

1 **Identifying unsafe behavior of construction workers: a dynamic**
2 **approach combining skeleton information and spatiotemporal features**

3 Han Wu^a, Yu Han^{b,*}, Meng Zhang^c, Bihonegn Dianarose Abebe^d, Molla Betelhem Legesse^e, Ruoyu Jin^f

4

5 ^a Graduate Student, Faculty of Civil Engineering and Mechanics, Jiangsu University, 301 Xuefu Road,
6 Zhenjiang, 212013, Jiangsu, China. Email: wu1024502851@163.com

7 ^{b,*} Professor, Faculty of Civil Engineering and Mechanics, Jiangsu University, 301 Xuefu Road, Zhenjiang,
8 212013, Jiangsu, China. Email: hanyu85@yeah.net

9 ^c Graduate Student,, Faculty of Civil Engineering and Mechanics, Jiangsu University, 301 Xuefu Road,
10 Zhenjiang, 212013, Jiangsu, China. Email: 2221923026@stmail.ujs.edu.cn

11 ^d Graduate Student,, Faculty of Civil Engineering and Mechanics, Jiangsu University, 301 Xuefu Road,
12 Zhenjiang, 212013, Jiangsu, China. Email: Dianaroseabebe27@gmail.com

13 ^e Graduate Student,, Faculty of Civil Engineering and Mechanics, Jiangsu University, 301 Xuefu Road,
14 Zhenjiang, 212013, Jiangsu, China. Email: betelhemlegesse2@gmail.com

15 ^f Associate Professor, School of Built Environment and Architecture, London South Bank University. 103
16 Borough Rd, London SE1 0AA, UK. Email: jinr@lsbu.ac.uk

17 **Abstract**

18 Vision-based methods for action recognition are valuable for the supervision of construction workers'
19 unsafe behaviors. However, existing methods are limited due to the lack of ability to extract worker action
20 information from video streams. Using spatiotemporal relationships between workers' skeletal points to
21 identify hazardous action remains a huge challenge for safety management of construction sites.. In this study,
22 an improved dynamic skeleton model, named Attention Module Spatial-Temporal Graph Convolutional
23 Neural Network (AM-STGCN) is built from the modality data of 2D skeleton points, and a combination of
24 designed human partitioning strategies and non-local attention mechanisms are adopted to extract global
25 information during worker movement to automatically identify unsafe behaviors on construction sites. The
26 method includes three basic modules, namely video data acquisition, workers' skeleton information extraction,
27 as well as recognition and classification of hazardous actions. The test accuracy reached 93.66% in the
28 laboratory, and 90.50% and 87.08% in typical working scenarios (i.e., high-altitude working scenarios with
29 close-up and far views) respectively. The promising test results indicated that the developed AM-STGCN
30 model could be more widely applied in wider construction scenarios, such as foundation excavation.

31 **Keywords**

32 Hazard scenario; Unsafe behavior; Construction safety; Skeleton modality data; Action recognition;
33 Dynamic model;

34 **1. Introduction**

35 The construction industry has been a global concern due to its high risk measured accident rates. The death
36 toll due to accidents in construction accounts for more than 20% of occupational deaths every year. According

37 to the statistics of the Health and Safety Executive of the UK, a total of 123 workers died due to accidents in
38 the UK in 2021, and 30 of them were related to the construction industry, accounting for 24.4% of the total
39 death toll of accidents (HSE 2022). In China, the death toll in the construction industry reached 794 in 2020
40 alone (MOHURD 2022). Site accidents are often the result of a combination of factors, the Heinrich Accident
41 Causation Theory stated that unsafe behavior of workers is the core cause of accidents (Heinrich 1941),
42 among them the lack of safety protection equipment and workers' hazardous actions are the main causes. To
43 facilitate safety inspection on construction sites, numerous studies on unsafe behavior of workers have been
44 conducted.

45 In recent years, computer vision technology has become one of the mainstream themes in construction
46 safety research (Fang et al. 2020). Machine vision technology based on color and contour feature extraction
47 (e.g., HOG) and deep learning-based target detection technology such as Faster-R CNN, YOLO, and SSD
48 were among the earlier major approaches on the detection of safety gear of workers. These methods were
49 used to evaluate the interaction between workers and safety helmets, such as space relationship, geometric
50 information and color feature (Park et al. 2015) or to evaluate workers' safety helmets and safety belts under
51 different operating conditions, e.g., scene, weather, light, etc. (Trabelsi et al. 2019; Fang et al. 2018). These
52 studies showed that the research on the inspection of workers' safety protective gear has achieved promising
53 outcomes, but the study of workers' hazardous actions had been still insufficient.

54 The essence of hazardous action is the change of skeleton joints, which is composed of the state of multiple
55 consecutive frames. Focusing on video clips allows a better and more accurate reflection of the features of
56 the action. The aforementioned detection methods for safety protection gear are only suitable for the detection

57 of static target states, and the use of these methods for the recognition of workers' hazardous actions will
58 significantly reduce the accuracy and increase the false alarm rate. At present, sensor technology is commonly
59 used to detect angular ratios and spatially varying signals between skeleton points to assess and classify the
60 behavior of construction workers. This requires the installation of sensors on each worker and piece of
61 equipment and is a heavy burden for construction site applications. Some researchers (Kim et al., 2016; Fang
62 et al., 2019) have applied vision-based methods to study the relationship between construction workers and
63 targets (e.g., machinery, equipment, materials, etc.) at a given moment, such as assessing the risk conditions
64 of workers who have just stepped into a danger zone or are in the blind spot of construction machinery.
65 However, these studies had not addressed workers' movements during non-conforming operations.
66 Combining long and short-term memory (LSTM) networks and other neural network approaches for time
67 series prediction of unsafe behaviors (Kong et al. 2021; Tang et al. 2020) has made good progress in
68 behavioral regulation. However, action analysis through temporal variation alone ignores spatial variation in
69 human posture, and this a challenge to tackle for recognition of complex construction actions.

70 The analysis method based on skeleton modality data for action recognition has achieved positive results
71 in extracting and application of motion information. It could extract information about workers' skeleton
72 points, analyze the spatial relationships between workers' skeleton points, and evaluate workers' behaviors.
73 Currently, many studies focus on combining the skeleton with sensors or vision devices, dividing the acquired
74 video into multiple static skeleton images to obtain the change information about workers' motion angle,
75 acceleration and skeleton length. Isolating the dynamic process of motion into static images to process the
76 information effectively extracts the spatial relationship of the skeleton at a single frame, but this approach

77 ignores the temporal relationship of the skeleton information in different frames. The lack of utilization of
78 spatiotemporal information in the dynamic process of workers' construction movements ultimately affects
79 the accuracy of motion recognition. In general, there are two limitations in the current recognition methods
80 for construction workers' actions, specifically: (1) The spatiotemporal information between the skeletons in
81 the construction activities is difficult to use, and a large amount of effective information would be lost; and
82 (2) the lack of effective detection of dynamic action processes affects the recognition accuracy of hazardous
83 actions.

84 Aiming to address these limitations, the main purpose of this study is to propose a deep learning method
85 that combines skeleton information and spatiotemporal features for the recognition of construction workers'
86 hazardous actions. Researchers established an improved dynamic skeleton model, which takes into account
87 the temporal and spatial relationship of adjacent skeleton joints. The model can be used for the analysis of
88 the dynamic process of construction workers' hazardous construction actions, and for achieving automatic
89 recognition and detection of hazardous actions. Unlike CNN, which segments all videos into frame-by-frame
90 pictures and inputs them into the network, this method directly inputs the skeleton and joint data of workers,
91 hence greatly reducing the number of parameters. It can be used in video surveillance systems to effectively
92 prevent on-site safety accidents.

93 **2. Literature review**

94 Traditionally, the recognition methods for construction workers' hazardous actions can be divided into two
95 categories in terms of implementation, namely sensor-based and vision-based methods.

96 **2.1 Sensor-based action recognition**

97 Sensor-based action recognition methods typically focus on workers' gestures and fall postures to observe
98 changes in limbs during worker's action. Cheng et al. (2013) obtained position parameters and chest posture
99 data from sensors mounted on workers over a sequence of time to identify workers' actions during activity.
100 Fang and Dzeng (2014) mounted workers' vests and helmets with motion sensors and brainwave sensors to
101 detect workers' fall risk by correlating changes in these two externally transmitted signals over time. Jebelli
102 et al. (2016) used an inertial measurement unit to record the characteristics of changes in sensor parameters
103 over this time period during a worker's fall to comprehensively assess the fall risk of rebar workers. Akhavian
104 and Behzadan (2016) captured workers' body movements by using embedded accelerometers and gyroscopic
105 sensors and simulated various types of construction activities in the laboratory through time-series changes
106 in parameters. As the use of sensors for construction action recognition requires manual processing of large
107 amounts of data, inflexible methods, complex operations and an unprotected user experience, researchers
108 have been gradually combining machine learning and deep learning methods with sensors for construction
109 action recognition.

110 For example, Gong et al. (2022) adopted a machine learning approach to analyze data from wearable
111 sensors over multiple time periods to identify and classify construction behavior based on parametric
112 temporal features. Bangaru et al. (2021) combined wearable EMG and MU sensors with ANN artificial neural
113 networks to perform data mining in the form of time series on sensor parameters installed on multiple parts
114 of the workers' body to achieve automatic recognition of their' construction actions, and the test results
115 showed good robustness. Ogunseju et al. (2021) used an Inception v1 network to acquire time-series data
116 signals from wearable sensors on workers' lower arms to identify and classify worker actions such as

117 carpentry. However, even though machine learning and deep learning methods improved the speed and
118 accuracy of detection, the results were still with a large amount of data and graphics. Manual analysis of the
119 data was required to define a range of parameter values for the action features. This recognition process
120 ignored the spatial characteristics and feature associations of worker actions. In addition, it demanded a high
121 level of knowledge from managers and was not practical for on-site safety supervision.

122 **2.2 Vision-based action recognition**

123 Vision-based action recognition method is a popular research trend in construction safety in recent years.
124 It mainly analyzes workers' construction actions by collecting construction images and videos. Using a single
125 frame picture to recognize action cannot effectively obtain coherent time information in the process of
126 hazardous actions, hence often leading to misjudgment (Guo and Lai 2014). Using RGB video as the research
127 object could obtain the spatial and temporal information of workers' limbs, and that could significantly
128 improve the recognition accuracy (Zhang 2019; Zhao 2019).

129 **2.2.1 Human action recognition**

130 For vision-based methods, action recognition often requires the extraction of action features. Manual
131 feature extraction methods are the main way to extract action features. Feature descriptors such as HoG, HOF
132 and MBH are introduced into the iDT algorithm to obtain the trajectory of feature points and to describe
133 human behavior. But generating local descriptors to describe human behaviors by manually extracting
134 spatiotemporal interest points will take longer computation time and lose more valuable information in videos
135 (Laptev 2005).

136 Deep learning methods have emerged due to its excellent extraction features and inspection efficiency in

137 image and video processing. There are three main types of methods: two-stream CNN methods, 3D-CNN
138 (3D convolutional neural network) methods and skeleton-based methods. Simonyan and Zisserman (2014)
139 divided the neural network into two parts, one for capturing the spatial features of images and the other for
140 analyzing the temporal information contained in videos. Since then, many scholars have improved this
141 method (Lan et al. 2017; Zhou et al. 2018). For example, Wu et al. (2015) proposed a method based on LSTM
142 and CNN, they applied CNN to perform feature extraction on video clips, then used LSTM to classify long-
143 term span temporal features, and finally extracted multiple manually defined different action features.
144 However, the result was the textual output form of the corresponding action. In addition, the time domain
145 information in the dual-stream network was all derived from the inter-frame optical flow, which was not good
146 for grasping information for a long time, and could be easily affected by many factors, e.g., background, light,
147 and shadow, etc. Compared to 2DCNN, 3DCNN has one more dimension for capturing temporal information,
148 so that long-term information in the action process can be effectively utilized. Tran D et al. (2015) proposed
149 a C3D architecture based on 3D CNN, which could capture spatiotemporal information for human action
150 recognition and improve the recognition accuracy greatly. However, due to a large number of parameters, it
151 was a heavy burden for the actual application effect.

152 Since the skeleton modality data is not affected by the above-mentioned factors and the connection effect
153 between skeleton points can visually represent the action information, it is more suitable for action
154 recognition. However, RNN and CNN networks treat the skeleton point data input by RGB video as a long-
155 term sequence or 2D matrix to extract features, hence making it difficult to understand the connection
156 information between human skeleton joints (Li et al. 2018; Si et al. 2019), resulting in poor recognition of

157 actions. As one type of graph neural network, GCN analyzes data by using generalized topological graph
158 structure, and is good at processing the relationship between such non-Euclidean data and modeling nodes,
159 which is suitable for the extraction of human skeleton information (Yan et al. 2018).

160 **2.2.2 Hazardous action recognition method in construction**

161 For the study of construction workers' hazard actions, previous studies usually focused on the detection,
162 location and tracking of workers. Memarzadeh et al. (2013) detected construction workers and equipment by
163 analyzing construction activities in videos using directional gradients and color histograms. Kim et al. (2016)
164 combined computer vision with fuzzy inference method and augmented reality technology to monitor
165 workers' contact with dangerous areas and evaluate workers' behavioral safety conditions when they worked
166 nearby heavy equipment.

167 In recent years, some scholars have focused their research on the process of workers' actions. Yang et al.
168 (2016) extracted the location of workers under continuous time series by dense trajectory method to identify
169 workers' behavior with positive results under MBH descriptors. Ding et al. (2018) integrated CNN and LSTM
170 networks to achieve automatic recognition of workers' unsafe behaviors by extracting visual features in a
171 video stream. It actually recognized actions by obtaining different time series information of multiple key
172 points in the video through LSTM and being unified by CNN after extracting data features. However, this
173 approach ignored the spatial feature changes in images of different frames attention, and therefore the result
174 generated was a textual description matching the action category.

175 To capture the spatial features of workers' movement more clearly, Escorcía et al. (2012) adopted an RGB-
176 D camera to collect motion skeleton data of workers under construction in a building and then used a

177 discriminative classifier to detect workers' actions. Similarly, Han and Lee (2013) extracted 3D skeleton data
178 of workers from videos and effectively used spatial features in the images to identify hazardous actions.
179 However, the spatial features of consecutive frames in the video are often redundant, and the use of 3D
180 convolution introduces repetitive spatial features. Such a form of action recognition actually discriminates
181 from a particular 3D skeleton pose by ignoring the changing relationship between the skeletons during the
182 action, and has low accuracy in recognizing actions with similar postures. Yu et al. (2017) applied a static
183 recognition method based on the image skeleton and tested the accuracy of worker's climbing action by
184 changing the values of the joint parameters, avoiding the redundancy of the parameters. Manually recording
185 the parameter changes in joint angle values was cumbersome and was limited to static skeleton information.
186 The detection accuracy for the three hazardous actions was only 81.44%, which would be further reduced
187 when used for the identification of multiple unsafe action types at construction sites. Guo et al. (2018) also
188 simplified the dynamic skeleton movement process to a static process in order to achieve real-time detection
189 of unsafe actions. They described the static pose by a few parameters for action recognition, without
190 considering the relationship between the skeleton information in time and space, and that resulted in a high
191 false alarm rate for the action measurement.

192 Based on skeleton modality data, this research is devoted to dealing with the spatiotemporal connection
193 between the overall skeleton data of workers, in order to achieve the dynamic detection of workers' hazard
194 actions and to facilitate the application in actual construction scenarios. A spatiotemporal graph convolutional
195 neural network (ST-GCN) is proposed by adopting a deep learning method for automatically extracting
196 worker skeleton information and identifying the dynamic process of construction workers' hazardous actions.

197 **3. Methodology**

198 The methodology to implement a deep learning model for dynamic recognition of workers' hazardous
199 actions consists of four steps, namely: (1) extraction of worker's skeleton points; (2) selection of action
200 recognition algorithms; (3) data collection and model building; and (4) recognition of dynamic process of
201 actions. The research is intended to provide a general methodological basis for subsequent studies of workers'
202 hazardous construction actions. Fig. 1 shows the method flow designed in this study, with each step including
203 specific implementation details.

204 **3.1 Openpose-based skeleton extraction method**

205 Due to the complex environmental factors of construction sites, workers' bodies are often obscured by
206 construction materials, equipment or structures. One of the current mainstream methods for extracting the
207 skeleton requires a top-down approach, which first detects the human beings by target detection algorithm,
208 and then detects the key points of single human skeleton. However, this kind of method is difficult to perform
209 worker's skeleton identification and extraction when the construction worker's body is more than 30%
210 occluded. Unlike most top-down methods for extracting workers' skeleton joints, the Openpose network uses
211 a bottom-up structure to extract workers' skeleton joints, which first detects each skeleton joint of workers
212 and then connects all the identified skeleton points to generate a complete image of the worker's skeleton.
213 This way of extracting workers' skeleton can effectively reduce the reliance on personnel detectors, improve
214 the timeliness of skeleton point extraction, and enable the recognition of key points of workers' skeleton in
215 the case of multi-person construction. Hence this method is suitable for applications in construction sites with
216 a large number of personnel movements. Therefore, the Openpose network was designed for skeleton

217 extraction of construction workers in this study.

218 VGG-16 was used as a pre-base network in this study, which was able to perform feature extraction and
219 generate feature maps for the prediction and connection of skeleton points via two channels. The specific
220 implementation process of extraction of workers' skeleton point using Openpose pose estimation network is
221 as follows: first, the CPM operation method is adopted to predict the skeleton joint points of all workers in
222 the video collected at the construction sites, and then to detect the heatmap of the skeleton points of the
223 workers (Wei et al. 2016). Each joint generates a corresponding Gaussian peak, and the location of the peak
224 is the worker's skeleton joint. After completing the worker skeleton point prediction, the isolated skeleton
225 points of the workers are connected by regressing the PAFs.

226 **3.2 ST-GCN-based construction action recognition method**

227 Construction action is composed of multiple pose graph structures such as the action of a worker climbing
228 on scaffolding, and it is difficult to effectively identify specific action categories only by considering the
229 spatial location between skeleton points (Zhou et al. 2020). The spatiotemporal graph convolutional neural
230 (ST-GCN) network changed the form of spectral-based convolution of GCN networks. The network adds a
231 temporal convolution module, it combines the location information of skeleton joints and temporal
232 information and introduces the graph convolution in the spatiotemporal domain for capturing the variation
233 patterns among nodes (Yan et al. 2018). In this way, it is easy to identify the types of workers' actions with
234 large changes in spatial location of the skeleton. Therefore, the ST-GCN network was adopted as the action
235 recognition used for this study. Fig. 2 describes the main structure of the ST-GCN introduced.

236 Different from CNN network that performs the sampling method and assign weights to the convolution

237 principle, researchers in this study replaced nodes with image pixel points in ST-GCN networks. Then the
 238 sampling function $p(v_{ti}, v_{tj})$ was used to represent the distance between the first-order neighboring nodes
 239 involved in the convolution process, where v_{ti} is a point in a sequence of joint points, v_{tj} denotes an adjacent
 240 node, and the weight function $w(v_{ti}, v_{tj})$ is applied to represent the weight vector of the nodes and their
 241 neighbors. After the weighted average of the standard normalized $Z_{ti}(v_{tj})$, the updated graph convolution
 242 equations were expressed as Eq. (1) and Eq. (2).

$$243 \quad f_{out}(v_{ti}) = \sum_{v_{tj} \in B(v_{ti})} \frac{1}{Z_{ti}(v_{tj})} f_{in}(p(v_{ti}, v_{tj})) \cdot w(v_{ti}, v_{tj}) \quad (1)$$

$$244 \quad Z_{ti}(v_{tj}) = |\{v_{tk} | l_{ti}(v_{tk}) = l_{ti}(v_{tj})\}| \quad (2)$$

245 The sampling function and weight function mentioned above were designed for the spatial graph structure
 246 only, without considering the temporal factor. Therefore, researchers defined the spatiotemporal graph by
 247 recomputing the label grouping mapping function. The equation of the spatiotemporal graph structure is
 248 shown in Eq. (3).

$$249 \quad l_{ST}(v_{qj}) = l_{ti}(v_{tj}) + (q - t + \lceil \Gamma/2 \rceil \times K) \quad (3)$$

250 Where $l_{ti}(v_{tj})$ is the label map for the single frame case at v_{ti} , Γ is the temporal kernel size, and l_{ST}
 251 represents the labeling map.

252 3.3 Algorithm adjustment and optimization

253 3.3.1 Human partitioning strategy adjustment

254 Considering the complexity of workers' operation actions, the action process often involves not only the
 255 location changes of the limbs, but also the changes of the torso which are equally important. Researchers set
 256 the domain span of the nodes to 2 based on the Spatial partitioning strategy and re-divided the skeleton point

257 neighborhood into 3 subsets, namely root nodes, centripetal groups and centrifugal groups. The new
258 partitioning strategy could expand the extraction range of workers' skeleton features, which could extract
259 features of key points and improve the accuracy rate for construction workers' construction actions. The new
260 partitioning strategy is shown in Fig. 3 where the root node is shown in purple, the skeleton point near the
261 center of gravity of the skeleton and adjacent to the root node (green) is the centripetal group, and the
262 centrifugal group is the neighboring nodes far from the centre of gravity of the skeleton (yellow).

263 Researchers then assigned weights to the skeleton points of each region according to the new partitioning
264 strategy. The new weight assignment is shown in Eq. (4), where r_j represents the distance from the skeleton
265 point j to the centre of gravity of the worker's body, and r_i is the average distance from the center of gravity
266 to the skeleton point.

$$267 \quad l_{ti}(v_{tj}) = \begin{cases} 0 & \text{if } r_j = r_i \\ 1 & \text{if } r_j < r_i \\ 2 & \text{if } r_j > r_i \end{cases} \quad (4)$$

268 In Eq.(4), the centre of gravity is the average coordinate of all joints in a body (black cross in Fig. 3), and
269 0 indicates that no weight is assigned to joints where $r_j=r_i$ (i.e., no change in the joint during movement), 1
270 infers that less weight is assigned to joints where the $r_j < r_i$ (i.e., closer to the centre of gravity), 2 denotes
271 more weight is assigned to joints where the $r_j > r_i$ (i.e., farther from the centre of gravity).

272 3.3.2 Non-local attention mechanism

273 The impact of different body parts on the accuracy of a worker's action recognition during construction
274 varies. Workers rely primarily on the limb parts of the body during construction work, while parts such as
275 the head and neck are not as involved and provide little effective information for movement recognition.

276 Furthermore, the relationships between skeletal joints during construction actions are not restricted to
277 adjacent joints. For example, for many action processes such as probing and climbing, the connection
278 between the joints of the arms, legs and the trunk cannot be ignored. However, in the original ST-GCN
279 network, the perceptual domain of the convolution operation was the neighboring nodes of the root node,
280 which was only used to capture local features of the action process, such as joint changes at the calf and joint
281 changes at the arm. Such a feature extraction approach cannot simultaneously analyze the joint changes at
282 the calf and arm during an action to identify the type of action (Simonyan Zisserman 2015). Despite the
283 adaptation of the human partitioning strategy in the previous section, there was a skill gap in extracting the
284 motor features of architectural actions.

285 To solve this problem, the researchers introduced a non-local neural network module to optimize the action
286 recognition network. The non-local network is usually embedded in vision models as a simple and efficient
287 general-purpose module that can improve the classification accuracy of images and videos (Wang et al. 2018;
288 Kong et al. 2019). Based on this module, researchers modified the original ST-GCN network and designed a
289 new dynamic skeleton model based on a non-local attention mechanism for hazard action recognition. The
290 new model was named attention module spatiotemporal graph convolutional neural (AM-STGCN) network.
291 The workflow of this new model was that it first focused on features of all joints, rather than only local
292 features of certain joints. It changed the way in which local information about workers' actions was extracted
293 to one in which global information is extracted. After analyzing the global information relationships between
294 the nodes, more effective features were obtained for key regions according to the human body partitioning
295 strategy. The network structure of AM-STGCN is shown in Fig. 4, where the model consists of nine layers

296 of spatiotemporal graph convolution operators. The first three layers had 64 output channels, the middle three
297 layers had 128 output channels, and the last three layers had a total of 256 output channels. Each layer
298 included spatial convolution operations (Conv S) and temporal convolution operations (Conv T), and residual
299 connections were added on each layer. Three attention modules were added to the temporal convolution
300 (Conv T) in the third layer of ST-GCN network in order to achieve optimal performance in the recognition
301 of workers' construction actions.

302 **3.4 The running process of the model**

303 After designing and improving the above method, a complete operation flow chart was constructed. The
304 data input of AM-STGCN model was skeleton sequence information, so in order to realize the classification
305 of workers' behavior, it would be necessary to combine Openpose skeleton point extraction algorithm and
306 AM-STGCN spatiotemporal graph convolutional neural network to jointly construct the framework of
307 construction workers' hazardous action recognition model.

308 The model utilized video streams as input, and first used Openpose for worker skeleton data extraction to
309 establish the spatiotemporal dimensional information in the human skeleton data and to construct a human
310 skeleton sequence map. Subsequently, the extracted worker skeleton information was passed into AM-
311 STGCN for learning and training of behavioral states. Finally, a SoftMax classifier was implemented for
312 behavioral result output. The specific operation flow is shown in Fig. 5.

313 **4. System building**

314 The main purpose of this study was to apply a new dynamic skeleton method to detect the action
315 characteristics of workers at construction sites. The focus of the study shifted from the static characteristics

316 of workers to the dynamic characteristics. This study was not limited to the detection of single-frame targets,
317 but also investigated the spatiotemporal relationships through video sequences for worker construction safety.
318 Therefore, all the model constructions in this study used video sequences as the data source.

319 **4.1 Datasets collection**

320 **4.1.1 Action selection for recognition**

321 In order to achieve dynamic recognition of workers' construction action processes in real construction
322 scenes, a new construction dataset based on the weights of the existing general scene of dataset was
323 established. Therefore, researchers selected the high-altitude scaffolding scenario, where fall-at-height
324 accidents are common, as a practical construction application case to study. This scenario can be adopted as
325 the basis for the study of the full-scene hazardous actions. Statistics on the causes of work-at-height injuries
326 and deaths point to unauthorized climbing, probing, leaning and crouching movements made by workers on
327 scaffolding as the main causes of accidents (HSE 2022; MOHUD 2022). Researchers selected some of the
328 hazardous actions for identification. However, it should be noted that the expansion of data types could be
329 implemented in the future. Videos were obtained for seven types of actions, including normal walking,
330 running operation, lean-over operation, scaffold climbing, hazard crossing operation, sitting on scaffolding,
331 and material handling. Fig. 6 shows some partial video intercepted segments of various types of construction
332 actions.

333 **4.1.2 Video acquisition and data processing**

334 Handheld cameras and drones were used to capture the types of workers' construction actions. Many
335 studies based on image data have shown that the difference in shooting conditions would have an impact on

336 the recognition effect (Jegham et al. 2020; Wen et al. 2022). Therefore, video clips were collected that
337 reflected the entirety of the construction operation and the camera or drone angle met the requirements of
338 multiple perspectives. Images and videos were collected under different weather and lighting conditions, and
339 effects of other factors (e.g., far view, close-up view, single-target, multi-objective, etc.) were considered to
340 avoid affecting the training effect. 2

341 The acquired videos needed to be processed in a uniform format. Researchers expanded the number of
342 samples and enhanced the data by horizontal mirror flip, by using toolkit to crop the videos into action
343 sequences of 5s each. Each action sequence contained at least one action category. A total of 8,000 action
344 sequences were obtained for skeleton extraction, containing 12,215 worker targets. The number of each
345 category of action sequences was basically kept balanced. The dataset was divided into a training set, a
346 validation set and a test set in the ratio of 6:3:1, where the test set was the original video without skeletal
347 point annotation. Table 1 shows the number of datasets for each type of action in the training set.

348 **4.1.3 Skeleton extracting and data labeling**

349 Different from the dataset format for static target annotation, the skeleton point sample data needs to extract
350 the human skeleton information of each frame image from the video. Referring to the format of the kinematics
351 skeleton dataset, a total of 18 key points of information on worker skeleton were extracted, and then, JSON
352 files of different action categories were generated through built-in data transformation algorithm module for
353 transmission to the AM-STGCN model. The file format of normal walking is shown in Fig. 7, where
354 "frame_index" is the frame index representing the skeleton data of a specific frame, and "skeleton" is the
355 skeleton joint point information of the frame. Finally, the json file was converted into npy and pkl format to

356 form a skeleton sample dataset.

357 **4.2 Recognition of workers' skeleton**

358 **In the data preparation stage, one second of video was divided into 30 frames.** So for a complete action of
359 video sequence, more than one hundred frames were generated. In another word, for one action, more than
360 100 skeleton data as shown in Fig. 7 would be generated. In these data, the extracted skeleton information of
361 each frame was different from the previous frame. By extracting the 2D coordinate information of the 18
362 nodes of the worker's skeleton frames, the extraction degree of the workers' skeleton information could be
363 maximized to ensure the accuracy of the spatiotemporal sequence when passed into ST-GCN network for
364 analysis. Fig. 8 shows the effect of extracting information about experimenter's skeleton at a certain frame
365 during the execution of the three movements (e.g., Frame 89 for the sitting on scaffolding; Frame 143 for the
366 scaffold climbing; and Frame 67 for the lean-over operation).

367 **4.3 System testing**

368 After extracting the human skeleton, the Pytorch platform was used for training the dynamic recognition
369 model. The platform used Windows 10 64 as the operating system, with a built-in NVIDIA GeForce RTX
370 1050ti graphics card and an Intel i7 processor. At the same time, CUDA and other operators were installed to
371 accelerate the model on the graphics processor (GPU). **Before the model training,** The training parameters
372 were configured according to the hardware device performance of the training platform. The specific
373 parameter settings are shown in Table 2, the learning rate was reduced with step decay during training to
374 enhance the training effect (Feng and Li 2018).

375 In order to save training time and improve the training effect, migration learning was used in the training
376 process, and the weights of ST-GCN fully trained in the Kinetics dataset were loaded as the initial training
377 weights (Weiss et al. 2016). The training is set for 300 epochs, and the loss value was calculated for each
378 completed epoch. The variation of the loss values is shown in Fig. 9 (a). It can be seen that the training loss
379 value decreased rapidly in the first 30 rounds, when the training reached 40 rounds, the learning rate
380 decreased to 0.01 to continue the training. As can be seen from Fig 9 (b), after 30 rounds of training, the Top1
381 accuracy exceeded 90% and the accuracy started to converge. After 300 rounds of training, the learning rate
382 decreased to 0.00001 and the loss value also decreased from 1.19 to about 0.14 converging. The Top1
383 accuracy could be stably maintained at about 91%, indicating that the model for hazardous action recognition
384 was completely trained.

385 **4.4 Recognition of workers' action**

386 Scaffolding climbing in the laboratory was selected as an example. Fig.10 illustrates the test process of an
387 unsafe climbing action in the laboratory, and the video results were generated from the modified ST-GCN
388 model. Figure 10(a) shows the worker skeleton information extraction using Openpose for a 5s video
389 sequence containing a worker target. It learned the process of a worker making a complete climbing
390 movement and subsequently extracted the skeleton information changes of the worker's legs and waist during
391 leg lifts and drops. Fig. 10(b) shows a visual representation of workers' skeleton extraction, which is a
392 continuous frame of human skeleton map with a duration of 5s, indicating that the video stream with skeleton
393 information is fed into AM-STGCN for model learning and action classification. Fig. 10(c) and Fig. 10(d)
394 show the key point extraction and action recognition results after using the human body partitioning strategy

395 and the non-local neural network attention module. After the whole action sequence was made, the classifier
396 was used to evaluate the action type and to output the final result.

397 **4.5 Model validation and analysis**

398 After the model was trained, divided sample test set was adopted to test the model. The essence of worker
399 action recognition was to classify types according to the set action object, and the classification task
400 commonly used Accuracy (A), Precision (P) and Recall (R) as evaluation indicators. The calculation process
401 of precision rate and recall rate is shown in Eq. (5) and Eq. (6).

$$402 \quad \textit{Precision} = \frac{T_p}{T_p + F_p} \quad (5)$$

$$403 \quad \textit{Recall} = \frac{T_p}{T_p + F_N} \quad (6)$$

404 In the equations, T_p represents the number of workers whose actions are correctly identified, F_p represents
405 the number of workers whose incorrect actions are mistakenly considered correct, and F_N denotes the number
406 of workers whose correct construction actions are evaluated to be wrong. The specific laboratory test results
407 indicators are shown in Table 3.

408 As can be seen from Table 3, the model did not miss any recognition of actions, and for the sample test set,
409 the overall accuracy of worker action recognition reached 93.39%. Among the different actions tested, the
410 method achieved high recognition recall for scaffold climbing and sitting on scaffolding, by reaching 95.48%
411 and 96.18% respectively. The recognition recall for the other five actions was lower, but the recall was around
412 93%. This could be caused by the fact that scaffold climbing and sitting on scaffolding had a distinct limb
413 performance during the movement and more characteristic changes in skeletal information. In contrast, for
414 actions such as normal walking, the temporal and spatial information of the skeleton did not change

415 significantly during the action of a few seconds, and there was partial identity of skeletal information during
416 the action. Overall, the models could achieve promising action recognition results in a laboratory setting.

417 **4.6 Case study**

418 To verify the practicality of the method, application and testing work were carried out in combination with
419 real construction scenarios. Three construction projects in Zhenjiang China were selected as real-world test
420 sites to obtain different test videos. The high-altitude scaffolding actions of workers in the close-up view and
421 far view were acquired for testing, and the action sequences included single-person and multi-person targets.
422 A total of 3,000 action sequences were acquired for method testing in both close-up and far views respectively.

423 For each action, the model extracted the skeleton information every 1 frame and combined the skeleton
424 information of 10 frames to complete the result output once. After outputting multiple skeleton information
425 in this manner, the action classification was finally completed by recognizing the results of the whole action
426 sequence. Fig. 11 shows the recognition results of workers' work video collected on site. The recognition
427 effectiveness of the method was then evaluated using the accuracy A, precision P and recall R, and the average
428 recognition time T for each frame of the video.

429 The obtained videos of construction workers from high-altitude scaffolding were used for method testing
430 and statistical analysis. A total of 1,500 videos of workers' operations with 2,137 worker targets were selected
431 for the close-up view test, and a confusion matrix was introduced to evaluate the model effects (Yang et al.,
432 2016). The test results are shown in Table 4 and Figure 12. The method can recognize all the workers' targets
433 and evaluate the actions in the close-up view. There was no omission in the identifying of actions, and the
434 average recognition time for a single frame video was 127.61ms. The classification results for different

435 actions varied slightly, with the highest recall rates for scaffold climbing actions and sitting on scaffolding,
436 which was generally consistent with laboratory. In terms of classification results, the highest number of
437 actions were misclassified as normal walking and material handling, while the lowest number of actions were
438 misclassified as scaffold climbing and sitting on scaffolding. Compared to the test results in the laboratory,
439 the overall recognition accuracy of the method in the close-up view has decreased, but it can still reach
440 90.50%, indicating that the method can better recognize hazardous actions of high-altitude scaffolding work
441 in the close-up view.

442 A total of 1,500 videos were selected for the far-view test, containing 2,291 worker targets. The effect of
443 worker action recognition in far view is shown in Table 5 and Fig. 13. In the test of the far view, due to the
444 smaller size of the worker targets within the video, the feature acquisition ability of the worker skeleton
445 information was reduced, resulting in a decrease in the recognition effect of the method compared to that in
446 close-up view. The average recognition time of a single frame action was 132.54 ms. The model still had the
447 highest recall for scaffold climbing and sitting on scaffolding, with 90.40% and 89.78% respectively. The
448 model continued to have the lowest number of actions misclassified as these two action types when
449 identifying other action types. For all other action types, the recall rate decreased to varying degrees, which
450 was roughly the same as the test results in the close-up view. The average accuracy in the far view was able
451 to maintain at 87.08%. There was no omission in the action recognition, indicating that the method was also
452 robust for high-altitude scaffolding in the far view.

453 **5. Discussion**

454 **5.1 Accuracy variation analysis**

455 The difference in accuracy may be due to the fact that the spatiotemporal information of the skeleton in
456 scaffold climbing and sitting on scaffolding was more obvious. But in actions such as normal walking and
457 running operation, the change in skeletal information during the action was approximately the same, and that
458 caused recognition errors. In order to analyze the influence of the spatiotemporal information extracted by
459 the dynamic skeleton model on the classification of different action types, two types of actions with high
460 recognition accuracy and two types of actions with low recognition accuracy were selected for comparison.
461 The effect of the spatiotemporal features of the skeleton on the accuracy of action recognition was analyzed
462 by comparing the weight parameters assigned to the body parts in different frames. The comparison results
463 are shown in Table 6 and Table 7.

464 From Tables 6 7, it can be seen that for scaffold climbing and sitting on scaffolding, where the skeleton
465 features vary significantly, the weight parameters assigned to each body part in different frames varied
466 significantly. The skeleton features were easily observed during the continuous 5s movements, so the
467 recognition accuracy was higher. For the normal walking and running operation movements, the differences
468 of weights assigned to the body parts in the different frames of the two action types were less significant,
469 resulting in similar skeletal features in several action types and hence leading to recognition errors. Overall,
470 the improved dynamic skeleton model had high recognition accuracy for complex construction actions, and
471 the recognition accuracy for actions with low complexity was also higher than previous methods.

472 **5.2 Performance evaluation**

473 In this study, a new dynamic skeleton model (AM-STGCN) was designed to identify hazardous actions of
474 construction workers. The model analyzed the spatial characteristics of workers' skeletons between different

475 frames by two convolutional modules, namely spatial convolution and temporal convolution. Researchers
476 extracted the global information of key joints in human partitioning strategies and non-local neural network
477 modules to identify complex worker actions. To verify the performance of the improved dynamic skeleton
478 model algorithm, researchers selected three models for comparisons to the improved method, including the
479 baseline model ST-GCN network(i); ST-GCN network adjusted by human body partitioning strategy only
480 (ii); and ST-GCN network modified by non-local attention mechanism only(iii). The experiments used the
481 same data and training parameters. The comparison results are shown in Table 8. For each worker action type,
482 the improved algorithm outperformed the baseline model and the other partially improved methods measured
483 by recognition accuracy.

484 **5.3 Potential limitation**

485 In terms of the overall effect, the method could achieve the recognition of the dynamic process of
486 construction workers' hazardous actions under different shooting viewpoints. Combining the results of action
487 recognition under two different viewpoints, researchers found that in the scaffold climbing recognition, the
488 method did not miss the action target, but there were cases of misjudgment. The accuracy decreased with the
489 increasing number of recognized targets and the reduction of the target size. Overall, the recognition effect
490 of scaffold climbing and sitting on scaffolding was promising. It was also found that there were mainly the
491 following reasons for misjudgment: (1) the influence of occlusions; (2) the interaction between actions.

492 The high-altitude scaffolding scenario selected for this study was complex. The junction parts of horizontal,
493 vertical and diagonal bars of scaffolding would form a complex interwoven structure of bars. When
494 construction workers were at multiple scaffolding junctions and their body parts are covered by large areas,

495 it would be difficult for the method to integrate whole-body skeleton information for action recognition
496 (Sahoo et al. 2022). In the video captured under the far view, the method sometimes misidentified the body
497 parts of construction workers as actions, as shown in Fig. 14.

498 Actions were composed of a series of consecutive behavioral gestures. There were situations where
499 different action processes had partially similar gestures, and other actions might also be interspersed when
500 engaging in specific actions, resulting in method misclassification (Vasconez et al. 2021). The detection case
501 shown in Fig. 15 was generated from a video sequence containing multiple actions, researchers added a fast
502 process of hazard crossing to the normal walking process. The model first judged the feature as a normal
503 walking but then misjudged it as a lean-over operation in the second half of the video sequence. Similarly,
504 the rare cases where a fast normal walking process was added to the lean-over operation process also caused
505 misjudgment.

506 Although the recognition result was based on the integration of the recognition results of a large number
507 of stage frames, a misjudgment in a single video frame had little impact on the overall recognition result. It
508 has been reported that for the detection of construction workers' hazardous actions, a very small number of
509 misjudgments may lead to serious injuries and fatalities (Pinto et al. 2011). Therefore, the occurrence of
510 miscalculations needed to be avoided as much as possible. In terms of the causes of miscalculation of the
511 method, factors such as obstacles and occlusions in the recognition of skeleton features by the method could
512 affect the extraction of skeleton features (Li and L 2022), while in the analysis and classification of actions,
513 the correlation and interpolation between actions could also disrupt the recognition and judgment of the
514 method for specific actions (Yang 2018).

515 In addition to selecting near and far views as the study scenes, researchers also selected single and multi-
516 person targets as the sample data set. In the test process, researchers found that the single-person action
517 recognition effect was slightly better than the multi-person action recognition effect, but the difference was
518 marginal. This might be due to the fact that most of the selected multi-person construction targets were two
519 workers, and the difference in the number of targets was not obvious. In fact, there were usually many worker
520 targets in a construction work area. To further investigate the effect of the number of workers' actions on the
521 action feature extraction ability and action recognition effectiveness, it is necessary to select construction
522 videos containing more workers' actions to study the variation of method performance in the future.

523 **6. Conclusions**

524 Researchers proposed a framework for recognizing workers' hazardous actions by fusing skeleton
525 extraction and spatiotemporal features. Openpose (i.e., skeleton point extraction network) and ST-GCN (i.e.,
526 spatiotemporal graph convolutional neural network) were designed to jointly build a dynamic skeleton model,
527 which could analyze the spatiotemporal relationship of workers' skeletons and automatically recognize
528 construction workers' hazardous actions. In order to achieve an enhanced performance of the model method
529 for construction site application, researchers made algorithm adjustments and built a dataset of real-life
530 construction scenes based on the existing public dataset. The high-altitude scaffolding scene was used as a
531 test case and tested under several challenging situations such as viewpoint change (close-up view and far
532 view) and target change (single target and multiple targets). The results showed that the method recognition
533 accuracy reached 90.50% and 87.08%, respectively, and the single frame recognition time could be controlled
534 between 127~133ms with good robustness.

535 Compared to previous studies on construction workers' hazardous actions, the contributions of this research
536 are as follows: (1) by introducing a dynamic skeleton model to analyze the spatiotemporal relationship
537 between skeleton points during workers' construction actions, researchers effectively utilized the information
538 of workers' movement characteristics, which complemented the defects of previous research methods that
539 ignored dynamic action information; (2) combining the Openpose skeleton extraction algorithm and the
540 improved ST-GCN, an operational framework for hazardous action recognition was constructed to automate
541 and visualize the process of hazardous action recognition of construction workers. This provided a technical
542 basis for managers to check workers' hazardous actions through real-time monitoring in the future; and (3)
543 through the adjustment of partitioning strategy and the addition of attention mechanism, the method enabled
544 the extraction of global features of workers' skeleton information, which effectively improved the accuracy
545 of action recognition.

546 Given the high accuracy of the method during testing, researchers believe that the proposed method has
547 good application prospects. The method provides managers a new perspective on construction site safety
548 management, rather than just focusing on the status of workers wearing safety gear or the state of worker-
549 object (e.g., material, equipment, area) interaction. This detection method for capturing hazardous
550 construction actions can be used in parallel with target detection methods such as safety protective equipment
551 to capture both dynamic and static hazards present on construction site, expanding the scope of site safety
552 detection and assisting in screening for more types of unsafe behavior events to improve worker safety. In
553 addition, the unsafe information captured could be used for job training and warning education to better
554 regulate construction workers' behavior. In the follow-up work, the feature extraction algorithm needs to be

555 further improved to reduce the influence of occlusions and different actions, to better understand the
556 spatiotemporal relationships between skeleton points, and to improve the recognition accuracy under
557 complex influencing factors. It is also necessary to improve the parameters of the dynamic skeleton model to
558 enhance the recognition speed of the model, and to expand the number of action types and worker targets to
559 enhance the applicability of the model scenes for future real-time monitoring of general scenes.

560 **Acknowledgement**

561 This research is supported by the National Natural Science Foundation of China (Grant No. 72071097),
562 MOE (Ministry of Education in China) Project of Humanities and Social Sciences (Grant No.20YJAZH034),
563 and Foundation of Jiangsu University (Grant No. SZCY-014).

564 **Data Availability Statement**

565 Some or all data, models, or code that support the findings of this study are available from the
566 corresponding author upon reasonable request.

567 **References**

- 568 Akhavian, R., and Behzadan, A.H. 2016. "Smartphone-based construction workers' activity recognition and
569 classification." *Automation in Construction*, 71, 198-209.
- 570 Bai, S., Kolter, J.Z., and Koltun, V. 2018. "An Empirical Evaluation of Generic Convolutional and Recurrent
571 Networks for Sequence Modeling." *Computer Science*, 1803.01271.
- 572 Bangaru, S.S., Wang, C., Busam, S.A., and Aghazadeh, F. 2021. "ANN-based automated scaffold builder
573 activity recognition through wearable EMG and IMU sensors." *Automation in Construction*, 126,
574 103653.
- 575 Cheng, T., Teizer, J., Migliaccio, G.C., and Gatti, U.C. 2013. "Automated task-level activity analysis through
576 fusion of real time location sensors and worker's thoracic posture data." *Automation in Construction*,
577 29, 24-39.
- 578 Detrembleur C, Van den Hecke A, and Dierick F. (2000). "Motion of the body centre of gravity as a summary
579 indicator of the mechanics of human pathological gait." *Gait & posture*, 12(3), 243-250.
- 580 Ding, L., Fang, W., Luo, H., Love, P.E.D., Zhong, B., and Ouyang, X. 2018. "A deep hybrid learning model

581 to detect unsafe behavior: Integrating convolution neural networks and long short-term memory.”
582 *Automation in Construction*, 86, 118-124.

583 Escorcia, V., Dávila, M.A., Golparvar-Fard, M., and Niebles, J.C. 2012. “Automated vision-based recognition
584 of construction worker actions for building interior construction operations using RGB-D cameras.”
585 *Construction Research Congress 2012: Construction Challenges in a Flat World*, 879-888.

586 Fang, W., Ding, L., Love, P.E.D., Luo, H., Li, H., and Peña-Mora, F. 2020. “Computer vision applications in
587 construction safety assurance.” *Automation in Construction*, 110, 103013.

588 Fang, W., Ding, L., Luo, H., and Love, P.E.D. 2018. “Falls from heights: A computer vision-based approach
589 for safety harness detection.” *Automation in Construction*, 91, 53-61.

590 Fang, W., Zhong, B., Zhao, N., Love, P.E.D., Luo, H., and Xue, J. 2019. “A deep learning-based approach for
591 mitigating falls from height with computer vision: Convolutional neural network.” *Advanced
592 Engineering Informatics*, 39, 170-177.

593 Fang, Y.-C., and Dzeng, R.-J. 2014. “A Smartphone-based Detection of Fall Portents for Construction Workers.”
594 *Procedia Engineering*, 85, 147-156.

595 Feng, Y., and Li, Y. 2018. “An overview of deep learning optimization methods and learning rate attenuation
596 methods.” *Hans Journal of Data Mining*, 8(4), 186-200.

597 Gong, Y., Yang, K., Seo, J., and Lee, J.G. 2022. “Wearable acceleration-based action recognition for long-term
598 and continuous activity analysis in construction site.” *Journal of Building Engineering*, 52, 104448.

599 Guo, G., and Lai, A. 2014. “A survey on still image based human action recognition.” *Pattern Recognition*,
600 47(10), 3343-3361.

601 Guo, H., Yu, Y., Ding, Q., and Skitmore, M. 2018. “Image-and-skeleton-based parameterized approach to real-
602 time identification of construction workers’ unsafe behaviors.” *Journal of construction engineering and
603 management*, 144(6), 04018042.

604 Han, S., and Lee, S. 2013. “A vision-based motion capture and recognition framework for behavior-based
605 safety management.” *Automation in Construction*, 35, 131-141.

606 Heinrich, H.W. 1941. “Industrial Accident Prevention. A Scientific Approach (Second Edition).”
607 McGraw-Hill, New York. <https://ajph.aphapublications.org/doi/pdf/10.2105/AJPH.22.1.119-b>

608 HSE. 2022. “Work-related fatal injuries in Great Britain.” Accessed August 3, 2022. [https://www.hse.gov.
609 uk/statistics/fatals.html](https://www.hse.gov.uk/statistics/fatals.html).

610 Jebelli, H., Ahn, C.R., and Stentz, T.L. 2016. “Fall risk analysis of construction workers using inertial
611 measurement units: Validating the usefulness of the postural stability metrics in construction.” *Safety
612 Science*, 84, 161-170.

613 Jegham, I., Ben Khalifa, A., Alouani, I., and Mahjoub, M.A. 2020. “Vision-based human action recognition:
614 An overview and real-world challenges.” *Forensic Science International: Digital Investigation*, 32,
615 200901.

616 Kim, H., Kim, K., and Kim, H. 2016. “Vision-based object-centric safety assessment using fuzzy inference:

617 Monitoring struck-by accidents with moving objects.” *Journal of Computing in Civil Engineering*,
618 30(4), 04015075.

619 Kong, T., Fang, W., Love, P.E.D., Luo, H., Xu, S., and Li, H. 2021. “Computer vision and long short-term
620 memory: Learning to predict unsafe behaviour in construction.” *Advanced Engineering Informatics*,
621 50, 101400.

622 Kong, Y., Li, L., and Zhang, K. 2019. “Attention module-based spatial–temporal graph convolutional networks
623 for skeleton-based action recognition.” *Journal of Electronic Imaging*, 28(4), 043032.

624 Lan, Z., Zhu, Y., Hauptmann, A.G., and Newsam, S. 2017. “Deep local video feature for action recognition.”
625 *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, 1-7.

626 Laptev, I. 2005. “On space-time interest points.” *International journal of computer vision*, 64(2), 107-123.

627 Li, J., Li, B., Xu, J., Xiong, R., and Gao, W. 2018. “Fully connected network-based intra prediction for image
628 coding,” *IEEE Transactions on Image Processing*, 27(7), 3236-3247.

629 Li, Z., and Li, D. 2022. “Action recognition of construction workers under occlusion.” *Journal of Building*
630 *Engineering*, 45, 103352.

631 Memarzadeh, M., Golparvar-Fard, M., and Niebles, J.C. 2013. “Automated 2D detection of construction
632 equipment and workers from site video streams using histograms of oriented gradients and colors.”
633 *Automation in Construction*, 32, 24-37.

634 MOHURD. 2022. “Circular on the production safety accidents of housing and municipal engineering in 2020.”
635 Accessed 19 June 2022. [https://www.mohurd.Gov.cn/gongkaifdzdgnr/tzgg/202006/20](https://www.mohurd.Gov.cn/gongkaifdzdgnr/tzgg/202006/20200624_246031.html)
636 [200624_246031.html](https://www.mohurd.Gov.cn/gongkaifdzdgnr/tzgg/202006/20200624_246031.html)

637 Park, M.-W., Elsafty, N., and Zhu, Z. 2015. “Hardhat-Wearing Detection for Enhancing On-Site Safety of
638 Construction Workers.” *Journal of Construction Engineering and Management*, 141(9), 04015024.

639 Pinto, A., Nunes, I.L., and Ribeiro, R.A. 2011. “Occupational risk assessment in construction industry-
640 Overview and reflection.” *Safety Science*, 49(5), 616-624.

641 Sahoo, S.P., Modalavalasa, S., and Ari, S. 2022. “DISNet: A sequential learning framework to handle occlusion
642 in human action recognition with video acquisition sensors.” *Digital Signal Processing*, 131, 103763.

643 Si, C., Chen, W., Wang, W., and Tan, T. 2019. “An attention enhanced graph convolutional lstm network for
644 skeleton-based action recognition.” *Proceedings of the IEEE/CVF conference on computer vision and*
645 *pattern recognition*, 1227-1236.

646 Simonyan, K., and Zisserman, A. 2015. “Very Deep Convolutional Networks for Large-Scale Image
647 Recognition.” *Computer Science*, 1409.1556.

648 Tang, S., Golparvar-Fard, M., Naphade, M., and Gopalakrishna, M. 2020. “Video-based motion trajectory
649 forecasting method for proactive construction safety monitoring systems.” *Journal of Computing in*
650 *Civil Engineering*, 34(6), 04020041.

651 Trabelsi, R., Jabri, I., Melgani, F., Smach, F., Conci, N., and Bouallegue, A. 2019. “Indoor object recognition
652 in RGBD images with complex-valued neural networks for visually-impaired people.”

653 *Neurocomputing*, 330, 94-103.

654 Tran, D., Bourdev, L., Fergus, R., Torresani, L., and Paluri, M. 2015. "Learning spatiotemporal features with
655 3d convolutional networks." *Proceedings of the IEEE international conference on computer vision*,
656 4489-4497.

657 Vasconez, J.P., Admoni, H., and Auat-Cheein, F. 2021. "A methodology for semantic action recognition based
658 on pose and human-object interaction in avocado harvesting processes." *Computers and Electronics in*
659 *Agriculture*, 184, 106057.

660 Wang, X., Girshick, R., Gupta, A., and He, K. 2018. "Non-local neural networks." *Proceedings of the IEEE*
661 *conference on Computer Vision and Pattern Recognition*, 7794-7803.

662 Wei, S.E., Ramakrishna, V., Kanade, T., and Sheikh, Y. 2016. "Convolutional pose machines." *Proceedings of*
663 *IEEE Conference on Computer Vision and Pattern Recognition*, 4724-4732.

664 Weiss, K., Khoshgoftaar, T.M., and Wang, D. 2016. "A survey of transfer learning." *Journal of Big data*, 3(1),
665 1-40.

666 Wen, J., Shen, Y., and Yang, J. 2022. "Multi-view gait recognition based on generative adversarial network."
667 *Neural Processing Letters*, 54(3), 1855-1877.

668 Wu, Z., Wang, X., Jiang, Y.G., Ye, H., and Xue, X. 2015. "Modeling spatial-temporal clues in a hybrid deep
669 learning framework for video classification." *Proceedings of the 23rd ACM international conference*
670 *on Multimedia*, 461-470.

671 Yan, S., Xiong, Y., and Lin, D. 2018. "Spatial temporal graph convolutional networks for skeleton-based action
672 recognition." *Thirty-second AAAI conference on artificial intelligence*, 180.07455.

673 Yang, J., Shi, Z., and Wu, Z. 2016. "Vision-based action recognition of construction workers using dense
674 trajectories." *Advanced Engineering Informatics*, 30(3), 327-336.

675 Yang, J. 2018. "Enhancing action recognition of construction workers using data-driven scene parsing."
676 *Journal of Civil Engineering and Management*, 24(7), 568-580.

677 Yu, Y., Guo, H., Ding, Q., Li, H., and Skitmore, M. 2017. "An experimental study of real-time identification
678 of construction workers' unsafe behaviors." *Automation in Construction*, 82, 193-206.

679 Zhang, P., Lan, C., Xing, J., Zeng, W., Xue, J., and Zheng, N. 2019. "View adaptive neural networks for high
680 performance skeleton-based human action recognition." *IEEE Transactions on Pattern Analysis and*
681 *Machine Intelligence*, 41(8), 1963-1978.

682 Zhao, R., Xu, W., Su, H., and Ji, Q. 2019. "Bayesian hierarchical dynamic model for human action recognition."
683 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 7733-7742.

684 Zhou, B., Andonian, A., Oliva, A., and Torralba, A. 2018. "Temporal relational reasoning in videos."
685 *Proceedings of the European conference on computer vision*, 803-818.

686 Zhou, K., Wu, T., Wang, C., Wang, J., and Li, C. 2020. "Skeleton Based Abnormal Behavior Recognition
687 Using Spatio-Temporal Convolution and Attention-Based LSTM." *Procedia Computer Science*, 174,
688 424-432.