

1 Advancing healthcare practice and education via data sharing: Demonstrating the utility of open data by training an
2 artificial intelligence model to assess cardiopulmonary resuscitation skills

3 *Merryn D. Constable,¹ *Francis Xiatian Zhang,² Tony Conner,³ Daniel Monk,³ Jason Rajsic,¹ Claire Ford,³ Laura
4 Jillian Park,³ Alan Platt,³ Debra Porteous,³ Lawrence Grierson,⁴ Hubert P. H. Shum²

- 5 1. Department of Psychology, Northumbria University, UK.
6 2. Department of Computer Science, Durham University, UK.
7 3. Department of Nursing and Midwifery, Northumbria University, UK.
8 4. Department of Family Medicine, McMaster University, Canada.

9 **MDC and FXZ declare equal contributions.*

10

11 Corresponding Author:
12 Merryn Constable
13 Northumberland Building
14 College Lane
15 Newcastle upon Tyne
16 UK, NE1 8SG
17 merryn.constable@northumbria.ac.uk
18 Word count: 8280

19
20 ORCIDs:
21 MDC: 0000-0001-5149-5670
22 FXZ: 0000-0003-0228-6359
23 AC: 0000-0001-9762-5573
24 DM: 0000-0002-1358-0425
25 JR: 0000-0002-3771-5216
26 CF: 0000-0003-1566-2704
27 LP: 0000-0002-2427-7659
28 AP: 0000-0001-5646-8671
29 DP: 0000-0001-8952-9119
30 LG: 0000-0003-0739-5976
31 HPHS: 0000-0001-5651-6039

32 Abstract

33 Health professional education stands to gain substantially from collective efforts toward building video databases of
34 skill performances in both real and simulated settings. An accessible resource of videos that demonstrate an array of
35 performances – both good and bad - provides an opportunity for interdisciplinary research collaborations that can
36 advance our understanding of movement that reflects technical expertise, support educational tool development, and
37 facilitate assessment practices. In this paper we raise important ethical and legal considerations when building and
38 sharing health professions education data. Collective data sharing may produce new knowledge and tools to support
39 healthcare professional education. We demonstrate the utility of a data-sharing culture by providing and leveraging a
40 database of cardio-pulmonary resuscitation (CPR) performances that vary in quality. The CPR skills performance
41 database (collected for the purpose of this research, hosted at UK Data Service's ReShare Repository) contains
42 videos from 40 participants recorded from 6 different angles, allowing for 3D reconstruction for movement analysis.
43 The video footage is accompanied by quality ratings from 2 experts, participants' self-reported confidence and
44 frequency of performing CPR, and the demographics of the participants. From this data, we present an Automatic
45 Clinical Assessment tool for Basic Life Support that uses pose estimation to determine the spatial location of the
46 participant's movements during CPR and a deep learning network that assesses the performance quality.

47 *Keywords:* Healthcare Professional Skills, Nursing Skills, Competency-Based Education, Deep Learning,
48 Pose Estimation, Healthcare Skills Databases.

49

50 Advancing healthcare practice and education via data sharing: Demonstrating the utility of open data by
51 training an artificial intelligence model to assess cardiopulmonary resuscitation skills

52 Healthcare professionals, across all disciplines, must master many movement-based technical skills to
53 ensure positive outcomes for patients and avoid injury to themselves. Improper patient lifting techniques can cause
54 harm to both patient and practitioner; inaccurate intubation can damage vocal cords, and inefficient surgeries are
55 linked with poor operative results. Due to the importance of quality care, training programs for healthcare
56 professionals around the world have started to shift curricula towards an education paradigm known as competency-
57 based education (CBE), which eschews time-based schedules of learner progression in favour of systems of
58 matriculation that depend on direct observation of learner competence across a pre-defined set of professionally
59 relevant activities (Frank et al., 2010; Harden, 2007).

60 The subjective nature of competency-based assessment provides a key challenge. Where some healthcare
61 disciplines [e.g. robotic-assisted surgery (El-Sayed et al., 2024)] have devoted more research toward understanding
62 how objective measurements of movement relate to expertise and the clinical outcomes of patients, most disciplines
63 have lacked the technology to easily extract such information to determine how and when movement patterns matter
64 for the patient and practitioner. Innovations in computer science that provide an opportunity to study technical skills
65 in simulated and real environments without needing specialised recording devices provide this much-needed
66 opportunity for healthcare skills more generally. Thus, empirically validated competency targets could be
67 established by analysing video datasets that can be easily accumulated during training and practice across multiple
68 healthcare disciplines.

69 A second challenge of CBE is that it is human resource intensive. CBE requires frequent formative and
70 summative assessments to be implemented by healthcare educators who may have contemporaneous patient care
71 and educational commitments. With appropriate objective competency thresholds established, automated feedback
72 and assessment systems could provide the opportunity for self-directed deliberate practice, lessening the amount of
73 formative feedback required from a coach and allowing students who need more practice time to have that
74 opportunity. These formative assessment techniques could also flag when someone may meet the required
75 competency and is ready to be assessed for progression. Such tools would also represent cost savings given that they

76 would reduce the human resource requirements of training, which is of rising concern within the educational sector
77 (Castillo et al., 2019).

78 Yet, the above benefits cannot be achieved without substantial investment in accumulating and sharing
79 relevant data. Thus, this paper represents a call for collective effort from healthcare institutions (both educational
80 and providing) toward amassing skills performance data that can be used to determine clear competency thresholds
81 that healthcare students, educators, and professionals can target to enhance patient and practitioner outcomes. With
82 recent innovations, the computational barrier to processing and analysing such complex data no longer exists;
83 instead, the barrier is access. The surgical field has a much longer history of evaluating human movement-based
84 skills on a kinematic level and has begun to look at implementing pose estimation and artificial intelligence
85 techniques to better understand surgical skill and enhance education (Constable, Shum, et al., 2024; Likosky et al.,
86 2021). As such, surgical tool and procedure datasets (Bouget et al., 2017; Srivastav et al., 2018) are currently being
87 collated to accelerate the development of specialised pose estimation algorithms within the discipline. We suggest
88 that all healthcare disciplines could benefit from such an agenda and call for a concerted effort across healthcare and
89 science, more broadly, to develop policies and practices that provide the means to develop technologies to support
90 healthcare trainees, professionals and patients alike.

91 By creating repositories of skills performance data alongside other factors of importance (e.g.
92 demographics, educational level, patient factors), collaborative efforts from researchers in the fields of computer
93 science, data science, human movement and healthcare education can begin to meet the first challenge of
94 establishing objective, validated, and measurable competency-based thresholds. Furthermore, established
95 competency thresholds and videos of skills performances will subsequently assist computer scientists in building
96 rigorous skills assessment tools. With the two above challenges in mind, the present work aims to demonstrate how
97 the accumulation of skills performance data sets can be combined with innovations in computer science
98 (specifically, pose estimation and deep learning for the classification of expertise) to study technical competencies
99 within healthcare professional education and provide automated means of assessing skills performance for the
100 educational setting.

101 In computer vision, pose estimation refers to tracking movement from videos or pictures. Here, we focus on
102 deep learning techniques which involve 'training' an Artificial Neural Network (ANN) on annotated videos or stills

103 that indicate the relevant objects or body parts to be tracked (Ionescu et al., 2014; Lin et al., 2014). After training,
104 the ANN can identify the relevant body parts or objects and, in turn, 'poses' within videos without human assistance.
105 The kinematic parameters of a given movement can be calculated from body part locations and the frame rate of the
106 video. Although these techniques require more powerful hardware than a standard computer, acquiring and setting
107 up such tracking is more accessible than using specialised motion tracking cameras because it can be done with
108 traditional video cameras. Further, pose estimation algorithms and toolboxes are advancing rapidly, with multi-
109 person pose estimation possible with pre-trained networks at near real-time speeds (Huang et al., 2020), meaning
110 that real-time feedback could be provided.

111 The kinematic data obtained from pose estimation can be used to levy an assessment of movement or to
112 support formative development. The kinematic data could be used directly or by classifier algorithms to provide
113 further assessment. For example, optimal posture during CPR requires the practitioner's shoulders to be directly
114 above the patient. Deviation from this optimal posture could be relayed back to the trainee to allow them to explore
115 and feel their approach. Pose estimation data could complement instructor observations or data from computerised
116 manikins, allowing trainees to refine psychomotor techniques, improving compressions and protecting practitioners
117 from injury. Such applications of kinematic feedback have been repeatedly demonstrated in high-stakes sports
118 training environments (Giblin et al., 2016; Glazier, 2021) and surgical training environments (Judkins et al., 2008)
119 and are consistent with fundamental teaching and learning theory (Ericsson, 2004; Platt et al., 2021).

120 Classifier algorithms can extend on the above by making decisions along a given dimension. Here, neural
121 networks are trained on relevant data that could be used to decide on competence (e.g., pose estimation data and
122 evaluation data). Classifier algorithms have been demonstrated as useful and highly accurate for both assessment
123 purposes and for highlighting aspects of the skill for improvement (e.g. suturing, (Ismail Fawaz et al., 2018)). For
124 example, if competency is of interest, the ANN must be trained in example performances to learn patterns
125 representing good, adequate, or poor performance. The neural network can then identify the competency level
126 displayed in new videos based on learned patterns. Nevertheless, it is vital to be cautious in the development of these
127 algorithms: the decisions risk being biased if the training data does not accurately represent the to-be-assessed data
128 or the intended purpose (Veale & Binns, 2017); decisions may be based on parameters that co-occur with the

129 classified groups but are not meaningful to the decision-making process. Considering the quality, depth and breadth
130 of the training data can assist in protecting against such concerns.

131 While such tools are more commonly applied directly to tasks that are weighted toward psychomotor ability,
132 these tools could also be used to assess non-technical skills. For example, situational awareness is critical in many
133 healthcare contexts but is challenging to teach and assess. Simulation-based educational interventions yield better
134 outcomes in situational awareness training (Walshe et al., 2019). In high-fidelity clinical simulations, final-year
135 nursing and paramedicine students perceived that using eye-tracking technology combined with video debriefing
136 assisted in their development of situational awareness (O'Meara et al., 2015). This intervention required participants
137 to wear eye-tracking glasses, which may be challenging to implement for many training programmes. However, recent
138 advances in pose estimation show that human attention can be tracked and modelled within a task space with
139 information about head pose and orientation. Of course, this is less precise than eye-tracking. Nevertheless, this
140 technique has been demonstrated to be a viable method of assessing concentration loss, collaborative attention and
141 stress levels more generally for industrial applications (Lagomarsino et al., 2022), suggesting that pose estimation
142 techniques could assist in the tracking and understanding situational awareness for healthcare applications.

143 **Considerations in implementing pose estimation and classifier algorithms**

144 Such techniques have limitations, especially in real healthcare settings, which are often busy and complex
145 environments. It is possible that occluded points will not be estimated or will be estimated with lower accuracy. Using
146 multiple cameras (Kocabas et al., 2019) will alleviate this issue. A range of gap-filling strategies, including those
147 employing ANNs (Kanazawa et al., 2018), can also be used to estimate missing data. If high precision is still needed,
148 a hybrid approach could be used with specialist simulators combined with a visual approach. For example, manikins
149 that track compression depth and rate (e.g. QPCR manikins from Laerdal) could be used simultaneously with video
150 data, which allows posture to be tracked. With further advances in computer vision techniques, accuracy thresholds
151 across all relevant dimensions may reach a level where a hybrid approach may not be needed for high-precision cases.

152 Just as humans make mistakes, algorithms can too. A recent systematic review evaluating the use of
153 machine learning for classifying surgical expertise indicated typical accuracy rates of over 80% (Lam et al., 2022).
154 Accuracy rates are likely to rapidly improve as appropriate video data is obtained for development purposes;
155 nevertheless, even with accuracy rates at 80%, trust and acceptance of the use of classifier algorithms for assessment

156 would be enhanced with a 'human-in-the-loop' approach (Enarsson et al., 2022) that provides a means of cross-
157 checking. With human oversight, classifier algorithms could be used to track progression against milestones for
158 healthcare trainees and flag when a student is ready to be formally assessed against competency thresholds by a
159 human assessor. Importantly, the human in the loop must be participatory; decisions should not blindly follow the
160 classifier's recommendation and decisions should be monitored if this strategy is used (Kazim et al., 2021).

161 Considering potential bias and how the algorithm makes decisions is ethically crucial in developing
162 classifier algorithms. If the data used to train the algorithms is biased, then this bias will be evident in any decisions.
163 This point highlights the importance of accessing wide and diverse data sets. Of course, it can be challenging to
164 determine if a data set is biased for the purposes for which it is being used. Classifier algorithms that output a
165 meaningful description of the decisional parameters [Explainable Artificial Intelligence, (Taylor & Taylor, 2020)],
166 combined with a human-in-the-loop approach, can assist in mitigating any potential bias. Unfortunately, classifier
167 algorithms often disadvantage underrepresented groups (Holstein et al., 2019). Thus, carefully considering the
168 training data is essential to ensure the algorithms' fair and equitable decisions (Corbett-Davies et al., 2023; Veale &
169 Binns, 2017).

170 In some cases, it may be reasonable to develop algorithms that ignore protected characteristics to ensure
171 that the algorithms cannot learn the systematic biases present in society; however, in many cases, this may eliminate
172 important information which impacts conclusions made (Hajian & Domingo-Ferrer, 2013). For example, in the case
173 of physical disability, movement patterns may be markedly different. Ignoring that characteristic of the performer
174 may lead to improper classification of competence, particularly when they have found a viable movement pattern
175 that deviates from the sample norm. If an algorithm considers such diverse information, then the outcome will likely
176 be fairer (Veale & Binns, 2017). Fairness and equity will require careful consideration of the context. Ensuring that
177 systematic discrimination is avoided during implementation must be a high priority (Hagendorff, 2019). For
178 thorough reviews and strategies for anti-discrimination in machine learning, see Veale & Binns (2017) and
179 Hagendorff (2019).

180 **Data acquisition, storage and use considerations**

181 Beyond the demonstrated application, opening doors to movement analysis in real healthcare settings
182 provides opportunities to understand how movement patterns relate to patient outcomes in a given environment

183 through data mining. Such a pursuit could lead to data-driven support for change to environments and policies that
184 support practitioners. Furthermore, it has been challenging in some healthcare professions to link competency
185 milestones with patient outcomes (Kendrick et al., 2023). Large datasets that speak to the trajectory of expertise,
186 patient outcomes, and practitioner injury would allow the profession to develop competency thresholds informed by
187 empirical evidence where necessary. Yet, data availability is currently limited, and initiatives to accumulate and
188 share such data require careful consideration.

189 Where the potential benefit of collecting and evaluating video performance data to enhance healthcare
190 appears substantial, there are also significant barriers. Installing cameras to monitor patient care (e.g. operating room
191 black boxes) is becoming more commonplace; however, healthcare professionals have raised concerns over data
192 safety and litigation. Cultural factors seem to play a role in such concerns. For example, Canadian healthcare
193 professionals are more concerned (Gordon et al., 2022) than Danish healthcare professionals, who indicated
194 relatively high opinions toward monitoring initiatives (Strandbygaard et al., 2022). Where it is most certainly
195 essential to consider acceptance to maintain an environment of trust, it is also important to note that regardless of
196 perception, video data most often supports healthcare professionals from a legal standpoint (van Dalen et al., 2019)
197 and thus is most likely to offer protection in a litigious environment.

198 Legal policies around recording healthcare professionals, students and patients will likely differ
199 considerably between governing bodies. Still, there has been a considerable cultural shift toward prioritising the data
200 privacy of individuals and ensuring that personal data is protected, with the General Data Protection Regulation
201 (GDPR) being one of the most comprehensive examples globally. Within the (evolving) legal framework,
202 institutions should develop strong policies to ensure that video footage and any associated data regarding outcomes
203 is recorded, stored, and used ethically and legally.

204 Data minimisation is one principle that may arise in both legal and ethical frameworks globally (e.g.
205 Europe's GDPR and California's CPRA) that allows for a balance between the processing of personal data and data
206 privacy (Goldsteen et al., 2022). This principle requires that the minimum amount of data be collected and
207 processed. However, this practice could hinder finding important patterns between clinical practice, patient and
208 practitioner factors, and patient outcomes without careful consideration. Identifying patterns that result in
209 incremental gains to patient outcomes is important in healthcare situations. Where identifying patterns in such large

210 datasets may not have been previously possible, machine learning opens the doors to such pursuits, and the
211 importance of a particular variable may not be known ahead of time.

212 'Big Data' seems at odds with principles of data minimisation, and indeed, recently the UK's Information
213 Commissioner's Office has indicated that data minimisation should be applied at both training and inference stages
214 of machine learning (Kazim et al., 2021) with global regulatory bodies also providing guidelines for ethical use. In
215 response, new methods are being developed to minimise the data required whilst maintaining accuracy or providing
216 evidence that the minimum amount of data was used to achieve the aims (Goldsteen et al., 2022). As the use of
217 artificial intelligence technology becomes more common, discussion of how to balance data privacy with scientific
218 advancement will likely become a particularly hot topic from an ethical standpoint. Regardless, adherence to the
219 Declaration of Helsinki (World Medical Association, 2013), also requires that consent be obtained for the use of any
220 identifiable human data.

221 Privacy rights can be protected by depersonalising data, which is already a commonly implemented ethical
222 practice. It is not typical for science to be interested in the identity attached to data; therefore, depersonalisation
223 would rarely impact the potential scientific gain. At a minimum, depersonalisation can be achieved by removing
224 obviously identifying information such as names or faces. Artificial Intelligence (AI) tools can also assist with this,
225 as we have used in the present paper. However, it is important to consider that richer data sets may provide
226 information that could be combined to identify a participant. For example, rich patient data, or even kinematic data
227 that is only spatially or temporally based, could be backward engineered to identify the source. While this is
228 unlikely, given that the motivation to do so would be low, it does pose a risk that should be carefully assessed.

229 Where a patient is concerned, consent, confidentiality, anonymity, and protection needs should be carefully
230 considered from both an ethical and legal standpoint, as the data could be considered particularly sensitive. Video
231 data of procedures need not always be added to a patient's medical record if the video is collected solely for quality
232 improvement, and the video would not be used to inform patient care (van Dalen et al., 2019).

233 **Moving toward a culture of data sharing**

234 In sum, collective efforts toward accumulating skills performance data alongside relevant demographic or
235 patient data have the potential to advance healthcare professional education substantially. In doing so, empirically

236 validated competency targets can be established, and computer scientists can leverage the data to develop high-
237 quality, robust tools to facilitate healthcare professionals and trainees to learn and maintain healthcare skills.

238 Establishing a culture of data sharing does face several challenges, although policies are rapidly evolving to
239 support such initiatives. As policy is being developed globally, at a local level, healthcare and educational
240 institutions may wish to establish their collaborative policies within current legal and ethical frameworks to advance
241 research around healthcare professional skills education. In parallel, the educational and research community could
242 benefit from establishing consensus on what may be considered informative data and guidelines for data
243 organisation to facilitate access. By carefully considering the type of data required and its organisation, the field will
244 balance the opportunity that exploratory work can bring with data minimisation principles. Established skills
245 performance databases will then invite interdisciplinary researchers to engage with the institutions that have a stake
246 in the outcomes (healthcare and educational institutions) to advance the field.

247 Here, we demonstrate the process in practice with an established interdisciplinary research team of
248 healthcare professional educators, human movement scientists, and computer scientists. We provide a database of
249 CPR skills and performances of varied experts who have been assessed for performance quality by two experts.
250 Then, using computer vision and machine learning, we leverage this data to demonstrate the possibility of an
251 Automatic Clinical Assessment tool for Basic Life Support.

252

253 **Automatic Clinical Assessment for Basic Life Support**

254 Early recognition of a cardiac event and quick application of CPR with high-quality chest compressions is
255 advocated internationally (Berg et al., 2023; Merchant et al., 2020; Olasveengen et al., 2021; Resuscitation Council
256 UK, 2021). Indeed, there is consensus within the evidence that high-quality CPR improves outcomes for patients in
257 cardiac arrest (Gates et al., 2015). Thus, basic life support (BLS) education is essential to healthcare professional
258 education. It is also a primary feature of first aid training delivered to individuals who are not healthcare
259 professionals or trainees. High-quality chest compressions are reflected in hand and elbow position, compression
260 depth, rate and recoil, alongside consideration of the angle of compression force application and rescuer safety
261 (Resuscitation Council UK, 2021). Feedback during training has been found to significantly enhance compression
262 quality (Baldi et al., 2017). Studies using Kinect depth cameras and pose estimation techniques show promising

263 tracking and feedback provision capabilities (Lins et al., 2019; Xie et al., 2020). Indeed, real-time feedback from
264 Kinect can significantly improve chest compression quality for rescuers who weigh below 71 kilograms (Wang et
265 al., 2018). The body weight limit here may be attributed to a higher quality baseline in those with higher body
266 weights. However, such a limit also indicates the importance of accessing diverse datasets. Other Kinect-based
267 studies have demonstrated comparable benefits to real-time visual feedback for skills development in CPR
268 (Semeraro et al., 2013). Although these Kinect-based studies demonstrate the utility of providing feedback for
269 training purposes, Kinect does require specialist cameras and sensors, which may limit uptake.

270 Pose estimation can be performed using a computer and any camera. Initial work comparing expert ratings
271 and evaluations from such computer vision techniques has shown promising results. In a study of arm angle (and
272 chest-to-chest distance between team members), pose estimation was thought to be more precise in estimating arm
273 angle than experts (Weiss et al., 2023). Here, we seek to demonstrate how deep learning techniques can provide an
274 automatic assessment of CPR technique against a comprehensive set of metrics that assess both the quality of
275 movement concerning the CPR performance and the postural safety of the performer.

276 Alongside this paper, we provide a CPR performance data set comprising a range of competencies for use
277 to advance research that understands technical competency and builds tools to support the development of such
278 competencies. The data set includes demographic information, self-ratings of confidence and frequency of
279 performance, and two expert evaluations of the performance. The database contains video data of CPR from
280 multiple angles with a checkerboard allowing for 3D reconstruction.

281

282

Methods

283 Participants

284 Participants were recruited on three different days. Participants were recruited from Northumbria
285 University's Department of Nursing and Midwifery on the first day via the researchers' networks. Recruitment
286 resulted in 22 participants with varied expertise, ranging from complete novices who had never performed CPR
287 before to individuals with extremely high levels of expertise in CPR (trained professionals and educators who
288 regularly perform CPR). On the second day, 20 students who attended a skills event held by the Department were
289 recruited on a voluntary basis as an opportunity to practice CPR and contribute to research. They were of varied skill

290 levels, with some having previously undergone training and others not; all were students of the Department. On the
 291 third day, the recording session was set up to coincide with a first-year training session; 10 first-year students in the
 292 Department of Nursing and Midwifery and one non-student in the Department were recruited. Thus, data from 53
 293 participants was collected. Authors [MDC, XZ, TC, DM, JR, CF, LJP, AP] participated. All participants provided
 294 informed consent and indicated how they would like their data to be used (Video or evaluative data available only to
 295 the research team/available in a safeguarded science repository for scientific use). For the present paper, we have
 296 used data from all participants; where participants consented (n = 40) videos with faces digitally obscured have been
 297 uploaded to UK Data Service Reshare (Constable, Zhang, et al., 2024). Participants who did not consent to their data
 298 being used outside the research team have not been included in the repository. Researchers or educational
 299 professionals may access the repository in a safeguarded manner subject to adhering to the terms and conditions of
 300 the repository. To gain access researchers must email the data controller (MDC) stating their status as a researcher
 301 and their intention for the data; they will then be granted access. Northumbria's Ethics System (No. 44602) approved
 302 the research, and all research was performed per the Declaration of Helsinki.

303 The average age of participants was 33.60 years (SD = 13.00); and 14 were men, and 39 were women (self-
 304 declared), see Table 1 for age by gender. Participants self-reported confidence in performing CPR, ranging the full
 305 spectrum of possible responses from Very Confident to Very Unconfident (5-point Likert scale), with the median
 306 response being 'Somewhat confident'. Self-reported frequency also ranged the full spectrum of possible responses
 307 from Very Frequently to Very Infrequently (5-point Likert scale), with the data being skewed toward infrequent
 308 performance (median response = Very Infrequently). The skew in the data reflects that both students and clinicians
 309 use CPR skills relatively infrequently and thus require regular refresher training (Oermann et al., 2011).

310 *Table 1. Average age (Standard deviation in parentheses) by self-identified gender.*

Gender	Male	Female
Age	42.57(15.02)	30.38 (10.67)

311

312 **Data Protection**

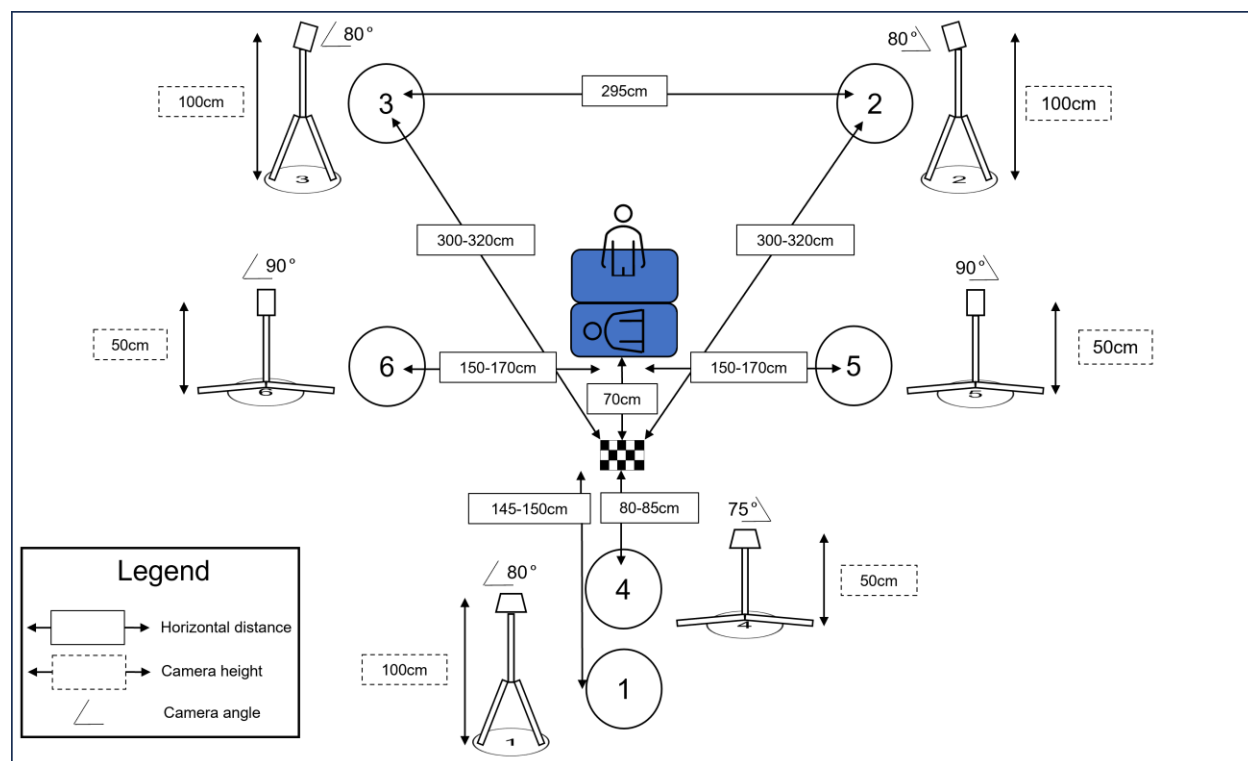
313 We obtained informed consent from participants who were able to indicate how they would like their data
 314 to be used and stored. Further, the data has been depersonalised by digitally obscuring their faces. Access is granted

315 with safeguarding protections such that users must agree to the terms and conditions of the repository. Importantly,
316 these terms and conditions require users to be registered, to only use the data for research or learning purposes, and
317 to maintain the confidentiality of the participants. Any participants who indicated that they did not want their data to
318 be used outside of the research team have not been included in the data set available to researchers external to the
319 research team. Furthermore, participants could elect to have only their videos or evaluations shared should they
320 wish.

321 **Recordings**

322 Each person was recorded using 6 Go-Pro Cameras. The first camera was set up to have a wide frontal view (see
323 Figure 1). Cameras 2 and 3 were placed behind the participant, offset to the right and left. Camera 4 was placed in
324 front of Camera 1 to provide a lower and closer frontal view. Cameras 5 and 6 were placed perpendicularly to the
325 direction the participant was facing in line with the participant. A checkerboard was placed in front of a QCPR
326 manikin and in view of all six cameras as a common landmark. The CPR space was defined for the participant with
327 two foam mats. One for them to kneel on, the manikin was placed on the other.

328 Participants were asked to perform 4 sets of 30 chest compressions for each recording with a short pause in
329 between to rest. Participants were asked at the beginning and end of the task to clap. This clap was used to calibrate
330 cameras in time.



331

332 *Figure 1.* The recording space. Circles depict the location of the cameras. The checkerboard was placed in front of
 333 the task space with one foam mat for the manikin and one foam mat for the participant. Approximate distances
 334 between cameras are provided, although there was some slight variation between days.

335 **Ratings**

336 An evaluative framework (see supplementary material) was developed in consultation with the BLS experts
 337 on the team (AC, DM). Both experts have been teaching BLS for over 20 years in clinical and educational settings,
 338 and the UK Resus Council recognises both as Advanced Life Support Instructors, representing exceptional expertise
 339 in the field. Given that the present data was collected within the UK educational system, the evaluative dimensions
 340 were informed by guidelines from the Resuscitation Council UK (2021). Additional evaluative dimensions were
 341 included to reflect good posture and technique taught to maintain endurance, reduce the likelihood of injury, and
 342 prevent fatigue. Each evaluative dimension represented a factor that would be currently instructed in the educational
 343 setting; nevertheless, in practice, each factor is not equally important for patient outcomes as indicated by the
 344 International Liaison Committee on Resuscitation's recommendations, which are updated yearly based on
 345 cumulative science (Berg et al., 2023).

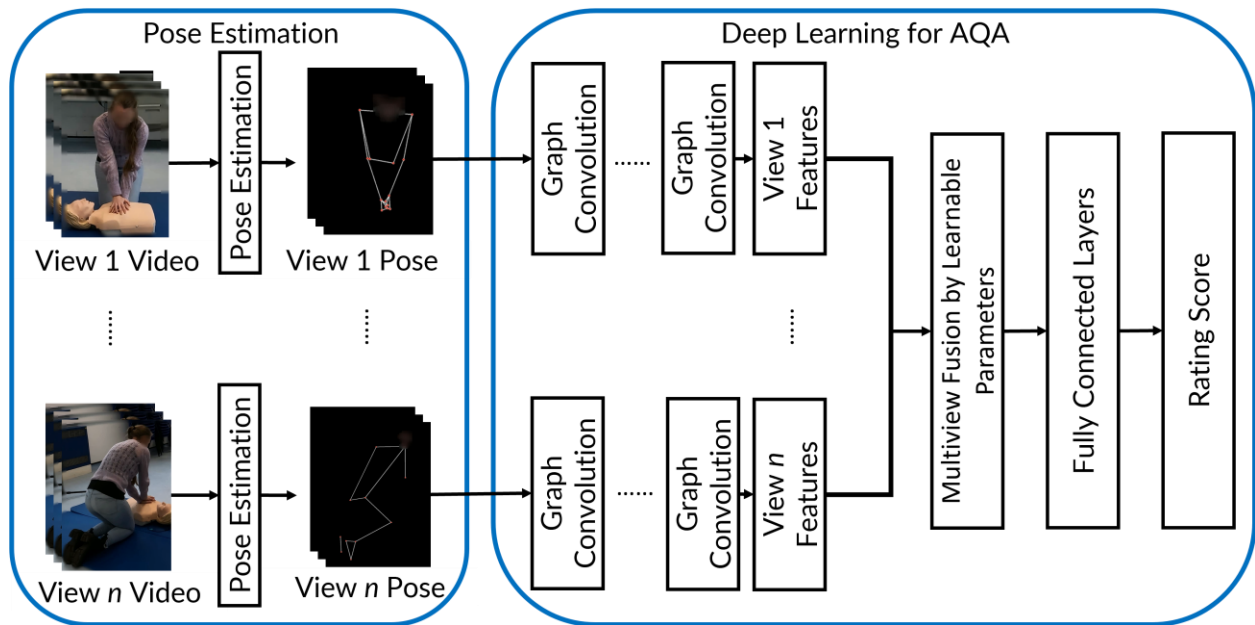
346 For each cycle (4 per participant), one point was provided for good form in the following evaluative dimensions:
347 Hand Position (Centre of the chest, within one average hand margin), Arm Position (Straight Arms, minimal flex in
348 elbows or minimal variability in elbow joint angle), Shoulder Position (Over patient, the line from centre of patient
349 to shoulders should be perpendicular), Depth of Compressions (5-6cm), Rate of Compressions (100-120 per
350 minute), Release (complete recoil of chest, hands return to neutral start point). A metronome was used to assist in
351 evaluating the rate of compressions. Experts also coded for incorrect form (see supplementary material – evaluative
352 checklist), such that if there were multiple ways a participant could exhibit poor form, that will be reflected in the
353 data (e.g. Depth of compressions could be either too shallow or too deep). Overall ratings (Excellent, Good,
354 Borderline, Poor, Unacceptable) for each cycle were also provided.

355 The two expert raters initially rated alone and then resolved any discrepancies to provide an agreed rating.
356 To determine rater agreement when raters were rating alone we calculated weighted Cohen's kappa for overall
357 ratings for each cycle (rated: Unacceptable, Poor, Borderline, Good, Excellent). Overall, raters were in moderate
358 agreement for Cycles 2, 3, and 4 when they rated alone, $\kappa_s = 0.550, 0.567, 0.518$, respectively. Agreement was poor
359 for Cycle 1, $\kappa_s = 0.204$, potentially reflecting inconsistencies in performance during the initial cycle, which could
360 reflect a 'practice' run.

361 **Automatic Clinical Assessment**

362 Our methodology systematically assesses CPR techniques using human motion data, mirroring expert
363 evaluations while leveraging the advantages of automation. Our approach to assessing clinical techniques includes
364 two main components: markerless pose estimation and a deep learning network designed specifically for Automatic
365 Action Quality Assessment (AQA), as shown in Figure 2. The first step in our framework is to use markerless pose
366 estimation to capture the 3D positions of a participant's joints from different angles in the video. This process is notable
367 because it does not rely on physical sensors or markers attached to the participant. Instead, it directly analyses the
368 video frames to identify and track the movements of the joints. Following this, the pose information that has been
369 extracted is input into a deep learning network. This network is trained to assess the quality of CPR performance
370 against predefined criteria as same as those used in manual expert assessments. The network produces ratings for
371 various aspects of the CPR technique, thus providing an objective, automated evaluation of the participant's skill level.

372



373

374 *Figure 2.* Overview of the framework for our automatic AQA. The process begins with extracting joint motion
 375 information from raw video footage captured from multiple viewpoints, using pose estimation techniques. Following
 376 this, deep learning algorithms analyse the spatial-temporal features of the extracted data from each angle. The
 377 system then integrates these features across different views to accurately predict the performance ratings for each
 378 assessed skill.

379 **Markerless Pose Estimation**

380 We employed MediaPipe (Lugaresi et al., 2019), a framework that enables machine learning models to
 381 interpret and analyse human motion video data. This approach allowed us to capture the movements of various
 382 joints, outputting their positions in a three-dimensional space (X , Y , and Z coordinates) relative to a standard'
 383 world coordinate system,' which provides a consistent frame of reference for movement analysis. The results of our
 384 pose estimation are organised in a structured format, denoted as $P_k \in R^{T \times N \times C}$, where k denotes the camera
 385 viewpoint, T denotes the total number of video frames (or the duration of the video), N denotes the count of
 386 distinct joints tracked, and C denotes various data features for each joint, including their spatial coordinates and the
 387 confidence level of these estimations. An illustrative example of how we visualise this pose estimation data can be
 388 seen in Figure 3.



389

390 *Figure 3.* An example of the estimation visualisation. Key joint landmarks relevant to the clinical technique are
 391 clearly captured.

392 **Deep Learning for Automatic Action Quality Assessment**

393 **The Graph Representation of Human Skeleton**

394 After obtaining the pose estimation outputs, denoted as P_k , which detail the positions of various joints over
 395 time, we conceptualise the human skeleton as a graph structure. This graph-based representation, $G(V, E)$, allows
 396 our neural network to incorporate the anatomical and biomechanical constraints inherent to human movement (Feng
 397 et al., 2022). In this graph, the set of nodes V represents the joints, indexed from 1 to N , where each node v_i
 398 corresponds to a specific joint. The edges E , represent the connections (i.e. bones) between these joints, such as
 399 bones or ligaments, defined within a set S that specifies which pairs of joints are connected, which promotes the
 400 network to represent the physical structure of the human body.

401 **The Multiview Spatial Temporal Graph Convolutional Network**

402 After defining the graph representation based on the pose estimation results, we developed a multiview
 403 neural network tailored for Automatic Action Quality Assessment (AQA), using the estimated human poses as its
 404 foundation. To achieve this, we adopted Spatial Temporal Graph Convolutional Networks (STGCN) (Yan et al.,

2018), configuring the network with a five-layer architecture to robustly capture and model the dynamics of the pose estimation data from each camera viewpoint. This approach incorporates the skeleton graph structure dynamics of the human body to constrain our deep learning model for joint motion analysis. The detailed operation for the input pose estimation P_k as follow:

$$H_k = \Lambda^{-1/2} A \Lambda^{-1/2} P_k W_k,$$

where H_k denotes the learned hidden feature for camera viewpoint k , W_k denotes the learnable parameters in our STGCN. A and Λ denote the adjacency matrix and its normalised form, respectively. The value of A is defined such that $A_{i,j} = 1$ if $v_i v_j \in E$, which introduces the semantics of our defined human skeleton graph into deep learning.

Then, the learned features H_k from each view are then fused via a learnable parameter. Finally, a two-layer fully connected neural network is applied to estimate the final rating for each item score related to clinical technique quality.

416 **Optimisation**

In the optimisation phase, our primary aim is to enhance the accuracy of our network's predictions. To this end, we utilise the Mean Absolute Error (MAE) as our metric of choice. MAE is a straightforward yet effective measure that calculates the average of the absolute differences between the predicted values by our network and the rating values ('ground truth') after clinical experts' agreement. This metric is particularly useful for our item scores, as it clearly indicates how close our predictions are to reality, on average, without being influenced by the direction of errors. The formula for our loss function, which incorporates MAE, is given by:

$$L = \sum_{q=1}^Q |\hat{y}_q - y_q|$$

where Q denotes the different item scores related to clinical technique quality, \hat{y}_q and y_q denote the predicted score and the ground truth score for each item, respectively. Our goal during training is to minimise this loss, which means reducing the average absolute error between our predictions and the expert ratings, thereby aligning our network's assessments more closely with the expert evaluations.

For the optimisation process, we employ the Adam optimiser, a widely used optimisation algorithm known for its effectiveness in handling sparse gradients and automatically adjusting the learning rate. This choice promotes

430 a more efficient and robust training process, with an initial learning rate set at 0.01 and a weight decay of 0.1, to
431 gradually improve our model's performance by iteratively adjusting its parameters in a direction that minimises the
432 loss function.

433 **Evaluation**

434 We evaluate our automatic AQA framework's performance using Mean Absolute Error (MAE), employing
435 a fivefold cross-validation approach. In each validation cycle, 80% of the data was used as the training set, while the
436 remaining 20% was used as the test set, ensuring different train and test sets for each iteration. The training epoch is
437 set to 100. Our evaluation not only compares our automated method's predictions to the final expert-agreed scores
438 but also examines the alignment of initial individual expert annotations with these consensus ratings.

439 The comparison involves calculating the MAE between our method's predictions and the expert-agreed
440 scores and, similarly, between individual expert scores before consensus and the final agreed scores. This approach
441 highlights our method's potential accuracy in relation to initial expert assessments. MAE was selected as the primary
442 metric due to its interpretability, robustness, and consistency. It is a common metric for skill assessment in the biomedical
443 engineering domain (Anastasiou et al., 2023; Wagner et al., 2023) It allows us to directly quantify the average error in our
444 model's predictions compared to expert ratings. Thus, we use MAE across both optimisation and evaluation phases,
445 facilitating a clearer comparison of our model's accuracy relative to expert assessments.

446 For a fair comparison, we align our automated assessments with expert evaluations by focusing on data from
447 cameras 1, 4, and 5, which the experts predominantly used. This strategy ensures that our method is evaluated from
448 the most relevant perspectives for accurate clinical technique assessment. Cycle 1 was included in the analysis to assess
449 the performance of our system in scenarios where human raters have difficulty reaching consensus. The low agreement
450 among experts in Cycle 1 highlights the complexity of certain CPR assessments and underscores the importance of having
451 an automated system that can provide consistent evaluations. By including Cycle 1, we ensure that our system is tested not
452 only on straightforward cases where human agreement is high but also on more challenging cases where human agreement
453 is low.

454 We implemented our method with PyTorch 1.10.1 and trained the models using one Nvidia GeForce GTX
455 2080 Ti GPU. For further reproduction and implementation of our research, our code and step-by-step deployment
456 instructions can be found on our GitHub page: <https://github.com/FrancisXZhang/CPR>.

457 **Results and Discussion of AQA**

458 In our analysis, Table 1 illustrates the MAE comparisons between the automatic AQA framework and the
 459 initial scores given by two human evaluators. Specifically, the MAE values represent the average difference between
 460 the scores assigned by our AQA system or each evaluator and the final agreed-upon ground truth scores established
 461 through expert consensus. The full score for each evaluated item is 4. In most cases, the error margin of our framework
 462 remains below 1, underscoring our automated methods' potential accuracy and applicability. When comparing our
 463 method with manual assessments, we found that our automatic AQA consistently exhibits significantly lower error in
 464 evaluating hand, arm, and shoulder positions. This may be attributed to our framework's reliance on precise pose
 465 information, providing a more objective assessment of the participant's posture. Our AQA exhibited higher error rates
 466 in the compression depth and compression rate items. This discrepancy could be because these two items require
 467 assessing interactions between the participant and the dummy (Kılıç et al., 2018), something our framework currently
 468 does not capture. This demonstration focused solely on the pose information of the participant and did not incorporate
 469 visual interaction data between humans and objects; further work could establish the importance of considering such
 470 interactions. It is also noteworthy that the expert raters used a metronome to assist in their rate judgements, which may
 471 account for a higher than typical expert-assessed accuracy rate in this dimension.

472 Table 1 Mean Average error (Human Experts vs. AQA framework)

Item	Evaluator 1	Evaluator 2	AQA
Hand Position	1.62	1.08	0.33
Arm Position	0.70	0.15	0.07
Shoulder Position	0.40	0.34	0.13
Depth	0.49	0.30	0.69
Rate	0.89	0.11	1.67
Compression Release	1.04	0.98	1.00
Total	3.96	2.69	2.98

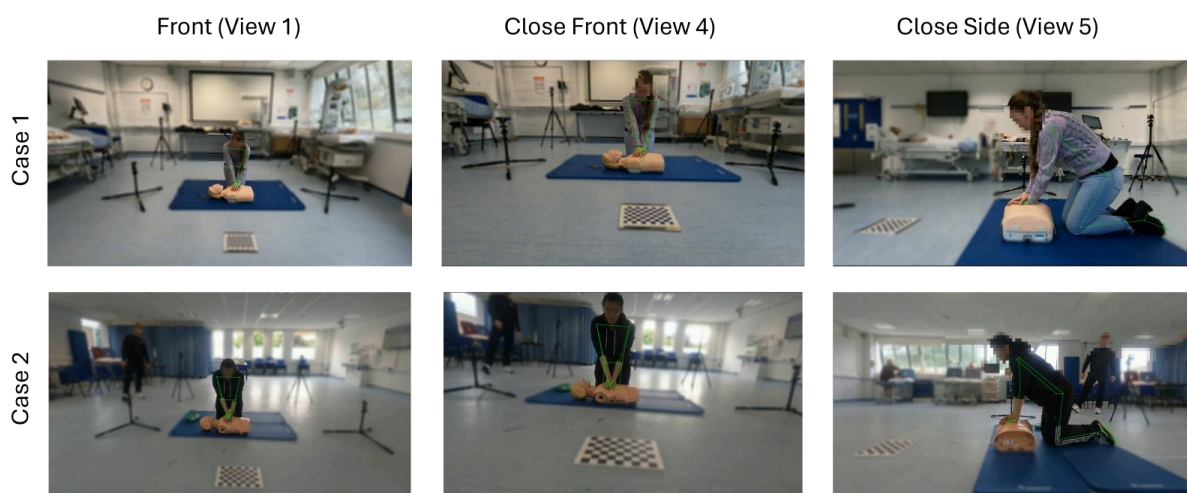
473

474 We acknowledge that the size of our dataset, comprising 53 samples, might appear limited for training deep learning
 475 models. However, we employed a fivefold cross-validation approach, promoting robust evaluation by using 80% of
 476 the data for training and 20% for testing in each fold, which helps in assessing the model's performance
 477 comprehensively and mitigating overfitting. Our dataset size is comparable to those used in other research within the
 478 domain of healthcare training systems, such as the 10 cases used in Liao et al. (2020) demonstrating the feasibility of
 479 using similar dataset sizes. Moreover, the results indicate that the tool performs well as compared to human raters,
 480 confirming the reliability of the measure.

481

482

483



484

485 *Figure 4. Qualitative exemplars for comparing our model with human raters. Our model's performance*
 486 *aligns with the final rating after agreement, whether in the human raters' initial agreed score or their conflicted*
 487 *score.*

488 To further compare our model with the human raters, qualitative research for individuals by our model
 489 compared to the raters' evaluations is shown in Figure 4. To better demonstrate the efficiency of our work, we
 490 included both cases of conflicts and agreements in our exemplars. As the main potential advantage of our system is
 491 based on pose estimation, these exemplars are mainly focused on pose-related scores.

492 In Figure 4, Case 1, both human raters gave high scores for the subject's hand, arm, and shoulder positions
493 (score from human raters = 4, maximum rating = 4). Similarly, our model also gave a score of 4 for these pose-
494 related aspects. In Figure 4, Case 2, although the reviewers gave a score of 4 for the arm and shoulder pose of the
495 participant, one rater gave a 0 for the hand position while the other gave a 4 before they reached an agreement. The
496 reason for the 0 was 'Too far toward the feet', but the rater who gave a 4 thought it was 'Just in the center of the
497 chest'. After double-checking the video, they ultimately agreed to give the subject a score of 4, as the hand pose was,
498 indeed, 'In the center of the chest (within one average hand margin)'.

499 One major reason for this conflict is that the raters typically make their judgments based on a single camera
500 view. However, the information from one camera can be limited due to different participants' angles towards the
501 camera and varying initial postures, which sometimes affects their judgment. For instance, in Case 2, the ratings are
502 highly focused on the close front view. The participant's habit of performing CPR vertically makes their hand
503 compression position appear farther from their body for the human raters. Our method, which fuses information
504 from multiple camera views, overcomes this limitation. Even though our pose estimation may occasionally show
505 misdetections (pose estimation visualisation for the full video can be found in the Supplementary Material),
506 summarising the motion information from multiple views makes our rating results more robust for pose-related
507 scores.

508 The present work sought to illustrate how accumulating skills performance datasets in real or simulated
509 settings could provide a foundation for understanding healthcare skills and building educational technology to
510 support healthcare professionals and educators. A fundamentally interdisciplinary approach (in this case, with a
511 strong emphasis on computer science) in the study of technical skills may assist in ensuring healthcare professional
512 competence. Specifically, we have shown how innovations in computer vision can be leveraged to provide (1) Data
513 from real settings that can be scientifically assessed (e.g. the spatial location of each joint in cartesian coordinate
514 space alongside estimation confidence) and (2) an assessment of performance quality based on both objective
515 measurements and learned parameters from expert raters. This assessment technique demonstrated comparable
516 accuracy in overall assessment relative to our expert raters, indicating the validity of the approach. Furthermore, the
517 margin of error from this assessment technique was typically below 1, indicating that accuracy between defined
518 thresholds of competency was excellent. However, specific to performance features, our assessment technique

519 sometimes outperformed or underperformed relative to the raters. Our intention with this work is to call for
520 foundations to be put in place that allow for the collection and sharing of diverse healthcare technical skills
521 performances that provide the opportunity for interdisciplinary collaborations to enhance the efficacy and efficiency
522 of healthcare professional education.

523 Manikins and simulators that provide automated feedback on a range of critical measures to facilitate the
524 acquisition of an appropriate skill threshold exist (e.g. QCPR manikins (Laerdal)). However, they are often limited in
525 the type of feedback they can provide, highly specific to a given skill or class of skills, or expensive. The computer
526 vision approach demonstrated in the present work has the potential to provide automated and targeted feedback for a
527 range of skills that can be assessed visually with low-cost video cameras and computers that would already be present
528 within an educational setting. Furthermore, the fact that a range of skills could be assessed with the relatively low-
529 cost set-up and without the need for specialist simulators represents an economic benefit. The flexibility of the
530 computer vision approach also makes it ideal for assessing complex skills performance in highly realistic simulations.
531 Indeed, self-training programmes exist using items that can be found around the home for a makeshift manikin
532 (Wanner et al., 2016), it is possible that this technology could be implemented using a webcam to provide feedback
533 to the trainee at home.

534 In future research, we plan to consider human-object interaction in our model. Incorporating interactions
535 with other humans or objects into our model necessitates a comprehensive approach. First, we would expand our
536 dataset to include scenarios involving human-human and human-object interactions, ensuring a wide range of contexts
537 and activities. Second, more detailed annotations would be required to label these interactions accurately, such as the
538 actual physical contact between humans and objects to make training closer to real-world conditions (Zhou et al.,
539 2023). Third, our model architecture would need modifications to handle the additional complexity, such as integrating
540 temporal-based pose estimation for more consistent motion information capture (Zhou et al., 2023). By addressing
541 these steps, we aim to significantly enhance the model's utility in more realistic and dynamic environments, ultimately
542 improving its applicability for various educational and training purposes.

543

544

545 Credit Statement

546 Conceptualisation: MDC, LG, HPHS, DP, AP. Data Curation: MDC, AC, FXZ, DM. Formal Analysis:
547 FXZ. Investigation: MDC, AC, FXZ, JR, DM, CF, LP. Methodology: MDC, LG, HPHS, DP, AP, AC, FXZ, DM.
548 Project Administration: MDC, JR. Resources: MDC, HPHS, DP, AP, AC. Software: FXZ. Supervision: MDC,
549 HPHS, JR. Visualisation: FXZ, JR. Writing – Original Draft. MDC, FXZ. Writing – Review & Editing. MDC, LG,
550 HPHS, DP, AP, CF, AC, FXZ, JR, DM, LP.

551

552 Acknowledgements

553 The authors would like to thank Stuart Barker, Mollie O'Donnell, and Oliver Barnett for their assistance in data
554 collection and recruitment of participants. The authors would also like to thank the technicians at the Department of
555 Nursing and Midwifery, Northumbria University for their assistance setting up the recording space.

556 Competing Interests

557 The authors declare no competing interests.

558 Funding

559 This paper was funded by Seed Funding from Northumbria University awarded to MDC, LG, HPHS, AP and DP.

560 Data Availability Statement

561 Participants consented to how they would like their data to be used and which data could be made available. Video
562 and Evaluation data are provided for those participants who consented to the sharing of their data, and for further
563 research use. The data has been anonymised such that faces have been blurred from the data set. This data is
564 protected such that an application must be made and the applicant must indicate that they will be using the data for
565 research or educational purposes, and preserve the anonymisation. The analyses provided are on the basis of the full
566 data set collected. The authors are able to provide the full data set for verification purposes only and this data set
567 should not be shared due to ethical constraints and to adhere to the participants' wishes.

568

569

References

- 570
571 Anastasiou, D., Jin, Y., Stoyanov, D., & Mazomenos, E. (2023). Keep Your Eye on the Best:
572 Contrastive Regression Transformer for Skill Assessment in Robotic Surgery. *IEEE*
573 *Robotics and Automation Letters*, 8(3), 1755–1762.
574 <https://doi.org/10.1109/LRA.2023.3242466>
- 575 Baldi, E., Cornara, S., Contri, E., Epis, F., Fina, D., Zelaschi, B., Dossena, C., Fichtner, F.,
576 Tonani, M., Maggio, M. D., Zambaiti, E., & Somaschini, A. (2017). Real-time visual
577 feedback during training improves laypersons' CPR quality: A randomised controlled
578 manikin study. *Canadian Journal of Emergency Medicine*, 19(6), 480–487.
579 <https://doi.org/10.1017/cem.2016.410>
- 580 Berg, K. M., Bray, J. E., Ng, K.-C., Liley, H. G., Greif, R., Carlson, J. N., Morley, P. T.,
581 Drennan, I. R., Smyth, M., Scholefield, B. R., Weiner, G. M., Cheng, A., Djärv, T.,
582 Abelairas-Gómez, C., Acworth, J., Andersen, L. W., Atkins, D. L., Berry, D. C., Bhanji,
583 F., ... Yamada, N. K. (2023). 2023 International Consensus on Cardiopulmonary
584 Resuscitation and Emergency Cardiovascular Care Science With Treatment
585 Recommendations: Summary From the Basic Life Support; Advanced Life Support;
586 Pediatric Life Support; Neonatal Life Support; Education, Implementation, and Teams;
587 and First Aid Task Forces. *Circulation*, 148(24), e187–e280.
588 <https://doi.org/10.1161/CIR.0000000000001179>
- 589 Bouget, D., Allan, M., Stoyanov, D., & Jannin, P. (2017). Vision-based and marker-less surgical
590 tool detection and tracking: A review of the literature. *Medical Image Analysis*, 35, 633–
591 654. <https://doi.org/10.1016/j.media.2016.09.003>

- 592 Castillo, J., Gomar, C., Rodriguez, E., Trapero, M., & Gallart, A. (2019). Cost minimisation
593 analysis for basic life support. *Resuscitation*, *134*, 127–132.
594 <https://doi.org/10.1016/j.resuscitation.2018.11.008>
- 595 Constable, M. D., Shum, H. P. H., & Clark, S. (2024). Enhancing surgical performance in
596 cardiothoracic surgery with innovations from computer vision and artificial intelligence:
597 A narrative review. *Journal of Cardiothoracic Surgery*, *19*(1), 94.
598 <https://doi.org/10.1186/s13019-024-02558-5>
- 599 Constable, M. D., Zhang, F. X., Connor, T., Monk, D., Rajsic, J., Ford, C., Park, L. J., Barker, S.,
600 Platt, A., Porteous, D., Grierson, L., & Shum, H. P. H. (2024). *Cardiopulmonary*
601 *Resuscitation Performance: Video, Demographic and Evaluation Data, 2023* [Data
602 Collection]. UK Data Service. <https://doi.org/10.5255/UKDA-SN-857038>
- 603 Corbett-Davies, S., Gaebler, J. D., Nilforoshan, H., Shroff, R., & Goel, S. (2023). *The Measure*
604 *and Mismeasure of Fairness* (arXiv:1808.00023). arXiv.
605 <https://doi.org/10.48550/arXiv.1808.00023>
- 606 El-Sayed, C., Yiu, A., Burke, J., Vaughan-Shaw, P., Todd, J., Lin, P., Kasmani, Z., Munsch, C.,
607 Rooshenas, L., Campbell, M., & Bach, S. P. (2024). Measures of performance and
608 proficiency in robotic assisted surgery: A systematic review. *Journal of Robotic Surgery*,
609 *18*(1), 16. <https://doi.org/10.1007/s11701-023-01756-y>
- 610 Enarsson, T., Enqvist, L., & Naarttijärvi, M. (2022). Approaching the human in the loop – legal
611 perspectives on hybrid human/algorithmic decision-making in three contexts. *Information*
612 *& Communications Technology Law*, *31*(1), 123–153.
613 <https://doi.org/10.1080/13600834.2021.1958860>

- 614 Ericsson, K. A. (2004). Deliberate practice and the acquisition and maintenance of expert
615 performance in medicine and related domains. *Academic Medicine: Journal of the*
616 *Association of American Medical Colleges*, 79(10 Suppl), S70-81.
617 <https://doi.org/10.1097/00001888-200410001-00022>
- 618 Feng, L., Zhao, Y., Zhao, W., & Tang, J. (2022). A comparative review of graph convolutional
619 networks for human skeleton-based action recognition. *Artificial Intelligence Review*, 1–
620 31.
- 621 Frank, J. R., Snell, L. S., Cate, O. T., Holmboe, E. S., Carraccio, C., Swing, S. R., Harris, P.,
622 Glasgow, N. J., Campbell, C., Dath, D., Harden, R. M., Iobst, W., Long, D. M.,
623 Mungroo, R., Richardson, D. L., Sherbino, J., Silver, I., Taber, S., Talbot, M., & Harris,
624 K. A. (2010). Competency-based medical education: Theory to practice. *Medical*
625 *Teacher*, 32(8), 638–645. <https://doi.org/10.3109/0142159X.2010.501190>
- 626 Gates, S., Quinn, T., Deakin, C. D., Blair, L., Couper, K., & Perkins, G. D. (2015). Mechanical
627 chest compression for out of hospital cardiac arrest: Systematic review and meta-analysis.
628 *Resuscitation*, 94, 91–97. <https://doi.org/10.1016/j.resuscitation.2015.07.002>
- 629 Giblin, G., Tor, E., & Parrington, L. (2016). The impact of technology on elite sports
630 performance. *Sensoria: A Journal of Mind, Brain & Culture*, 12.
631 <https://doi.org/10.7790/sa.v12i2.436>
- 632 Glazier, P. S. (2021). Beyond animated skeletons: How can biomechanical feedback be used to
633 enhance sports performance? *Journal of Biomechanics*, 129, 110686.
634 <https://doi.org/10.1016/j.jbiomech.2021.110686>

- 635 Goldsteen, A., Ezov, G., Shmelkin, R., Moffie, M., & Farkash, A. (2022). Data minimisation for
636 GDPR compliance in machine learning models. *AI and Ethics*, 2(3), 477–491.
637 <https://doi.org/10.1007/s43681-021-00095-8>
- 638 Gordon, L., Reed, C., Sorensen, J. L., Schulthess, P., Strandbygaard, J., Mcloone, M.,
639 Grantcharov, T., & Shore, E. M. (2022). Perceptions of safety culture and recording in
640 the operating room: Understanding barriers to video data capture. *Surgical Endoscopy*,
641 36(6), 3789–3797. <https://doi.org/10.1007/s00464-021-08695-5>
- 642 Hagendorff, T. (2019). From privacy to anti-discrimination in times of machine learning. *Ethics*
643 *and Information Technology*, 21(4), 331–343. [https://doi.org/10.1007/s10676-019-09510-](https://doi.org/10.1007/s10676-019-09510-5)
644 5
- 645 Hajian, S., & Domingo-Ferrer, J. (2013). Direct and indirect discrimination prevention methods.
646 In *Discrimination and privacy in the information society: Data mining and profiling in*
647 *large databases* (pp. 241–254). Springer.
- 648 Harden, R. M. (2007). Outcome-based education—The ostrich, the peacock and the beaver.
649 *Medical Teacher*, 29(7), 666–671. <https://doi.org/10.1080/01421590701729948>
- 650 Holstein, K., Wortman Vaughan, J., Daumé, H., Dudik, M., & Wallach, H. (2019). Improving
651 Fairness in Machine Learning Systems: What Do Industry Practitioners Need?
652 *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, 1–
653 16. <https://doi.org/10.1145/3290605.3300830>
- 654 Huang, Y., Shum, H. P. H., Ho, E. S. L., & Aslam, N. (2020). High-speed multi-person pose
655 estimation with deep feature transfer. *Computer Vision and Image Understanding*, 197–
656 198, 103010. <https://doi.org/10.1016/j.cviu.2020.103010>

- 657 Ionescu, C., Papava, D., Olaru, V., & Sminchisescu, C. (2014). Human3.6M: Large Scale
658 Datasets and Predictive Methods for 3D Human Sensing in Natural Environments. *IEEE*
659 *Transactions on Pattern Analysis and Machine Intelligence*, 36(7), 1325–1339.
660 <https://doi.org/10.1109/TPAMI.2013.248>
- 661 Ismail Fawaz, H., Forestier, G., Weber, J., Idoumghar, L., & Muller, P.-A. (2018). Evaluating
662 Surgical Skills from Kinematic Data Using Convolutional Neural Networks. In A. F.
663 Frangi, J. A. Schnabel, C. Davatzikos, C. Alberola-López, & G. Fichtinger (Eds.),
664 *Medical Image Computing and Computer Assisted Intervention – MICCAI 2018* (pp.
665 214–221). Springer International Publishing. [https://doi.org/10.1007/978-3-030-00937-](https://doi.org/10.1007/978-3-030-00937-3_25)
666 [3_25](https://doi.org/10.1007/978-3-030-00937-3_25)
- 667 Judkins, T. N., Oleynikov, D., & Stergiou, N. (2008). Enhanced Robotic Surgical Training Using
668 Augmented Visual Feedback. *Surgical Innovation*, 15(1), 59–68.
669 <https://doi.org/10.1177/1553350608315953>
- 670 Kanazawa, A., Black, M. J., Jacobs, D. W., & Malik, J. (2018). End-to-End Recovery of Human
671 Shape and Pose. *2018 IEEE/CVF Conference on Computer Vision and Pattern*
672 *Recognition*, 7122–7131. <https://doi.org/10.1109/CVPR.2018.00744>
- 673 Kazim, E., Denny, D. M. T., & Koshiyama, A. (2021). AI auditing and impact assessment:
674 According to the UK information commissioner's office. *AI and Ethics*, 1(3), 301–310.
675 <https://doi.org/10.1007/s43681-021-00039-2>
- 676 Kendrick, D. E., Thelen, A. E., Chen, X., Gupta, T., Yamazaki, K., Krumm, A. E., Bandeh-
677 Ahmadi, H., Clark, M., Luckoscki, J., Fan, Z., Wnuk, G. M., Ryan, A. M., Mukherjee, B.,
678 Hamstra, S. J., Dimick, J. B., Holmboe, E. S., & George, B. C. (2023). Association of

- 679 Surgical Resident Competency Ratings With Patient Outcomes. *Academic Medicine:*
680 *Journal of the Association of American Medical Colleges*, 98(7), 813–820.
681 <https://doi.org/10.1097/ACM.00000000000005157>
- 682 Kılıç, D., Göksu, E., Kılıç, T., & Buyurgan, C. (2018). Resuscitation quality of rotating chest
683 compression providers at one-minute vs. Two-minute intervals: A mannequin study. *The*
684 *American Journal of Emergency Medicine*, 36(5), 829–833.
- 685 Kocabas, M., Karagoz, S., & Akbas, E. (2019). Self-Supervised Learning of 3D Human Pose
686 Using Multi-View Geometry. *2019 IEEE/CVF Conference on Computer Vision and*
687 *Pattern Recognition (CVPR)*, 1077–1086. <https://doi.org/10.1109/CVPR.2019.00117>
- 688 Lagomarsino, M., Lorenzini, M., Balatti, P., Momi, E. D., & Ajoudani, A. (2022). Pick the Right
689 Co-Worker: Online Assessment of Cognitive Ergonomics in Human-Robot Collaborative
690 Assembly. *IEEE Transactions on Cognitive and Developmental Systems*, 1–1.
691 <https://doi.org/10.1109/TCDS.2022.3182811>
- 692 Lam, K., Chen, J., Wang, Z., Iqbal, F. M., Darzi, A., Lo, B., Purkayastha, S., & Kinross, J. M.
693 (2022). Machine learning for technical skill assessment in surgery: A systematic review.
694 *Npj Digital Medicine*, 5(1), Article 1. <https://doi.org/10.1038/s41746-022-00566-0>
- 695 Liao, Y., Vakanski, A., & Xian, M. (2020). A Deep Learning Framework for Assessing Physical
696 Rehabilitation Exercises. *IEEE Transactions on Neural Systems and Rehabilitation*
697 *Engineering*, 28(2), 468–477. <https://doi.org/10.1109/TNSRE.2020.2966249>
- 698 Likosky, D., Yule, S. J., Mathis, M. R., Dias, R. D., Corso, J. J., Zhang, M., Krein, S. L.,
699 Caldwell, M. D., Louis, N., Janda, A. M., Shah, N. J., Pagani, F. D., Stakich-Alpirez, K.,
700 & Manojlovich, M. M. (2021). Novel Assessments of Technical and Nontechnical

- 701 Cardiac Surgery Quality: Protocol for a Mixed Methods Study. *JMIR Research*
702 *Protocols*, 10(1), e22536. <https://doi.org/10.2196/22536>
- 703 Lin, T.-Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., & Zitnick, C.
704 L. (2014). Microsoft COCO: Common Objects in Context. In D. Fleet, T. Pajdla, B.
705 Schiele, & T. Tuytelaars (Eds.), *Computer Vision – ECCV 2014* (Vol. 8693, pp. 740–
706 755). Springer International Publishing. https://doi.org/10.1007/978-3-319-10602-1_48
- 707 Lins, C., Eckhoff, D., Klausen, A., Hellmers, S., Hein, A., & Fudickar, S. (2019).
708 Cardiopulmonary resuscitation quality parameters from motion capture data using
709 Differential Evolution fitting of sinusoids. *Applied Soft Computing*, 79, 300–309.
710 <https://doi.org/10.1016/j.asoc.2019.03.023>
- 711 Lugaresi, C., Tang, J., Nash, H., McClanahan, C., Uboweja, E., Hays, M., Zhang, F., Chang, C.-
712 L., Yong, M. G., Lee, J., Chang, W.-T., Hua, W., Georg, M., & Grundmann, M. (2019).
713 *MediaPipe: A Framework for Building Perception Pipelines* (arXiv:1906.08172). arXiv.
714 <https://doi.org/10.48550/arXiv.1906.08172>
- 715 Merchant, R. M., Topjian, A. A., Panchal, A. R., Cheng, A., Aziz, K., Berg, K. M., Lavonas, E.
716 J., Magid, D. J., & null, null. (2020). Part 1: Executive Summary: 2020 American Heart
717 Association Guidelines for Cardiopulmonary Resuscitation and Emergency
718 Cardiovascular Care. *Circulation*, 142(16_suppl_2), S337–S357.
719 <https://doi.org/10.1161/CIR.0000000000000918>
- 720 Oermann, M. H., Kardong-Edgren, S. E., & Odom-Maryon, T. (2011). Effects of monthly
721 practice on nursing students' CPR psychomotor skill performance. *Resuscitation*, 82(4),
722 447–453. <https://doi.org/10.1016/j.resuscitation.2010.11.022>

- 723 Olasveengen, T. M., Semeraro, F., Ristagno, G., Castren, M., Handley, A., Kuzovlev, A.,
724 Monsieurs, K. G., Raffay, V., Smyth, M., Soar, J., Svavarsdottir, H., & Perkins, G. D.
725 (2021). European Resuscitation Council Guidelines 2021: Basic Life Support.
726 *Resuscitation*, *161*, 98–114. <https://doi.org/10.1016/j.resuscitation.2021.02.009>
- 727 O'Meara, P., Munro, G., Williams, B., Cooper, S., Bogossian, F., Ross, L., Sparkes, L.,
728 Browning, M., & McClounan, M. (2015). Developing situation awareness amongst
729 nursing and paramedicine students utilising eye tracking technology and video debriefing
730 techniques: A proof of concept paper. *International Emergency Nursing*, *23*(2), 94–99.
731 <https://doi.org/10.1016/j.ienj.2014.11.001>
- 732 Platt, A., McMeekin, P., & Prescott-Clements, L. (2021). Effects of the Simulation Using Team
733 Deliberate Practice (Sim-TDP) model on the performance of undergraduate nursing
734 students. *BMJ Simulation and Technology Enhanced Learning*, *7*(2), 66–74. Scopus.
735 <https://doi.org/10.1136/bmjstel-2019-000520>
- 736 Resuscitation Council UK. (2021). *Adult basic life support Guidelines*. Resuscitation Council
737 UK. [https://www.resus.org.uk/library/2021-resuscitation-guidelines/adult-basic-life-](https://www.resus.org.uk/library/2021-resuscitation-guidelines/adult-basic-life-support-guidelines)
738 [support-guidelines](https://www.resus.org.uk/library/2021-resuscitation-guidelines/adult-basic-life-support-guidelines)
- 739 Semeraro, F., Frisoli, A., Loconsole, C., Bannò, F., Tammaro, G., Imbriaco, G., Marchetti, L., &
740 Cerchiari, E. L. (2013). Motion detection technology as a tool for cardiopulmonary
741 resuscitation (CPR) quality training: A randomised crossover mannequin pilot study.
742 *Resuscitation*, *84*(4), 501–507. <https://doi.org/10.1016/j.resuscitation.2012.12.006>

- 743 Srivastav, V., Issenhuth, T., Kadkhodamohammadi, A., de Mathelin, M., Gangi, A., & Padoy, N.
744 (2018, August 24). *MVOR: A Multi-view RGB-D Operating Room Dataset for 2D and*
745 *3D Human Pose Estimation*. arXiv.Org. <https://arxiv.org/abs/1808.08180v3>
- 746 Strandbygaard, J., Dose, N., Moeller, K. E., Gordon, L., Shore, E., Rosthøj, S., Ottesen, B.,
747 Grantcharov, T., & Sorensen, J. L. (2022). Healthcare professionals' perception of safety
748 culture and the Operating Room (OR) Black Box technology before clinical
749 implementation: A cross-sectional survey. *BMJ Open Quality*, *11*(4), e001819.
750 <https://doi.org/10.1136/bmjopen-2022-001819>
- 751 Taylor, J. E. T., & Taylor, G. W. (2020). Artificial cognition: How experimental psychology can
752 help generate explainable artificial intelligence. *Psychonomic Bulletin & Review*.
753 <https://doi.org/10.3758/s13423-020-01825-5>
- 754 van Dalen, A. S. H. M., Legemaate, J., Schlack, W. S., Legemate, D. A., & Schijven, M. P.
755 (2019). Legal perspectives on black box recording devices in the operating environment.
756 *BJS (British Journal of Surgery)*, *106*(11), 1433–1441. <https://doi.org/10.1002/bjs.11198>
- 757 Veale, M., & Binns, R. (2017). Fairer machine learning in the real world: Mitigating
758 discrimination without collecting sensitive data. *Big Data & Society*, *4*(2),
759 2053951717743530. <https://doi.org/10.1177/2053951717743530>
- 760 Wagner, M., Müller-Stich, B.-P., Kisilenko, A., Tran, D., Heger, P., Mündermann, L., Lubotsky,
761 D. M., Müller, B., Davitashvili, T., Capek, M., Reinke, A., Reid, C., Yu, T., Vardazaryan,
762 A., Nwoye, C. I., Padoy, N., Liu, X., Lee, E.-J., Disch, C., ... Bodenstedt, S. (2023).
763 Comparative validation of machine learning algorithms for surgical workflow and skill

- 764 analysis with the HeiChole benchmark. *Medical Image Analysis*, 86, 102770.
765 <https://doi.org/10.1016/j.media.2023.102770>
- 766 Walshe, N. C., Crowley, C. M., O'Brien, S., Browne, J. P., & Hegarty, J. M. (2019). Educational
767 Interventions to Enhance Situation Awareness: A Systematic Review and Meta-Analysis.
768 *Simulation in Healthcare*, 14(6), 398. <https://doi.org/10.1097/SIH.0000000000000376>
- 769 Wang, J.-C., Tsai, S.-H., Chen, Y.-H., Chen, Y.-L., Chu, S.-J., & Liao, W.-I. (2018). Kinect-
770 based real-time audiovisual feedback device improves CPR quality of lower-body-weight
771 rescuers. *The American Journal of Emergency Medicine*, 36(4), 577–582.
772 <https://doi.org/10.1016/j.ajem.2017.09.022>
- 773 Wanner, G. K., Osborne, A., & Greene, C. H. (2016). Brief compression-only cardiopulmonary
774 resuscitation training video and simulation with homemade mannequin improves CPR
775 skills. *BMC Emergency Medicine*, 16(1), 45. <https://doi.org/10.1186/s12873-016-0110-5>
- 776 Weiss, K. E., Kolbe, M., Nef, A., Grande, B., Kalirajan, B., Meboldt, M., & Lohmeyer, Q.
777 (2023). Data-driven resuscitation training using pose estimation. *Advances in Simulation*,
778 8(1), 12. <https://doi.org/10.1186/s41077-023-00251-6>
- 779 World Medical Association. (2013). World Medical Association Declaration of Helsinki: Ethical
780 principles for medical research involving human subjects. *JAMA*, 310(20), 2191–2194.
781 <https://doi.org/10.1001/jama.2013.281053>
- 782 Xie, H., Luo, H., Lin, J., & Yang, N. (2020). A novel algorithm of fast CPR quality evaluation
783 based on kinect. *Journal of Algorithms & Computational Technology*, 14,
784 1748302620983661. <https://doi.org/10.1177/1748302620983661>

- 785 Yan, S., Xiong, Y., & Lin, D. (2018). Spatial temporal graph convolutional networks for
786 skeleton-based action recognition. *Proceedings of the AAAI Conference on Artificial*
787 *Intelligence*, 32(1).
- 788 Zhou, K., Chen, C., Ma, Y., Leng, Z., Shum, H. P. H., Li, F. W. B., & Liang, X. (2023). *A Mixed*
789 *Reality Training System for Hand-Object Interaction in Simulated Microgravity*
790 *Environments*. 167–176. <https://doi.org/10.1109/ISMAR59233.2023.00031>
791