using a t-SNE 2D projection. As shown in Figure 3, real samples are



**(a) HumanAct12 (N-Class) t-SNE visualization**
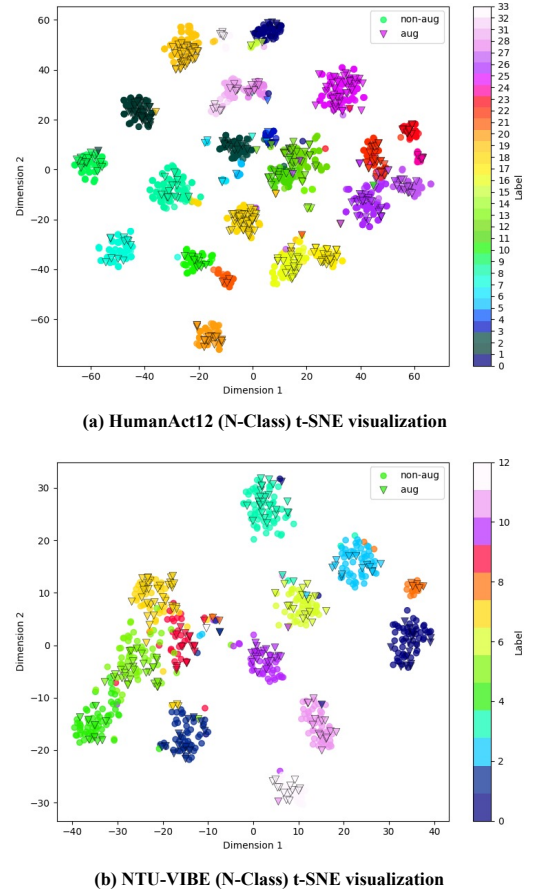


**(b) NTU-VIBE (N-Class) t-SNE visualization**

Fig. 3. t-SNE visualisations comparing real and synthetic skeleton samples of (a) HumanAct12 and (b) Refined NTU-RGBD (NTU-VIBE). Real samples are denoted by (∘), and synthetic samples are denoted by (▽).

shown as circles (∘), while synthetic samples are shown as inverted triangles (▽). We observe that the synthetic samples densely populate the existing class regions, effectively enlarging per-class coverage without shifting the class centroids, thereby preserving fidelity. At the same time, some synthetic points appear along the edges of the clusters, filling low-density areas that are under-represented in the real set; this expands intra-class boundaries while remaining label-consistent, yielding greater diversity and improved generalisation in few-shot regimes. To quantify this effect, we compute the *average within-class cluster covariance*, i.e., the trace of the per-class covariance matrix averaged across classes, which reflects intra-class dispersion. As shown in Table IV, covariance increases on both datasets when synthetic samples are added, indicating broader intra-class coverage while preserving label consistency. Overall, the visualisation and quantitative results indicate that our generator adds both quantity and diversity while maintaining clear inter-class separation.

### F. Reconstruction Evaluation

To evaluate skeleton generation quality, we compared our method with MDM [37] and T2M-GPT [43], using four metrics: FID (Fréchet Inception Distance), KID (Kernel

TABLE IV

COMPARISON OF AVERAGE WITHIN-CLASS CLUSTER COVARIANCE
BEFORE AND AFTER APPLYING DATA AUGMENTATION.

| Dataset | Original | +Aug | Difference |
|---|---|---|---|
| HumanAct12 | 10.926 | 11.490 | +0.564 |
| NTU-VIBE | 4.392 | 4.753 | +0.361 |

TABLE V

ABLATION STUDY OF STGCN++ RESULTS ON THE HUMANACT12
DATASET USING THE VALIDATION SET, TRAINED WITH 100% OF THE
DATA, UNDER DIFFERENT COMBINATIONS OF MODULES. RESULTS ARE
REPORTED AS *mean ± std* OVER 5 INDEPENDENT RUNS.

| Module | Condition | CLS Loss | Dropout | Refinement | STGCN++ Acc. |
|---|---|---|---|---|---|
| Baseline | ✗ | ✗ | ✗ | ✗ | 78.77 ±2.65 |
| | ✓ | ✗ | ✗ | ✗ | 80.62 ±1.58 |
| | ✓ | ✓ | ✗ | ✗ | 81.98 ±2.56 |
| | ✓ | ✓ | ✓ | ✗ | 80.77 ±2.67 |
| All (Ours) | ✓ | ✓ | ✓ | ✓ | **83.19** ±2.73 |

Inception Distance), Diversity, and Precision/Recall. Detailed explanations of these metrics were provided in the supplementary materials. As shown in Table IX, our method achieves the lowest FID and comparable KID, indicating that the generated motions are most similar to real data regarding overall distribution and visual coherence. The highest diversity among generative models demonstrates a strong ability to produce a wide range of motion styles rather than repetitive patterns. While T2M-GPT slightly outperforms Precision, suggesting high realism in individual samples, our method maintains a strong balance across all metrics. These results indicate that our approach generates realistic motions and captures a broader spectrum of plausible human movements, outperforming prior methods in fidelity and diversity.

## G. Ablation Study

*1) Evaluation of proposed modules:* We conduct an ablation study by incrementally adding each module to STGCN++ using HumanAct12 with 100% data usage: As shown in Table V, incorporating condition embedding and classification loss yielded an accuracy improvement of 2.35%. Although Sampling Dropout enhanced diversity, it

TABLE VI

SKELETON-BASED ACTION RECOGNITION PERFORMANCE
ON HUMANACT12 UNDER DIFFERENT DROPOUT RATES.
GREEN CELLS INDICATE THE BEST PERFORMANCE PER
COLUMN.

| Ratio | 100% | 95% | 90% | 75% |
|---|---|---|---|---|
| Dropout 0 | 83.85 | 84.62 | 80.77 | 83.85 |
| Dropout 0.1 | 82.31 | 82.31 | **83.85** | 80.00 |
| Dropout 0.2 | **86.15** | **86.15** | 80.00 | **84.62** |
| Dropout 0.5 | 84.62 | 84.62 | 83.08 | 81.54 |

also introduced fidelity degradation issues when used independently. However, this issue was effectively mitigated by the Generative Refinement Module, which enabled our method to maintain diversity and accuracy, as evidenced by the highest performance of 83.19% achieved using all proposed components.

*2) Dropout and Renoise Strategy Comparison:* We evaluated the impact of different dropout rates and GRM renoise values on model performance. As shown in Table VI, we tested four dropout settings during the sampling stage: no dropout, 0.1, 0.2, and 0.5. The results indicate that higher dropout values increase diversity but also lead to a loss of fine-grained details. A dropout rate of 0.2 achieved the best performance for downstream skeleton-based action recognition with STGCN++ during sampling, suggesting that introducing a moderate dropout value can enhance the diversity of generated skeletal motions while also improving recognition accuracy.

As shown in Table VII, we further examined the Generative Refinement Module (GRM) under different renoise values. Here, renoise refers to a threshold such that generated samples with deviations larger than this value from the original data are discarded. On the HumanAct12 dataset, using renoise values between 10 and 20 proved effective in filtering out distorted samples, thereby reducing their negative impact on downstream action recognition tasks. Notably, in the downstream main results, we adopt unified dropout and GRM values and report the average over five repeated runs, which yields more robust results across tasks. Although both sets of parameters require manual tuning, they exhibit a certain degree of robustness: even when the chosen values are not optimal, the model still surpasses the baseline in downstream performance.

*3) Comparison of Data Augmentation Methods:* We compared the effectiveness of different data augmentation methods on downstream tasks in Table VIII. Unlike conventional methods, which indiscriminately perturb motion data without considering semantic consistency, our conditional generative approach produces label-consistent and realistic motion variations, enhancing diversity and performance, particularly in low-data settings.

TABLE VII

SKELETON-BASED ACTION RECOGNITION PERFORMANCE
ON HUMANACT12 UNDER DIFFERENT GRM RE-NOISE
VALUES. GREEN CELLS INDICATE THE BEST PERFORMANCE
PER COLUMN.

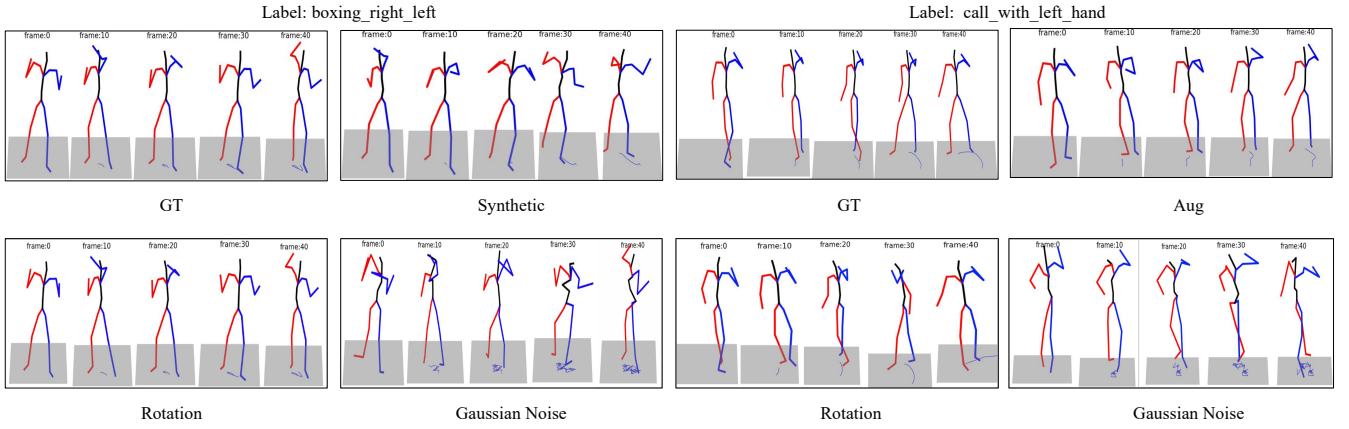| Ratio | 100% | 95% | 90% | 75% |
|---|---|---|---|---|
| Renoise 1 | 80.77 | 83.08 | 81.54 | 85.38 |
| Renoise 2 | 83.08 | 83.08 | **84.62** | 82.31 |
| Renoise 3 | 83.08 | 81.54 | 82.31 | 81.54 |
| Renoise 5 | 82.31 | 81.54 | 83.08 | 78.46 |
| Renoise 10 | 83.85 | **84.62** | 80.77 | 80.00 |
| Renoise 20 | **85.38** | 83.85 | 80.77 | **87.69** |

Fig. 4. Visualisation of our generated skeleton sequences conditioned on action labels. Our results demonstrate that our method generates diverse motion patterns while preserving label-specific semantics. Additional visualisations are provided in the supplementary file.

TABLE VIII

ACCURACY OF DIFFERENT DATA AUGMENTATION METHODS ON HUMANACT12 FOR SKELETON-BASED ACTION RECOGNITION.

| Ratio | 100% | 95% | 90% | 75% |
|---|---|---|---|---|
| W/O Augmentation | 75.69 | 77.78 | 75.00 | 73.61 |
| Gaussian Noise | 78.47 | 79.17 | 79.86 | 79.86 |
| Scaling | 79.17 | 77.08 | 75.00 | 72.22 |
| Rotating | 79.17 | 77.08 | 79.86 | 75.00 |
| Ours | **81.54** | **80.77** | **80.00** | **80.56** |

TABLE IX

COMPARISON OF SKELETON GENERATION QUALITY ON THE HUMANACT12 DATASET.

| Method | FID ↓ | KID ↓ | Diversity ↑ | Precision↑ | Recall↑ |
|---|---|---|---|---|---|
| Real data | 0.8398 | 0.0001 | 7.217 | 0.996 | 0.999 |
| MDM-orig [37] | 11.3120 | 0.0512 | 6.223 | 0.996 | 0.896 |
| MDM | 24.9942 | 0.1168 | 3.7580 | 0.998 | 0.390 |
| T2M-GPT[43] | 2.0362 | **0.0057** | 6.6808 | **0.999** | 0.980 |
| **Ours** | **1.3288** | 0.0170 | **6.8087** | 0.996 | **0.994** |

*4) Qualitative Results:* Figure 4 presents visualisations of conditional skeletal generation results. By analysing these qualitative results, we highlighted the diversity and fidelity of the generated samples. The results illustrated how our method preserved label-specific semantics while introducing subtle variations in key joints relevant to action recognition, in contrast to conventional skeleton data augmentation methods that mainly rely on simple geometric transformations. Moreover, the visualisations indicate that our generated data carries clear physical meaning, such as variations in movement speed and joint angles, since the 263-dimensional inputs encode rich physical attributes rather than mere skeleton points. Additional visualisations are provided in the supplementary material.

## V. CONCLUSION

We presented a conditional diffusion framework for skeleton-based action recognition, generating diverse, label-consistent motion sequences. Our approach significantly improves recognition performance, especially under limited data conditions, and consistently benefits a range of skeleton action recognition backbone architectures. We demonstrated how our design balances fidelity and diversity through extensive ablations, enabling scalable and controllable augmentation. This work represents a significant step forward in generative augmentation for structured human motion data, with practical implications for data-efficient learning in action recognition tasks and for reducing the cost of collecting large-scale skeleton data.

*a) Limitations and Future Work:* While our method generalises well across datasets and models, generation quality may degrade for rare or ambiguous actions under extreme label imbalance. As future work, we will investigate adaptive sampling and uncertainty-aware conditioning to enhance robustness and extend the model to multi-person interactions and longer motion sequences.

## REFERENCES

[1] A. Antoniou, A. Storkey, and H. Edwards. Data augmentation generative adversarial networks. In *International Conference on Artificial Neural Networks (ICANN)*, pages 594–603, 2018.

[2] A. Brock, J. Donahue, and K. Simonyan. Large scale GAN training for high fidelity natural image synthesis. In *International Conference on Learning Representations (ICLR)*, 2019.

[3] Z. Cao, G. Hidalgo Martinez, T. Simon, S. Wei, and Y. A. Sheikh. Openpose: Realtime multi-person 2d pose estimation using part affinity fields. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2019.

[4] X. Chen, B. Jiang, W. Liu, Z. Huang, B. Fu, T. Chen, and G. Yu. Executing your commands via motion diffusion in latent space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 18000–18010, 2023.

[5] Y. Chen, Z. Zhang, C. Yuan, B. Li, Y. Deng, and W. Hu. Channel-wise topology refinement graph convolution for skeleton-based action recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 13359–13368, 2021.

[6] M. Cormier, Y. Schmid, and J. Beyerer. Enhancing skeleton-based action recognition in real-world scenarios through realistic data augmentation. In *Proceedings of the IEEE/CVF Winter Conference on*

*Applications of Computer Vision (WACV) Workshops*, pages 290–299, 2024.

[7] R. Dabral, M. H. Mughal, V. Golyanik, and C. Theobalt. Mofusion: A framework for denoising-diffusion-based motion synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9760–9770, 2023.

[8] P. Dhariwal and A. Nichol. Diffusion models beat gans on image synthesis. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 8780–8794, 2021.

[9] H. Duan, J. Wang, K. Chen, and D. Lin. Pyskl: Towards good practices for skeleton action recognition. In *Proceedings of the 30th ACM International Conference on Multimedia (ACM MM)*, pages 7351–7354, 2022.

[10] C. Guo, S. Zou, X. Zuo, S. Wang, W. Ji, X. Li, and L. Cheng. Generating diverse and natural 3d human motions from text. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5152–5161, 2022.

[11] C. Guo, X. Zuo, S. Wang, S. Zou, Q. Sun, A. Deng, M. Gong, and L. Cheng. Action2motion: Conditioned generation of 3d human motions. In *Proceedings of the 28th ACM International Conference on Multimedia (ACM MM)*, pages 2021–2029, 2020.

[12] J. Ho, A. Jain, and P. Abbeel. Denoising diffusion probabilistic models. In *Conference on Neural Information Processing Systems (NeurIPS)*, page 6840–6851, 2020.

[13] J. Ho and T. Salimans. Classifier-free diffusion guidance. In *Deep Generative Models and Downstream Applications Workshop at the 35th Conference on Neural Information Processing Systems (NeurIPS)*, 2022.

[14] J.-F. Hu, W.-S. Zheng, J. Lai, and J. Zhang. Jointly learning heterogeneous features for rgb-d activity recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5344–5352, 2015.

[15] T. Huynh-The, C.-H. Hua, and D.-S. Kim. Encoding pose features to images with data augmentation for 3-d action recognition. *IEEE Transactions on Industrial Informatics (TII)*, 16(5):3100–3111, 2019.

[16] C. Ionescu, D. Papava, V. Olaru, and C. Sminchisescu. Human3.6m: Large scale datasets and predictive methods for 3d human sensing in natural environments. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 36(7):1325–1339, 2014.

[17] A. Jahanian, X. Puig, Y. Tian, and P. Isola. Generative models as a data source for multiview representation learning. In *International Conference on Learning Representations (ICLR)*, 2022.

[18] J. Kim, J. Kim, and S. Choi. Flame: Free-form language-based motion synthesis editing. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, pages 12345–12355, 2023.

[19] M. Kocabas, N. Athanasiou, and M. J. Black. Vibe: Video inference for human body pose and shape estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5253–5263, 2020.

[20] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Conference on Neural Information Processing Systems (NeurIPS)*, pages 1097–1105, 2012.

[21] X. Liu, Z. Feng, D. Kanojia, and W. Wang. DGFM: Full body dance generation driven by music foundation models. In *Audio Imagination: NeurIPS 2024 Workshop on AI-Driven Speech, Music, and Sound Generation*, 2024.

[22] Z. Liu, H. Zhang, Z. Chen, Z. Wang, and W. Ouyang. Disentangling and unifying graph convolutions for skeleton-based action recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 143–152, 2020.

[23] M. Loper, N. Mahmood, J. Romero, G. Pons-Moll, and M. J. Black. Smpl: A skinned multi-person linear model. *ACM Transactions on Graphics (TOG)*, 34(6):248:1–248:16, 2015.

[24] F. Meng, H. Liu, Y. Liang, J. Tu, and M. Liu. Sample fusion network: An end-to-end data augmentation network for skeleton-based human action recognition. *IEEE Transactions on Image Processing (TIP)*, 28(11):5281–5295, 2019.

[25] A. Nichol and P. Dhariwal. Improved denoising diffusion probabilistic models. In *International Conference on Machine Learning (ICML)*, page 8162–8171, 2021.

[26] J. Park, B. Kim, and J. Jeong. An analysis of synthetic data for improving performance of skeleton-based fall down detection models. In *5th International Conference on Big Data Analytics and Practices (IBDAP)*, pages 89–92, 2024.

[27] P. Provini, A. L. Camp, and K. E. Crandell. Emerging biological insights enabled by high-resolution 3d motion data: promises, perspectives and pitfalls. *Journal of Experimental Biology*, 226:jeb245138, 2023.

[28] A. Ramesh, P. Dhariwal, A. Nichol, C. Chu, and M. Chen. Hierarchical text-conditional image generation with clip latents, 2022.

[29] Z. Ren, S. Huang, and X. Li. Realistic human motion generation with cross-diffusion models. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 345–362, 2024.

[30] Z. Ren, Z. Pan, X. Zhou, and L. Kang. Diffusion motion: Generate text-guided 3d human motion by diffusion model. In *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5, 2023.

[31] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10684–10695, 2022.

[32] Y. Shafir, G. Tevet, R. Kapon, and A. H. Bermano. Human motion diffusion as a generative prior. In *International Conference on Learning Representations (ICLR)*, 2023.

[33] A. Shahroudy, J. Liu, T.-T. Ng, and G. Wang. Ntu rgb+d: A large scale dataset for 3d human activity analysis. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1010–1019, 2016.

[34] J. Shen, J. Dudley, and P. O. Kristensson. The imaginative generative adversarial network: Automatic data augmentation for dynamic skeleton-based hand gesture and human action recognition. In *Proceedings of the IEEE International Conference on Automatic Face and Gesture Recognition (FG)*, pages 1–8, 2021.

[35] J. Sohl-Dickstein, E. A. Weiss, N. Maheswaranathan, and S. Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *Proceedings of the 32nd International Conference on Machine Learning (ICML)*, pages 2256–2265, 2015.

[36] J. Song, C. Meng, and S. Ermon. Denoising diffusion implicit models. In *International Conference on Learning Representations (ICLR)*, 2021.

[37] G. Tevet, S. Raab, B. Gordon, Y. Shafir, D. Cohen-or, and A. H. Bermano. Human motion diffusion model. In *International Conference on Learning Representations (ICLR)*, 2023.

[38] B. Trabucco, K. Doherty, M. Gurinas, and R. Salakhutdinov. Effective data augmentation with diffusion models. In *International Conference on Learning Representations (ICLR)*, 2024.

[39] J. Tu, H. Liu, F. Meng, M. Liu, and R. Ding. Spatial-temporal data augmentation based on lstm autoencoder network for skeleton-based human action recognition. In *IEEE International Conference on Image Processing (ICIP)*, pages 3478–3482, 2018.

[40] D. Wu and L. Shao. Leveraging hierarchical parametric networks for skeletal joints based action segmentation and recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 724–731, 2014.

[41] C. Xin, S. Kim, Y. Cho, and K. S. Park. Enhancing human action recognition with 3d skeleton data: A comprehensive study of deep learning and data augmentation. *Electronics*, 13(4), 2024.

[42] S. Yan, Y. Xiong, and D. Lin. Spatial temporal graph convolutional networks for skeleton-based action recognition. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, pages 7444–7452, 2018.

[43] J. Zhang, Y. Zhang, X. Cun, S. Huang, Y. Zhang, H. Zhao, H. Lu, and X. Shen. T2m-gpt: Generating human motion from textual descriptions with discrete representations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5152–5161, 2023.

[44] Y. Zhou, X. Yan, Z.-Q. Cheng, Y. Yan, Q. Dai, and X.-S. Hua. Blockgcn: Redefining topology awareness for skeleton-based action recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10181–10191, 2024.

[45] W. Zhu, X. Ma, D. Ro, H. Ci, J. Zhang, J. Shi, F. Gao, Q. Tian, and Y. Wang. Human motion generation: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, pages 1–20, 2023.

[46] S. Zou, X. Zuo, Y. Qian, S. Wang, C. Guo, C. Xu, M. Gong, and L. Cheng. Polarization human shape and pose dataset. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 1–17, 2020.

[47] S. Zou, X. Zuo, Y. Qian, S. Wang, C. Xu, M. Gong, and L. Cheng. 3d human shape reconstruction from a polarization image. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 1–17, 2020.

[48] S. Zou, X. Zuo, S. Wang, Y. Qian, C. Guo, and L. Cheng. Human pose and shape estimation from single polarization images. *IEEE Transactions on Multimedia (TMM)*, 25(12):3560–3572, 2023.