037

038

039

043

044

055

057

058

059

060

AG-MAE: Anatomically Guided Spatio-Temporal Masked Auto-Encoder for **Online Hand Gesture Recognition**

Anonymous 3DV submission

Paper ID 60

Abstract

001 Hand gesture recognition plays a crucial role in the do-002 main of computer vision, as it enhances human-computer interaction by enabling intuitive, touch-free control and 003 004 communication. While offline methods have made significant advances in isolated gesture recognition, real-005 006 world applications demand online and continuous process-007 ing. Skeleton-based methods, though effective, face chal-800 lenges due to the intricate nature of hand joints and the diverse 3D motions they induce. This paper introduces 009 010 AG-MAE, a novel approach that integrates anatomical constraints to guide the self-supervised training of a spatio-011 temporal masked autoencoder, enhancing the learning of 012 3D keypoint representations. By incorporating anatomical 013 knowledge, AG-MAE learns more discriminative features 014 015 for hand poses and movements, subsequently improving on-016 line gesture recognition. Evaluation on standard datasets demonstrates the superiority of our approach and its po-017 018 tential for real-world applications. Code is available at: https://github.com/lambda-xyz-01/AGMAE. 019

1. Introduction 020

Online recognition of dynamic hand gestures plays an es-021 022 sential role in computer vision, human-computer interaction (HCI) and virtual reality (VR) applications, enabling 023 024 seamless, intuitive and natural interactions between users 025 and machines. Unlike traditional offline gesture recognition 026 systems [15, 20, 25], which focus on discrete segmented gestures, the framework of continuous dynamic gesture 027 recognition requires interpreting hand movements in a con-028 tinuous stream of data, enabling real-time interactions and 029 feedbacks. 030

031 Online gesture recognition raises significant challenges due to the intricate nature of hand movements and the di-032 verse range of motions (ROM) occurring within a contin-033 034 uous flow of non-segmented gestures. Unlike offline sce-035 narios, online recognition demands precise localization of



Figure 1. Hand models (a) without, (b) with anatomical constraints: joint angles (θ), bone lengths (l), finger curvature (κ).

gestures within this continuous flow, necessitating accurate identification of their start and end timings. Furthermore, real-world applications necessitate real-time processing, implying rapid inference algorithms without compromising accuracy. Ensuring high precision in classification 040 while minimizing false positives is crucial for ensuring a 041 natural and reliable interaction experience, particularly in 042 critical scenarios such as medical operations.

Recent advancements in self-supervised learning have shown promise in deriving discriminative representations 045 from unlabeled hand pose data [8, 23, 24, 42]. We ar-046 gue that this approach holds significant potential, particu-047 larly in the domain of dynamic and continuous hand ges-048 ture recognition. Therefore, we propose a method that com-049 bines the power of self-supervised learning with anatomi-050 cal constraints guidance to overcome the limitations inher-051 ent in traditional fully supervised approaches. By integrat-052 ing self-supervised skeletal learning and anatomical infor-053 mation during pre-training, we aim to extract rich and dis-054 criminative representations of hand poses. As illustrated in Figure 1, anatomical constraints such as bone length, bone 056 curvature and joint angles can be incorporated as a prior information into the learning model to ensure consistency in hand joint position estimation. Ultimately, improving the discrimination of learned 3D keypoint representations.

Existing few self-supervised learning methods [8, 23, 061 24] often prioritize model accuracy in matching ground 062

126

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

truth hand poses while neglecting anatomical correctness.
By incorporating anatomical constraints into the learning process, as evidenced in prior studies on hand pose estimation [26, 38] and hand tracking [1], we demonstrate improved model capabilities in learning richer representations of various hand poses and movements.

The following are the primary contributions of our work:

A spatio-temporal ViT-based model with Fourier embedding: We integrate Fourier feature embedding [40]
into a spatio-temporal vision transformer (ViT) model
to project spatial and temporal coordinates into a highfrequency domain. This enhancement captures intricate
spatial and temporal dependencies and nuanced patterns
in hand joint data, improving representation accuracy.

 Anatomical guidance for pre-training: We introduce anatomical constraints into the loss function to guide the pre-training of the masked autoencoder, ensuring anatomical consistency and learning discriminative features for various hand poses.

082 2. Related Work

083 2.1. Online Hand Gesture Recognition

Online gesture recognition methods extend beyond the 084 085 scope of offline methods, which primarily focus on discrete, 086 segmented gestures. In contrast, online gesture recogni-087 tion involves two key tasks: segmenting the continuous data stream to identify the start and end frames of each gesture, **088** and accurately labeling these gestures using prior informa-089 tion while minimizing delays and avoiding false positives. 090 091 Online recognition of hand gestures has been approached 092 through two main methods:

093 Full sequence-based methods: analyze an entire se-094 quence at once to detect gesture boundaries before forwarding the identified segmented candidates to a classification 095 096 module. They employ specialized heuristics based on ve-097 locity, energy, or trained networks to segment sequence and 098 subsequently classify each frame subset. Traditional methods utilized the Histogram of Oriented Gradient (HOG) al-099 gorithm in conjunction with an SVM classifier [33]. In con-100 trast, recent advancements have predominantly focused on 101 102 time-driven models. Köpüklü et al. [28] proposed a two-103 model hierarchical architecture based on lightweight CNNs. 104 Seg-LSTM [7] employs an LSTM with a specialized segmentation network, while the ST-GCN method [6] utilizes 105 106 an energy-based segmentation approach with additional adhoc rules. The 2ST-GCN method [2, 6, 17] integrates an 107 108 energy-based detection module with a fine-grained classi-109 fier for gesture/non-gesture discrimination.

Sliding window-based methods: perform continuous
and simultaneous detection and labeling, often using pretrained classifiers with fixed-size input subsequences and
sliding-window models. Sliding window techniques are

common, as shown by the [16] strategy, where a modified 114 DDNet [46] is trained with segmented and resampled ges-115 tures and randomly sampled non-gestural windows. Sim-116 ilarly, a modified version of DeepGRU [32] demonstrated 117 notable performance. TN-FSM [17] uses transform net-118 works to classify 10-frame windows, while Causal TCN 119 trains a temporal convolutional network on 20-frame win-120 dows labeled with gesture classes or non-gestures accord-121 ing to their intersections with the annotated ground truth 122 [16, 17]. In addition, OO-dMVMT [12] exploits multiple 123 temporal views of hand pose and movement to generate 124 complete gesture descriptions. 125

2.2. Skeleton-based Self-Supervised Learning

Self-supervised learning has primarily been successful in
image analysis, especially due to the emergence of masked
autoencoders (MAEs) [22], which have been proven suc-
cessful in a variety of applications [3, 9, 22]. Accordingly,
the field of skeletal data has recently seen a growing interest
in exploiting the potential of self-supervised learning.127
128
129

Contrastive learning methods [29, 34] apply momentum encoders for contrastive learning using single-stream skeleton sequences. Aiming for more generalized representations, AimCLR [21] implemented an extreme data augmentation strategy to increase the number of contrastive pairs and thus improve feature extraction. To prevent overfitting and improve feature generalization for action recognition, Ms2l [31] introduced a multitasking self-supervised framework that focuses on the extraction of joint representations via motion prediction and puzzle recognition.

MAE-based methods have received considerable attention. D-MAE [27] introduced a dual MAE focusing on token completion in a skeletal context, crucial for robust motion capture. Similarly, SkeletonMAE [44] proposed a graph-based MAE, emphasizing pre-training with skeleton sequences. Generative learning techniques such as LongT GAN [48] and P&C [39] emphasized encoder-decoder architectures to refine skeleton sequence representation.

Despite advances, self-supervised learning in hand gesture recognition, especially online, remains underexplored. Chen et al. [8] focused on 3D hand reconstruction, Sign-BERT [23] pre-trained hand-aware representations for sign language, while [24] pre-trained a MAE to encode separate individual hand poses without considering temporal correlation. Our work extends self-supervised skeleton learning to improve online gesture recognition, by incorporating spatio-temporal encoding and relying on prior knowledge and anatomical constraints to inform the learning process.

3. Methodology

We propose a comprehensive end-to-end framework for online hand gesture recognition, which is divided into two main phases. First, a spatio-temporal MAE (STMAE) is 164



Figure 2. (A) **Proposed AG-MAE:** a ratio of joints are masked in a given window, the unmasked joints are encoded by the encoder and then concatenated with the mask tokens and passed through the decoder to reconstruct the masked joints. (B) Masking strategies.

pre-trained to encode a sequence of skeletal hand gesture
frames into a robust feature representation. Subsequently,
a spatio-temporal graph convolutional network (STGCN) is
fine-tuned to classify gestures within a real-time data stream
using the learned representations. Figure 2-A illustrates
the architecture of our proposed bio-mechanically guided
spatio-temporal masked autoencoder (AG-MAE).

172 3.1. Pretraining

In pre-training, given an input window of hand poses $X \in$ 173 $\mathbb{R}^{W \times N \times 3}$, where W is the number of frames, N is the 174 number of hand joints, and 3 corresponds to the 3D coor-175 176 dinates (x, y, z), we first project the 3D joint coordinates into a higher-dimensional space \mathbb{R}^d using a Fourier embed-177 ding map, while incorporating positional encoding to main-178 tain spatial and temporal order. A ratio m_r (mask ratio) of 179 joints is masked according to one of the strategies in Fig-180 ure 2-B. Each joint is represented as a token of dimension 181 d in the Fourier embedding space. The unmasked joints are 182 processed by a ViT-based MAE encoder, mapping them to 183 a latent space \mathbb{R}^l . The encoded unmasked joints are con-184 catenated with the mask tokens and fed to the ViT-based 185 MAE decoder to reconstruct the masked joint coordinates, 186 producing $\tilde{X} \in \mathbb{R}^{W \times N \times 3}$. Reconstruction quality is evalu-187 ated using the mean squared error, along with the anatom-188 ical loss, which assesses the anatomical correctness of the 189 reconstructed hand poses. This process enhances the en-190 coding of hand poses window into a more discriminative 191 192 feature space, enhancing their utility for subsequent tasks 193 such as gesture spotting and classification.

Fourier Embedding. Fourier Feature Embedding (FFE) 194 improves the ability of the model to capture spatial and tem-195 poral relationships between hand joints. It projects spatial 196 and temporal coordinates into a high-frequency domain us-197 ing sine and cosine functions of varying frequencies. This 198 technique allows the model to discern nuanced patterns 199 in 3D keypoint motions [40]. Unlike linear embeddings, 200 which may overlook fine details, FFE preprocesses the in-201 put to capture higher-frequency details and intricate spa-202 tial dependencies, leading to more accurate representations 203 [24, 40]. The FFE embeds the 3D coordinates v(x, y, z)204 into a 256-dimensional space: 205

$$\gamma(v) = [a_1 \cos(2\pi b_1^T v), \ a_1 \sin(2\pi b_1^T v),
\dots, \dots, \dots, (1) 206
a_m \cos(2\pi b_m^T v), \ a_m \sin(2\pi b_m^T v)]^T$$

where b are the Fourier basis frequencies, and a are the corresponding Fourier series coefficients, resulting in a feature transformation with m distinct frequency components. 209

Positional Encoding. Positional encoding aims to pre-210 serve both spatial and temporal dimensions within the data. 211 Specifically, a spatial positional encoding is added to each 212 joint and maintained across all frames to retain the spatial 213 structure. Additionally, a temporal positional encoding is 214 applied to each frame, with the same encoding assigned 215 to all joints within a frame to ensure temporal consistency. 216 These encodings enable the model to effectively track and 217 correlate spatial and temporal relationships. 218



256

257

259

260

261

262

263

264

277

278

279

280

281

282

284

285

286

287

288



Figure 3. Hand anatomy constrains the biomechanics of hand motions, including joint angles (θ) and bone lengths (l).

219 Masking Strategy. Different masking strategies, illustrated in Figure 2-B, are employed to enhance the self-220 221 supervised learning model for characterizing hand poses.

- 222 Random spatial masking: is a joint-level masking strategy 223 that involves masking a given ratio of the same joints over time, *i.e.* the same set of random joints is masked in each 224 frame of the sequence. 225
- Random temporal masking: is a frame-level masking 226 strategy that involves masking a random number of 227 228 frames in the sequences, *i.e.* all joints of the hand are 229 masked in a given random set of frames.
- 230 • Random spatio-temporal masking: is a widely adopted and highly effective strategy in image- and skeleton-231 232 based self-supervised learning [22, 44] involving ran-233 domly masking a number of joints at both the frame- and 234 joint-level in the sequence.

3.2. Anatomical Constraints 235

The hand is anatomically constrained by biomechanical 236 limitations [11], allowing it to perform certain poses while 237 238 limiting its ROM. Each joint has a specific degree of free-239 dom (DoF) that defines its movement capabilities. For example, the index, middle, ring, and little fingers are consid-240 ered planar manipulators, meaning that their DIP, PIP, and 241 MCP joints move primarily in one plane since the DIP and 242 243 PIP joints only have 1 DoF for flexion (see Figure 3). The anatomical constraints can be categorized into two primary 244 245 categories: dynamic and static constraints.

Dynamic constraints can be subdivided into intrafinger 246 and interfinger constraints. Intrafinger constraints refer to 247 248 limitations on movement between different joints within the same finger [13]. For instance, Cobos et al. [11] outlined 249 several constraints, such as the requirement that to bend the 250 DIP joints, the PIP joints must also be bent for the index, 251 middle, ring, and little fingers, mathematically expressed as 252 $\theta_{DIP} = \frac{2}{3} \theta_{PIP}$. While these constraints are not rigid, indi-253 254 viduals generally adhere to them under normal conditions,

though there is some variation in the ability to control specific joints across individuals.

Interfinger constraints involve correlations between joints across different fingers, often resulting in coupled 258 movements among fingers [30]. For example, when the pinky finger bends, the ring finger also bends to a certain extent, reflecting a proportional relationship. However, variations exist among individuals regarding these constraints. Some constraints can be overcome, while others are inherent and cannot be explicitly represented in equations [38].

Static constraints define the normal ROM for hand 265 joints, setting limits on parameter values in models. These 266 constraints [11], provide crucial guidelines for understand-267 ing and modeling hand biomechanics. Despite individual 268 variations, static constraints play a significant role in defin-269 ing the anatomical capabilities of the hand. The limits for 270 each constraint can be obtained manually from measure-271 ments, from the literature (e.g., [10, 35]), or acquired in a 272 data-driven way from 3D annotations. Two main constraints 273 are commonly considered [11], as illustrated in Figure 3: 274 bone lengths, reflecting intra-finger constraints, and joint 275 angles, covering both intra- and inter-finger constraints. 276

For bone lengths, we define an interval $[b_i^{\min}, b_i^{\max}]$ for each bone *i* and penalize deviations if the length $||b_i||_2$, which corresponds to the Euclidean distance between the extremities of the bone at the joints, lies outside this interval. Mathematically, given a hand pose P, we define the bone length loss as:

$$\mathcal{L}_{BL}(P) = \frac{1}{N_b} \sum_{i=1}^{N_b} \mathcal{I}(\|b_i\|_2; b_i^{\min}, b_i^{\max})$$
(2) 283

where N_b is the number of bones, and $\mathcal{I}(||b_i||_2; b_i^{\min}, b_i^{\max})$ is an indicator function that penalizes the bone length $||b_i||_2$ if it falls outside the defined interval $[b_i^{\min}, b_i^{\max}]$.

For joint angles, each joint has its specific range of freedom. For instance, as cited by Cobos et al. [11]:

$$0^{\circ} \le \theta_{\text{MCP}} \le 90^{\circ}; \quad 0^{\circ} \le \theta_{\text{PIP}} \le 110^{\circ}; \quad 0^{\circ} \le \theta_{\text{DIP}} \le 90^{\circ}.$$
 289

We propose to compute the ranges of angles for each joint 290 based on the data. To compute the joint angles, we first need 291 to define a reference point relative to which the angles are 292 measured. The wrist joint appears to be the most suitable 293 as it has 0 DoF in the hand plane, and its movements are 294 minimal or almost negligible. 295

To constrain the angles, we consider each angle indepen-296 dently (e.g., θ_1 between the Index MCP and Middle MCP 297 in Figure 3) and penalize them if they lie outside the cor-298 responding interval. This corresponds to constraining them 299 within a box in a 2D space, where the endpoints are the 300 min/max limits. The angles are constrained to lie within this 301 structure by minimizing their distance to it. For angles, we 302 consider the angles between all pairs of joints in the hand, 303

331

372

373

374

375

376

377

378

even those within the same hand, leading to the followingdefinition of the loss with regards to angles:

$$\mathcal{L}_{JA}(P) = \frac{2}{N(N-1)} \sum_{i=2}^{N} \sum_{j=i+1}^{N} \mathcal{I}(\angle(\overrightarrow{J_w J_i}, \overrightarrow{J_w J_j}); a_{ij}^{\min}, a_{ij}^{\max}),$$
(3)

307 where $\angle(\overline{J_wJ_i}, \overline{J_wJ_j})$ is the angle between joint *i* and joint 308 *j* considering the wrist joint (J_w) as the vertex, $a_{ij}^{min} \in$ 309 A^{min} and $a_{ij}^{max} \in A^{max}$ are the minimum and maximum 310 angle values between joint *i* and *j*. The number of pairs 311 (J_w, J_i) and (J_w, J_j) where *i* and *j* are distinct joints in-312 dices is given by the binomial coefficient: $\binom{N}{2} = \frac{N(N-1)}{2}$. 313 Given the potential for error and inaccuracy in annotated

313 314 hand gesture datasets, the lack of a common hand kinemat-315 ical model across dataset and the lack of explicit equations 316 for some dynamic angles, we propose to rely primarily on 317 static constraints. These constraints offer a more accessible 318 and straightforward approach to defining ROMs for individual joints and the hand as a whole. The ROM is character-319 ized by its minimum and maximum values, determined by 320 the angles between the various joints, with the wrist joint 321 322 serving as the reference point, as well as by considerations of the distances between these joints. We argue that con-323 sidering the angles between all pairs of joints in the hand 324 325 provides the model with a rich feedback on the anatomical 326 correctness of the hand.

Integrating these constraints into the loss function encourages keypoint predictions that yield valid bone lengths
and valid angles, thus ensuring accurate hand anatomy. The
anatomical loss can be formulated as:

$$\mathcal{L}_A(P) = \mathcal{L}_{BL}(P) + \mathcal{L}_{JA}(P) \tag{4}$$

where \mathcal{L}_{BL} and \mathcal{L}_{JA} denote the bone length and joint angle losses, respectively.

334 3.3. Model Architecture

The AG-MAE model is designed to process temporal hand
skeleton data. It is based on an asymmetric encoder-decoder
architecture, both built upon the ViT model [41].

Encoder. The encoder is built to encode a given window of Fourier embedded hand non-masked tokens $X_{nmask}^F =$ $\gamma(X) \in \mathbb{R}^{N_{nmask} \times 256}$, where $X_{nmask} \in \mathbb{R}^{N_{nmask} \times 3}$ is the window of hand poses and N_{nmask} is the number of nonmasked joints across the window, into a latent space $X_{enc} \in$ $\mathbb{R}^{N_{nmask} \times d_l}$, where d_l is the latent space dimension.

The MAE encoder is implemented based on a ViT model with a depth of 6, featuring attention mechanisms in each layer. This architecture utilizes 8 heads for multi-head attention and incorporates feed-forward networks with a dimension of 512. The embedding dimension is set to 256, encoding each 3D hand joint coordinate into a 256-element vector ($d_l = 256$). **Decoder.** The MAE decoder is designed to complement 351 the encoder. It receives a complete set of tokens, which 352 includes encoded visible patches and mask tokens (see to 353 Figure 2). Mask tokens are shared, learned vectors that de-354 note the presence of a missing patch that needs to be pre-355 dicted. To ensure that mask tokens have location informa-356 tion, spatial and temporal positional embeddings are added 357 to all tokens in this set. Subsequently, the decoder attends 358 to this combined sets of tokens using attention mechanism 359 and predicts the coordinates of the missing joints. 360

The MAE decoder is utilized exclusively during pre-
training for the purpose of skeleton reconstruction, with
only the encoder being employed to generate hand poses
representations for downstream tasks.361
362
363

Loss. During pre-training, the reconstruction loss comprises the Mean Squared Error (MSE) loss alongside the anatomical loss, represented by the ROM constraints for each joint and finger calculated given the training data. The total loss, denoted as \mathcal{L} , is expressed as $\mathcal{L} = \mathcal{L}_{MSE} + \lambda \mathcal{L}_A$. 369 Here, $\mathcal{L}_{MSE} = \mathbb{E} \left[||X - \hat{X}||^2 \right]$, and \mathcal{L}_A signifies the anatomical loss, and λ denotes a weighting factor. 371

Minimizing this loss enables the MAE model to refine its predictions, striving to closely match the ground-truth coordinates while respecting anatomical correctness of hand skeleton. This iterative process facilitates the learning of discriminative representations across different hand poses within the latent space.

3.4. Fine-Tuning for Dynamic Recognition

To evaluate the effectiveness of our AG-MAE model 379 in learning discriminative hand pose representations, we 380 employ the spatio-temporal graph convolutional network 381 (STGCN) [45] as the backbone architecture for classifying 382 skeleton sequences. The STGCN excels at capturing tempo-383 ral relationships, allowing it to extract complex patterns in 384 sequential data. Additionally, it uses an edge-attention adja-385 cency matrix constructed with a learnable mask, enhancing 386 its ability to capture spatial dependencies. 387

Online Recognition.We implement a sliding-window-
based model to identify the boundaries of gestures (start
and end) and the gesture performed within the window of
frames. The start and end are defined as the transitions be-
tween gesture classes and the non-gesture class.388
390388
390
391390
391

The online model is based on an STGCN architecture 393 and incorporates a classification head to predict the gesture (including the non-gesture class) within the window. Following [12], two regression heads are integrated: one for 396 identifying the start and another for identifying the end of 397 any detected gesture. 398

431

433

434

435

436

437

438

439

440

399 A cross-entropy loss is applied to the gesture class out-400 put, and an MSE loss is applied to the two regression out-401 puts (start and end of the gesture, if any).

402 Offline Recognition. The offline model, on the other 403 hand, is trained on segmented, isolated gestures. We use an STGCN model with a single classification head to output a 404 405 gesture label for each segmented sequence. All sequences 406 are padded to the dataset maximum length for uniform pro-407 cessing, and training is conducted using cross-entropy loss.

For both online and offline settings, given a 3D hand 408 409 joint sequence, we utilize the pre-trained MAE encoder to extract the corresponding learned representations (la-410 411 tent space), which serve as the foundation for training the 412 STGCN models. No masking is applied during finetuning.

4. Experimental Setup 413

4.1. Evaluation Protocols and Metrics 414

Unlike offline evaluation, which focuses mainly on recog-415 416 nition accuracy, online evaluation relies on more in-depth 417 metrics to assess performance in real time, including:

Detection Rate (DR) measures the ratio of correctly de-418 419 tected gestures to the total number of gestures, considering temporal overlap with ground truth and duration consis-420 421 tency. A gesture is correctly detected if its temporal overlap exceeds 50% of the true interval, does not exceed twice the 422 423 actual duration, and matches the label.

Levenshtein Accuracy (LA) captures recognition accu-424 425 racy regardless of early or late detection. It's also known as 426 minimum edit distance, meaning it measures the minimum 427 number of single-label insertions, deletions, and substitutions needed to transform a set of labels into another. 428

Jaccard Index (JI) refers to the average relative overlap between ground truth and predicted labels, providing 430 insights into the alignment of detected gestures with the 432 ground truth gestures.

False Positive rate (FP) quantifies the ratio of false positive predictions to the total number of gestures, highlighting the ability of the model to minimize erroneous detections. Minimal false positives are desirable for robust gesture recognition systems.

Inference Time (IT) denotes the duration required for the model to perform inference and label a single frame, crucial for assessing real-time applicability.

441 Normalized Time to Detect (TNtD) quantifies the fraction of the sequence duration, from start to end, before 442 the system successfully detects the gesture. Normalization 443 aids in comparing detection performance across different 444 445 sequence lengths.

4.2. Datasets

The key characteristics of the datasets used for online eval-447 uation are given in Table 1. 448

Dataset	#S	#G	#J	#G/S	MeanT	StdT
SHREC21 [6]	180	17	20	3-5	77	61
IPN Hand [4]	4000	14	21	21	140	94
ODHG [14]	280	14	22	10	58	27

Table 1. Statistics for evaluation datasets: S (sequences), G (gestures), J (joints), G/S (continuous gestures per sequence), MeanT (average gesture duration), StdT (standard deviation).

SHREC'21: The SHREC'2021 Track dataset [6] meets 449 to practical application scenarios requiring real-time gesture 450 recognition within continuous hand movement sequences. 451 It includes 18 gesture classes categorized as static, coarse 452 dynamic, and fine dynamic gestures. Evaluation metrics in-453 clude DR, FP rate, JI and IT. 454

IPN Hand: The IPN Hand dataset [4] comprises over 4,000 gesture instances from 50 subjects. Each subject executed 21 gestures continuously, interspersed with random pauses, in a single video We use the provided training/test split for evaluation. Evaluation is based on LA and IT.

ODHG: Online Dynamic Hand Gesture (ODHG) [43] is the online version of the SHREC'17 track [14], providing 280 sequences of 10 non-segmented gestures occurring sequentially and performed by 28 subjects in a continuous online environment. Evaluation is based on LA and TNtD.

Due to the variability in hand models across datasets, particularly regarding the number of joints, we propose a dataset-specific approach for inferring anatomical constraints. Specifically, we derive the ranges for these constraints, namely the minimum and maximum values for bone lengths and joint angles, based on the training set.

4.3. Implementation Details

For the STMAE model, key hyperparameters include learn-472 ing rates. We employed the AdamW optimizer with a learn-473 ing rate of 2×10^{-4} and weight decay of 5×10^{-2} . The 474 learning rate is gradually reduced during training, with the 475 biomechanical loss weighting factor (λ) set to 1.0. The win-476 dow size is set to W = 16. Similarly, for both STGCN 477 models, we utilized the AdamW optimizer with a learning 478 rate of 1×10^{-3} and weight decay of 5×10^{-2} . The learning 479 rate undergoes gradual reduction throughout training. We 480 employed cross-entropy for training loss with label smooth-481 ing during fine-tuning, with a smoothing rate of 0.1. We use 482 a sliding window W = 16 as we found it to be an optimal 483 choice. All experiments are conducted using an NVIDIA 484 GeForce RTX 2080 GPU. 485

446

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

3DV 2025 Submission #60. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.



530

531

532

533

534

535



Figure 4. Left: example of an anatomically correct generated hand pose. Middle: example of a non-anatomically correct generated hand pose (thumb tip is extended beyond the normal range). Right: minimum and maximum bone lengths in the IPN Hand dataset.

486 5. Experimental Results

487 5.1. Ablation Studies

We conduct ablation studies to assess the effectiveness of
various components and enhancements in our method. All
experiments are conducted using the SHREC'21 dataset.

491 Masking Strategy. Our analysis of different masking 492 strategies provides valuable insights (Table 2). Random spatio-temporal masking with a ratio of 0.6 ($m_r = 60\%$) 493 proves to be the most effective, achieving 91.9% DR and 494 a minimum FP rate of 0.033, highlighting its effectiveness 495 496 for classification tasks. In contrast, random spatial mask-497 ing achieves the highest JI with a ratio of 0.7, highlighting its strength in detection tasks. However, the random 498 499 temporal masking shows comparatively lower performance, which can be attributed to its limited effectiveness in online 500 gesture recognition. This reduced performance is likely due 501 502 to the disruption of critical sequential patterns and temporal context that are essential for accurate real-time gesture 503 recognition. In particular, the random nature of temporal 504 masking can lead to masking of contiguous frames with-505 506 out intermediate information, disrupting the temporal flow. We argue that a guided temporal masking approach, such 507 508 as one informed by joint motion, may be more effective as 509 it reduces randomness and ensures that masking does not obscure important sequential information. 510

Masking strategy	Ratio	$\mathbf{DR}\uparrow$	$\mathbf{FP}\downarrow$	JI ↑
	0.5	90.3%	0.0490	0.6346
Random spatial	0.6	87.1%	0.0414	0.6228
I I	0.7	90.5%	0.0626	0.7257
	0.5	74.5%	0.5331	0.5024
Random temporal	0.6	80.8%	0.0516	0.5707
	0.7	81.6%	0.0753	0.5380
	0.5	88.3%	0.0694	0.5568
Random spatio-temporal	0.6	91.9 %	0.0330	0.6800
	0.7	88.3%	0.0406	0.6412

Table 2. Ablation study on masking strategy and ratio in pretraining phase (SHREC'21).

Feature Embedding. We compare FFE against learned 511 linear mapping using a fully connected layer (Table 3 -512 Line 1). Our experiments show that FFE significantly out-513 performs linear mapping. This enhancement is due to the 514 ability of FFE to capture intricate spatial and temporal rela-515 tionships among joints in skeletal data. By projecting input 516 coordinates into a high-frequency domain, FFE allows the 517 network to encode finer details, thereby improving the qual-518 ity of learned representations. 519

Anatomical Loss. The inclusion of anatomical con-520 straints significantly improves model performance, as evi-521 denced by improvements in all evaluation metrics (Table 3 522 - Line 2). The anatomical loss provides critical anatomical feedback during the pre-training phase that improves model 524 robustness without impacting inference time. This loss term 525 helps generate anatomically correct hand poses, reducing 526 misinterpretation of joint positions that could lead to confu-527 sion between gestures, especially in non-gesture frames that 528 involve random hand movements. 529

Figure 4 illustrates the differences between correct and incorrect hand poses; the latter is shown with an exaggerated thumb extension, which is effectively penalized by the anatomical loss. In addition, the anatomical constraints adapt to different hand shapes and sizes by defining bounding ranges for bone lengths and joint angles.

Method	DR↑	$\mathbf{FP}\downarrow$	JI ↑	IT (ms) \downarrow
AG-MAE <i>w/o</i> FE AG-MAE <i>w/o L</i> _A	83.1% 84.4%	0.082 0.065	0.573 0.571	0.41 0.63
AG-MAE	91.9%	0.033	0.680	0.63

Table 3. Ablation studies on different components of AG-MAE model (SHREC'21).

5.2. Comparison with State-of-the-Art Methods 536

Offline Evaluation.We first assess our approach in an of-
fline setting, focusing on segmented hand gesture sequences537from three distinct datasets:SHREC'21, IPN Hand, and539ODHG. The results, as detailed in Table 4, particularly emphasizing the critical role of self-supervised learning and541

3DV 2025 Submission #60. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

Method	Accuracy	Method	Accuracy	Method	Accuracy
DDNet [46]	87.8%	ResNeXt-101 [4]	86.3%	G Spotter[36]	95.3%
Stronger [16]	97.5%	Dist-Time [18]	87.5%	DSTA-Net [37]	97.0%
AG-MAE	98.5%	AG-MAE	93.7%	AG-MAE	93.6%
SHREC'21	Dataset.	IPN Hand Da	ataset.	ODHG Da	taset.

Table 4. Offline results on evaluation datasets.

pretraining in learning spatio-temporal representations of
hand skeleton data. Notably, our model achieves SOTA performance on the SHREC'21 and IPN Hand datasets.

Online Evaluation. For online evaluation, we adhere 545 to the proposed evaluation protocol and metrics for each 546 dataset. Table 5 shows the comparative performance of dif-547 548 ferent methods on the SHREC'21 dataset. Our approach achieves SOTA results in terms of DR and FP rate. No-549 tably, while group 4 of the original SHREC'21 paper also 550 551 uses the STGCN backbone, our model, augmented with a 552 masked autoencoder (MAE) for better representation learn-553 ing, achieves an improved recognition rate of 91.9% with a 554 notable reduction in false positives to 0.033.

However, we observe a notable decrease in the JI com-555 pared to the STGCN-based method from Group 4 (G4) 556 557 of the SHREC'21 paper. This difference may stem from Group 4's use of two separate models-one for detection 558 and one for classification-as well as their incorporation of 559 handcrafted similarity evaluations. Such handcrafted fea-560 tures can be highly effective for specific gestures, contribut-561 562 ing to their higher JI scores [6].

Method	Backbone	$\mathbf{DR}\uparrow$	$\mathbf{FP}\downarrow$	$\mathbf{JI}\uparrow$	$IT(ms)\downarrow$
G1 [Shrec21] [6]	Transformer	79.2%	0.257	0.603	1.36
G2 [Shrec21] [6]	CNN	48.6%	0.927	0.277	0.41
G3 [Shrec21] [6]	GRU	75.7%	0.340	0.619	3e-6
G4 [Shrec21] [6]	STGCN	89.9%	0.066	0.853	0.16
Stronger [16]	CNN	90.6%	0.347	0.740	0.10
G Spotter [36]	LSTM	90.3%	0.053	0.852	-
AG-MAE	STGCN	91.9%	0.033	0.680	0.63

Table 5. Online recognition results on SHREC'21 dataset.

Method	Modality	$\mathbf{LA}\uparrow$	IT (ms) \downarrow
ResNet50 [4]	RGB-Seg	33.27%	29.2
ResNet50 [4]	RGB-Flow	39.47%	43.1
ResNeXt-101 [4]	RGB-Seg	39.01%	39.9
ResNeXt-101 [4]	RGB-Flow	42.47%	53.7
TMMF [19]	RGB-Flow	68.12%	-
TSN-TSM [5]	RGB-Seg	65.27%	15.2
AG-MAE	3D keypoints	73.93%	19.4*

Table 6. Online evaluation on the IPN Hand dataset. (*) indicates that IT includes both keypoint extraction and inference times.

Table 6 demonstrates SOTA results of our model in terms of LA. Despite the relatively high reported inference time,

it is important to note that this includes the additional time 565 required for the extraction of 21 3D keypoints using Medi-566 aPipe [47]. Specifically, the 3D keypoints extraction con-567 tributes approximately 19.33 ms to the overall inference 568 time, while the inference time of our model alone is on the 569 order of $10^{-1}ms$. This suggests that the observed inference 570 time is primarily influenced by the keypoints extraction pro-571 cess rather than the model itself. 572

For the ODHG dataset, due to the lack of a standard evaluation split protocol, we follow the authors approach and573uation split protocol, we follow the authors approach and574employ a random k-fold split, allocating 70% of the data575for training and 30% for testing. Our model achieves an LA576of 82.0% and an NTtD of 0.34. The original paper reported577comparable results, with an LA of 82.2% and an NTtD of5780.21 using depth images.579

Limitations. Despite the notable performance of our 580 method, some limitations should be acknowledged. A key 581 limitation is the trade-off between DR, FP rate, and JI de-582 pending on the masking strategy and ratio employed (see 583 Table 2). This trade-off indicates that the optimal masking 584 strategy and ratio may depend on the specific application re-585 quirements, highlighting the need for a balanced approach 586 to achieve the best overall performance. Additionally, inte-587 grating the MAE with the STGCN backbone increases com-588 putational complexity, resulting in longer inference times. 589 This may constrain the practical deployment of the model 590 in real-time scenarios where processing speed is crucial. 591

6. Conclusion and Future Work

In this work, we introduce a novel framework for online hand gesture recognition combining self-supervised learning with anatomical constraints. By pre-training a spatiotemporal masked autoencoder with anatomical guidance, our approach transforms 3D hand keypoints into highly discriminative representations, enhancing performance in online gestire recognition. Comprehensive evaluations on SHREC'21, IPN Hand, and ODHG datasets shows that our method achieves SOTA results.

Future research will focus on refining adaptive masking strategies to further improve overall performance in online scenarios. Additionally, we will work on reducing model complexity to develop faster, more efficient models for deployment in resource-constrained environments.

563 564

610

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715

716

717

718

719

720

721

722

608 **References**

- [1] Andreas Aristidou. Hand tracking with physiological constraints. *The Visual Computer*, 34:213–228, 2018. 2
- [2] Danilo Avola, Marco Bernardi, Luigi Cinque, Gian Luca
 Foresti, and Cristiano Massaroni. Exploiting recurrent neural networks and leap motion controller for the recognition of
 sign language and semaphoric hand gestures. *IEEE Transac- tions on Multimedia*, 21(1):234–245, 2018. 2
 - [3] Hangbo Bao, Li Dong, Songhao Piao, and Furu Wei. Beit: Bert pre-training of image transformers. *arXiv preprint arXiv:2106.08254*, 2021. 2
 - [4] Gibran Benitez-Garcia, Jesus Olivares-Mercado, Gabriel Sanchez-Perez, and Keiji Yanai. Ipn hand: A video dataset and benchmark for real-time continuous hand gesture recognition. In 25th International Conference on Pattern Recognition, ICPR 2020, Milan, Italy, Jan 10–15, 2021, pages 4340– 4347. IEEE, 2021. 6, 8
 - [5] Gibran Benitez-Garcia, Lidia Prudente-Tixteco, Luis Carlos Castro-Madrid, Rocio Toscano-Medina, Jesus Olivares-Mercado, Gabriel Sanchez-Perez, and Luis Javier Garcia Villalba. Improving real-time hand gesture recognition with semantic segmentation. *Sensors*, 21(2):356, 2021. 8
 - [6] Ariel Caputo, Andrea Giachetti, Simone Soso, Deborah Pintani, Andrea D'Eusanio, Stefano Pini, Guido Borghi, Alessandro Simoni, Roberto Vezzani, Rita Cucchiara, et al. Shrec 2021: Skeleton-based hand gesture recognition in the wild. *Computers & Graphics*, 99:201–211, 2021. 2, 6, 8
 - [7] Fabio Marco Caputo, S Burato, Gianni Pavan, Théo Voillemin, Hazem Wannous, Jean-Philippe Vandeborre, Mehran Maghoumi, EM Taranta, Alaleh Razmjoo, Joseph J LaViola Jr, et al. Shrec 2019 track: online gesture recognition. In *Eurographics Workshop on 3D Object Retrieval*. The Eurographics Association, 2019. 2
 - [8] Yujin Chen, Zhigang Tu, Di Kang, Linchao Bao, Ying Zhang, Xuefei Zhe, Ruizhi Chen, and Junsong Yuan. Model-based 3d hand reconstruction via self-supervised learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10451–10460, 2021. 1, 2
 - [9] Yabo Chen, Yuchen Liu, Dongsheng Jiang, Xiaopeng Zhang, Wenrui Dai, Hongkai Xiong, and Qi Tian. Sdae: Selfdistillated masked autoencoder. In *European Conference on Computer Vision*, pages 108–124. Springer, 2022. 2
- [10] Fai Chen Chen, Silvia Appendino, Alessandro Battezzato,
 Alain Favetto, Mehdi Mousavi, and Francesco Pescarmona.
 Constraint study for a hand exoskeleton: human hand kinematics and dynamics. *Journal of Robotics*, 2013(1):910961,
 2013. 4
- [11] Salvador Cobos, Manuel Ferre, MA Sanchez Uran, Javier
 Ortego, and Cesar Pena. Efficient human hand kinematics for
 manipulation tasks. In 2008 IEEE/RSJ International Conference on Intelligent Robots and Systems, pages 2246–2251.
 IEEE, 2008. 4
- [12] Federico Cunico, Federico Girella, Andrea Avogaro, Marco
 Emporio, Andrea Giachetti, and Marco Cristani. Oo-dmvmt:
 A deep multi-view multi-task classification framework for
 real-time 3d hand gesture classification and segmentation. In

Proceedings of the IEEE/CVF Conference on Computer Vi-
sion and Pattern Recognition, pages 2745–2754, 2023. 2,
5665
666

- [13] Mark R Cutkosky et al. On grasp choice, grasp models, and the design of hands for manufacturing tasks. *IEEE Transactions on robotics and automation*, 5(3):269–279, 1989. 4
- [14] Q. De Smedt, H. Wannous, J.-P. Vandeborre, J. Guerry, B. Le Saux, and D. Filliat. 3d hand gesture recognition using a depth and skeletal dataset: Shrec'17 track. In *Proceedings of the Workshop on 3D Object Retrieval*, page 33–38, Goslar, DEU, 2017. Eurographics Association. 6
- [15] Naina Dhingra and Andreas Kunz. Res3atn-deep 3d residual attention network for hand gesture recognition in videos. In 2019 international conference on 3D vision (3DV), pages 491–501. IEEE, 2019.
- [16] Marco Emporio, Ariel Caputo, and Andrea Giachetti. Stronger: Simple trajectory-based online gesture recognizer. 2021. 2, 8
- [17] Marco Emporio, Ariel Caputo, Andrea Giachetti, Marco Cristani, Guido Borghi, Andrea D'Eusanio, Minh-Quan Le, Hai-Dang Nguyen, Minh-Triet Tran, Felix Ambellan, et al. Shrec 2022 track on online detection of heterogeneous gestures. *Computers & Graphics*, 107:241–251, 2022. 2
- [18] Graziano Fronteddu, Simone Porcu, Alessandro Floris, and Luigi Atzori. A dynamic hand gesture recognition dataset for human-computer interfaces. *Computer Networks*, 205: 108781, 2022. 8
- [19] Harshala Gammulle, Simon Denman, Sridha Sridharan, and Clinton Fookes. Tmmf: Temporal multi-modal fusion for single-stage continuous gesture recognition. *IEEE Transactions on Image Processing*, 30:7689–7701, 2021. 8
- [20] Mallika Garg, Debashis Ghosh, and Pyari Mohan Pradhan. Gestformer: Multiscale wavelet pooling transformer network for dynamic hand gesture recognition. In *Proceedings* of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 2473–2483, 2024. 1
- [21] 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, pages 762–770, 2022. 2
- [22] 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. 2, 4
- [23] Hezhen Hu, Weichao Zhao, Wengang Zhou, Yuechen Wang, and Houqiang Li. Signbert: pre-training of hand-modelaware representation for sign language recognition. In Proceedings of the IEEE/CVF international conference on computer vision, pages 11087–11096, 2021. 1, 2
- [24] Omar Ikne, Benjamin Allaert, and Hazem Wannous. Skeleton-based self-supervised feature extraction for improved dynamic hand gesture recognition. In 2024 IEEE 18th International Conference on Automatic Face and Gesture Recognition (FG). IEEE, 2024. 1, 2, 3
- [25] Omar Ikne, Rim Slama, Hichem Saoudi, and Hazem Wannous. Spatio-temporal sparse graph convolution network

786

787

788

789

790

791

792

793

794

795

796

797

798

799

800

801

802

803

804

805

806

807

808

809

810

811

812

813

814

815

816

817

818

819

820

for hand gesture recognition. In *The 18th IEEE Interna- tional Conference on Automatic Face and Gesture Recog- nition*, 2024. 1

- [26] Joseph HR Isaac, Muniyandi Manivannan, and Balaraman
 Ravindran. Single shot corrective cnn for anatomically correct 3d hand pose estimation. *Frontiers in Artificial Intelligence*, 5:759255, 2022. 2
- [27] Junkun Jiang, Jie Chen, and Yike Guo. A dual-masked auto-encoder for robust motion capture with spatial-temporal skeletal token completion. In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 5123–5131, 2022. 2
- [28] Okan Köpüklü, Ahmet Gunduz, Neslihan Kose, and Gerhard
 Rigoll. Real-time hand gesture detection and classification
 using convolutional neural networks. In 2019 14th IEEE international conference on automatic face & gesture recognition (FG 2019), pages 1–8. IEEE, 2019. 2
- [29] 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. 2
- [30] John Lin, Ying Wu, and Thomas S Huang. Modeling the constraints of human hand motion. In *Proceedings workshop on human motion*, pages 121–126. IEEE, 2000. 4
- [31] 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. 2
- [32] Mehran Maghoumi and Joseph J LaViola. Deepgru: Deep gesture recognition utility. In Advances in Visual Computing: 14th International Symposium on Visual Computing, ISVC 2019, Lake Tahoe, NV, USA, October 7–9, 2019, Proceedings, Part I 14, pages 16–31. Springer, 2019. 2
- [33] Eshed Ohn-Bar and Mohan Manubhai Trivedi. Hand gesture recognition in real time for automotive interfaces: A multimodal vision-based approach and evaluations. *IEEE transactions on intelligent transportation systems*, 15(6):2368–2377, 2014.
- [34] Haocong Rao, Shihao Xu, Xiping Hu, Jun Cheng, and Bin
 Hu. Augmented skeleton based contrastive action learning
 with momentum lstm for unsupervised action recognition. *Information Sciences*, 569:90–109, 2021. 2
- [35] Chr Ryf and A Weymann. The neutral zero method—a principle of measuring joint function. *Injury*, 26:1–11, 1995. 4
- [36] Junxiao Shen, John Dudley, George Mo, and Per Ola Kristensson. Gesture spotter: A rapid prototyping tool for key gesture spotting in virtual and augmented reality applications. *IEEE transactions on visualization and computer graphics*, 28(11):3618–3628, 2022.
- [37] Lei Shi, Yifan Zhang, Jian Cheng, and Hanqing Lu. Decoupled spatial-temporal attention network for skeleton-based action-gesture recognition. In *Proceedings of the Asian conference on computer vision*, 2020.
- [38] Adrian Spurr, Umar Iqbal, Pavlo Molchanov, Otmar Hilliges,
 and Jan Kautz. Weakly supervised 3d hand pose estimation
 via biomechanical constraints. In *European conference on computer vision*, pages 211–228. Springer, 2020. 2, 4

- [39] Kun Su, Xiulong Liu, and Eli Shlizerman. Predict & cluster: 781
 Unsupervised skeleton based action recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9631–9640, 2020. 2
- [40] Matthew Tancik, Pratul Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ramamoorthi, Jonathan Barron, and Ren Ng. Fourier features let networks learn high frequency functions in low dimensional domains. Advances in Neural Information Processing Systems, 33:7537–7547, 2020. 2, 3
- [41] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. 5
- [42] Ben Veldhuijzen, Remco C Veltkamp, Omar Ikne, Benjamin Allaert, Hazem Wannous, Marco Emporio, Andrea Giachetti, Joseph J LaViola Jr, Ruiwen He, Halim Benhabiles, et al. Shrec 2024: Recognition of dynamic hand motions molding clay. *Computers & Graphics*, page 104012, 2024.
- [43] Hazem Wannous and Jean-Philippe Vandeborre. Continuous hand gesture recognition using deep coarse and fine hand features. In *The 33rd British Machine Vision Conference– BMVC 2022*, 2022. 6
- [44] Hong Yan, Yang Liu, Yushen Wei, Zhen Li, Guanbin Li, and Liang Lin. Skeletonmae: Graph-based masked autoencoder for skeleton sequence pre-training. *arXiv preprint arXiv:2307.08476*, 2023. 2, 4
- [45] Sijie Yan, Yuanjun Xiong, and Dahua Lin. Spatial temporal graph convolutional networks for skeleton-based action recognition. In *Proceedings of the AAAI conference on artificial intelligence*, 2018. 5
- [46] Fan Yang, Yang Wu, Sakriani Sakti, and Satoshi Nakamura. Make skeleton-based action recognition model smaller, faster and better. In *Proceedings of the ACM multimedia asia*, pages 1–6. 2019. 2, 8
- [47] Fan Zhang, Valentin Bazarevsky, Andrey Vakunov, Andrei Tkachenka, George Sung, Chuo-Ling Chang, and Matthias Grundmann. Mediapipe hands: On-device real-time hand tracking. arXiv preprint arXiv:2006.10214, 2020. 8
- [48] Nenggan Zheng, Jun Wen, Risheng Liu, Liangqu Long, Jianhua Dai, and Zhefeng Gong. Unsupervised representation learning with long-term dynamics for skeleton based action recognition. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2018. 2
 821
 822
 823
 824
 825