



1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25

Thesis Title: Accuracy and Reliability of Remote
Shoulder Motion Capturing Methods: A
Systematic Review and Meta-Analysis

Volume: 1 of 1

Name: Pengchi Chen

Qualification thesis submission for: MSc by Thesis

The University of Hull & the University of York

Hull York Medical School

August 2025

26 **ABSTRACT**

27 **Background:**

28 The COVID-19 pandemic accelerated the demand for remote assessment tools in
29 rehabilitation, spotlighting the need for accurate and reliable technologies to
30 measure shoulder range of motion (ROM) outside of clinical environments. Emerging
31 tools such as smartphone apps, wearable sensors, and markerless motion capture
32 systems are increasingly being adopted, yet their accuracy and reliability compared
33 to reference standards remains unclear.

34

35 **Objective:**

36 To systematically evaluate the accuracy and reliability of existing remote shoulder
37 ROM measurement technologies, quantify measurement bias, and assess their
38 agreement with reference standards.

39

40 **Methods:**

41 A systematic review and meta-analysis was conducted on 26 studies evaluating
42 remote ROM measurement tools. Pooled mean bias (in degrees) was calculated as
43 the primary effect size for agreement, with reliability assessed using intraclass
44 correlation coefficients (ICCs). Subgroup analyses were performed by motion type,
45 technology category, population health status, and data acquisition method. Risk of
46 bias was assessed using the QUADAS-2 tool.

47

48 **Results:**

49 Remote measurement methods showed a small but consistent overestimation of
50 ROM compared to reference standards (pooled mean bias = 2.63°, 95% CI: 1.52°,
51 3.74°), particularly in flexion, internal rotation, and external rotation. No significant
52 bias was observed in abduction or extension. Both inertial measurement unit (IMU)
53 and non-IMU technologies demonstrated comparable levels of overestimation.
54 Pathological populations exhibited greater variability (bias = 4.33° vs. 2.37° in
55 healthy subjects). Self-measurements showed lower and non-significant bias
56 compared to assessor-guided methods. Reliability was generally high, especially for
57 test-retest assessments (ICCs > 0.90), though more variable in inter-rater and
58 pathological settings.

59

60 **Conclusion:**

61 Remote technologies for assessing shoulder range of motion show generally small
62 differences compared to reference standards, but performance is inconsistent, with
63 substantial heterogeneity and variable reliability across studies. Although average
64 error is often below reported minimal clinically important differences (MCID), the
65 wide variability observed suggests these tools may not reliably detect small but
66 clinically meaningful changes. As such, they may have value for tracking overall
67 trends in shoulder function, but their use in clinical decision-making—particularly in
68 pathological populations—should be approached with caution.

69 **Table of Contents**

70 **ABSTRACT** - 2 -

71 **TABLE OF CONTENTS**..... - 4 -

72 **LIST OF TABLES AND FIGURES** - 5 -

73 **ACKNOWLEDGEMENT & DECLARATION**..... - 6 -

74 **INTRODUCTION** - 8 -

75 **OBJECTIVES**..... - 11 -

76 **METHODS**..... - 13 -

77 **PROTOCOL AND REGISTRATION** - 13 -

78 **ELIGIBILITY CRITERIA** - 14 -

79 **DATABASE AND SEARCH STRATEGY** - 16 -

80 **STUDY SELECTION** - 16 -

81 **DATA EXTRACTION AND RISK OF BIAS ASSESSMENT** - 17 -

82 **STATISTICAL ANALYSIS** - 17 -

83 **ETHICS** - 19 -

84 **RESULTS** - 20 -

85 **STUDY SELECTION AND CHARACTERISTICS** - 20 -

86 **ACCURACY OF REMOTE SHOULDER ROM MEASUREMENT** 33

87 **PUBLICATION BIAS**..... 42

88 **RELIABILITY OF REMOTE SHOULDER ROM MEASUREMENT** 43

89 **RISK OF BIAS ASSESSMENT** 46

90 **DISCUSSIONS** 51

91 **CONCLUSION:** 57

92 **REFERENCES** 58

93 **APPENDIX 1: FULL DATABASE SEARCH**..... 62

94 **APPENDIX 2: SIGNALLING QUESTIONS USED FOR THE QUADAS-2 TOOL...** 64

95
96

97	List of Tables and Figures	
98	Tables	
99	Table 1. Overall study characteristics summary	21
100	Table 2. QUADAS-2 Risk of Bias Assessment	47
101		
102	Figures	
103	Figure 1. PRISMA Diagram for Study Selection	14
104	Figure 2. Accuracy by Motion Types	35
105	Figure 3. Accuracy by Technology Types	37
106	Figure 4. Accuracy in Healthy vs. Pathological Shoulders	39
107	Figure 5. Accuracy by Data Acquisition Methods	41
108	Figure 6. Funnel Plot of the Included Data	42
109	Figure 7. Visual Representation of Risk of Bias using QUADAS-2	46

110 **Acknowledgement & Declaration**

111 First and foremost, I would like to express my deepest gratitude to my supervisors,
112 Professor Amar Rangan and Dr. Peter Ellison, for their invaluable guidance and
113 expertise throughout the course of this project. Their mentorship has been
114 instrumental in shaping both this thesis and my development as a clinical academic.

115

116 I would also like to thank my colleague Miss Mia Prosser for her vital contributions as
117 second reviewer during the systematic review process. Her attention to detail and
118 thoughtful input were essential in strengthening the methodological rigour of the
119 work.

120

121 My sincere thanks to Professor Bob Phillips, who served as my TAP Chair, for his
122 ongoing support, critical feedback, and mentorship throughout my academic training.

123

124 I am also grateful to Dr. Emily Shoesmith, Dr. Elena Ratschen, Professor Mona
125 Kanaan, and Miss Ada Keding at the University of York for their teachings across
126 key modules including evidence synthesis, regression analysis, and measurement in
127 health and disease. Their instruction provided me with a solid methodological and
128 statistical foundation, not only for this thesis but for my future research endeavours.

129

130 I am grateful to the NIHR Academic Clinical Fellowship Programme, which provided
131 the opportunity and resources for me to undertake this research alongside my
132 clinical training.

133

134 Lastly, I would like to thank my family and friends for their unwavering support and
135 encouragement during the highs and lows of this academic endeavour.

136

137 I confirm that this work is original and that if any passage(s) or diagram(s) have been
138 copied from academic papers, books, the internet or any other sources these are
139 clearly identified by the use of quotation marks and the reference(s) is fully cited. I
140 certify that, other than where indicated, this is my own work and does not breach the
141 regulations of HYMS, the University of Hull or the University of York regarding
142 plagiarism or academic conduct in examinations. I have read the HYMS Code of

143 Practice on Academic Misconduct, and state that this piece of work is my own and
144 does not contain any unacknowledged work from any other sources.

145 **Introduction**

146 Shoulder range of motion (ROM) assessment is a fundamental component of
147 musculoskeletal and orthopaedic evaluations.⁴⁴ Accurate measurement of ROM is
148 essential for diagnosing shoulder conditions, monitoring treatment outcomes, and
149 guiding rehabilitation. Traditional shoulder ROM assessments are typically
150 performed in clinical settings using tools such as handheld goniometers or
151 inclinometers, requiring direct patient-clinician interaction.

152

153 A handheld goniometer is a mechanical device that measures joint angles by
154 aligning its arms with anatomical landmarks, typically requiring manual positioning by
155 a clinician. In contrast, inclinometers measure angular displacement relative to
156 gravity, allowing for more objective quantification of joint movement, particularly in
157 single-plane motions.²⁶ More recently, digital and remote measurement technologies
158 have been developed. Smartphone-based applications often utilise built-in
159 accelerometers and gyroscopes to function as digital inclinometers or employ
160 camera-based pose estimation algorithms to derive joint angles from images or
161 video.^{11, 35} Wearable inertial measurement units (IMUs) incorporate accelerometers,
162 gyroscopes, and sometimes magnetometers to track segment orientation and
163 movement in three-dimensional space.^{7, 38} Markerless motion capture systems use
164 computer vision techniques to estimate body joint positions without the need for
165 physical markers, typically through depth sensors or standard RGB cameras.^{7, 53}
166 These technologies differ in their underlying measurement principles, required user
167 input, and susceptibility to sources of error such as sensor drift, soft tissue artefact,
168 and camera positioning, which are important considerations when comparing their
169 accuracy and reliability to conventional clinical tools.

170

171 With the increasing interest in telerehabilitation and remote monitoring, there is a
172 growing body of research assessing the accuracy and reliability of these
173 technologies compared to traditional face-to-face ROM assessments. Accuracy
174 refers to the accuracy of these tools in measuring ROM compared to reference
175 standard, whereas reliability assesses their consistency across different conditions,
176 including test-retest reliability and inter-rater reliability.⁴⁹ Despite the desire to adopt

177 remote ROM measurement tools in routine clinical practice and research, a
178 comprehensive synthesis of the evidence on their accuracy and reliability is lacking.

179
180 The use of digital technologies such as smartphone goniometer apps and wearable
181 devices in ROM assessment has been widely explored, particularly before the
182 COVID-19 pandemic. Werner and colleagues (2014) validated a smartphone
183 clinometer application for measuring shoulder ROM.⁵⁴ While Mitchell and colleagues
184 (2014) and Johnson and colleagues (2015) assessed the reliability and validity of
185 smartphone-based goniometry in clinical settings.^{23, 28} Rigoni and colleagues (2019)
186 evaluated a wireless inertial motion capture device for ROM assessment.³⁷ Hayes
187 and colleagues (2015) examined traditional methods such as estimating joint angles
188 from patient photographs.¹⁸ These validation studies generally reported
189 measurement errors within approximately 5°–10° of reference standards and
190 intraclass correlation coefficients (ICCs) often exceeding 0.80, which are commonly
191 considered indicative of good reliability in musculoskeletal assessment.⁴⁹ In clinical
192 practice, acceptable measurement error for shoulder ROM is often interpreted in
193 relation to the minimum clinically important difference (MCID), which has been
194 reported to range between approximately 11° and 24° depending on the context and
195 method of assessment.²⁹

196
197 Therefore, while these early studies suggested that digital measurement tools could
198 achieve clinically acceptable levels of accuracy and reliability under controlled
199 conditions, they were largely conducted in face-to-face or supervised environments.
200 As such, the extent to which these performance benchmarks translate to remote or
201 unsupervised settings remains uncertain. Challenges persist in tele-rehabilitation
202 settings, including low sensitivities, poor inter-rater reliability, and variability in clinical
203 testing accuracy.^{5, 15} These limitations highlight the need for further research to
204 optimise assessment techniques.

205
206 The broader implications of remote rehabilitation have also been explored. Kane and
207 colleagues (2020) conducted a randomised controlled trial assessing post-operative
208 tele-follow-ups for rotator cuff repair patients, reporting comparable patient
209 satisfaction scores to in-person follow-ups; it did not, however, evaluate longitudinal
210 ROM outcomes, leaving uncertainties regarding the effectiveness of tele-

211 rehabilitation in monitoring recovery progress.²⁴ Similarly, Faber and colleagues
212 (2015) highlighted challenges in home-based rehabilitation, particularly in patients
213 with frozen shoulders, who often struggle to follow rehabilitation instructions. These
214 findings reinforce the need for structured and engaging tele-rehabilitation strategies
215 to enhance patient compliance.¹³ However, while poor adherence to rehabilitation
216 programmes is well recognised, the role of objective monitoring in addressing this
217 challenge remains less clearly defined. Reliable remote measurement of shoulder
218 range of motion (ROM) may provide clinicians with the ability to track patient
219 progress more accurately, identify suboptimal engagement or incorrect exercise
220 performance, and offer timely feedback or intervention. In this context, the availability
221 of accurate and reliable remote ROM assessment tools may support improved
222 patient engagement by enabling personalised feedback and accountability, thereby
223 complementing rehabilitation strategies rather than replacing them. This highlights
224 the importance of establishing whether such remote measurement technologies can
225 provide clinically acceptable levels of accuracy and reliability, which is the focus of
226 the present review.

227

228 Validated patient-reported outcome measures (PROMs) such as the Oxford
229 Shoulder Score (OSS), a 12-item joint-specific questionnaire, are widely utilised in
230 randomised controlled trials assessing treatment outcomes across populations with
231 various shoulder conditions.^{32, 36} While PROMs focus on pragmatic aspects of
232 patient outcomes, emphasising pain levels and the ability to perform daily activities
233 as indicators of improved shoulder function, there are concerns that they do not fully
234 reflect the range of motion achieved.²³ Following recommendation by the European
235 Society for Surgery of the Shoulder and Elbow, the Constant-Murley score has been
236 widely used, as it combines subjective patient-reported outcomes and objective
237 range of motion measurements, which contribute to 40% of the total score.⁹ While
238 range of motion is a primary outcome in shoulder assessment, it does not fully
239 capture functional movement quality. Muscle strength is an integral component of
240 shoulder function and is often assessed alongside range of motion in both clinical
241 and remote settings. Therefore, studies assessing strength were included where
242 they contributed to the evaluation of shoulder movement performance.

243

244

245 Digital app-based technologies have demonstrated potential in increasing patient
246 engagement following total joint arthroplasty, particularly for hip and knee
247 replacements; their application in shoulder rehabilitation, however, remains
248 unvalidated, representing a gap in the literature.²⁵ Advanced motion-detecting
249 technologies, such as depth-based skeleton tracking, have shown promise in remote
250 rehabilitation applications but remain impractical for routine clinical use due to
251 technological constraints.¹ Similarly, the accuracy of telehealth shoulder
252 examinations has been critically assessed, with findings suggesting the need for
253 improved methodologies and standardisation in remote assessments.⁵

254

255 In the context of this review, remote ROM measurement is defined as assessments
256 conducted without direct physical interaction between the patient and a healthcare
257 provider. This includes self-measurements performed by patients and clinician-
258 guided assessments conducted remotely via digital or telemedicine platforms. This
259 definition distinguishes the present review from previous work, such as systematic
260 reviews assessing the validity of digital ROM measurement devices in clinical
261 settings.⁴⁵ This review specifically investigates technologies that facilitate ROM
262 assessments remotely, with minimal or no physical contact between the patient and
263 clinician. This distinction is crucial as it frames the applicability of the findings toward
264 tele-rehabilitation, remote monitoring, and digital health integration rather than
265 merely in-clinic digital measurement advancements.

266

267 **Objectives**

268 The reserach objectives are summarised using a PICO-style framework:

269 Population (P): adults with healthy shoulders or musculoskeletal shoulder conditions
270 undergoing shoulder range of motion & strength assessment.

271 Intervention/Index test (I): remote shoulder ROM & strength measurement methods,
272 including smartphone applications, wearable sensors, and markerless motion
273 capture systems.

274 Comparator/Reference standard (C): conventional or reference assessment
275 methods, including goniometers, inclinometers, optical motion capture systems,
276 dynamometers and in-person clinician assessment where applicable.

277 Outcomes (O): accuracy, expressed as measurement bias or agreement with the
278 reference standard, and reliability, including test-retest reliability/intra-rater reliability,
279 inter-rater reliability and agreement among instruments.

280

281 Accordingly, this review aims to address the following questions:

282 1. What is the accuracy of remote shoulder ROM measurement methods compared
283 with reference standard assessments?

284 2. What is the reliability of remote shoulder ROM measurement tools, including test-
285 retest and inter-rater reliability where reported?

286 3. How does measurement performance vary according to motion type, technology
287 type, population, and data acquisition methods?

288

289 The wording of the thesis objectives has been refined from the original PROSPERO
290 registration to improve clarity and alignment with the final included evidence. These
291 revisions do not change the underlying review question, eligibility criteria, or core
292 outcomes, but rather provide clearer framing of the same review aims for the reader.

293 **Methods**

294 **Protocol and Registration**

295 This systematic review follows the Preferred Reporting Items for Systematic Reviews
296 and Meta-Analyses (PRISMA) guidelines. The protocol was registered on
297 PROSPERO (Registration Number: CRD42024564283). PRISMA diagram is
298 available as **Figure 1**.

299

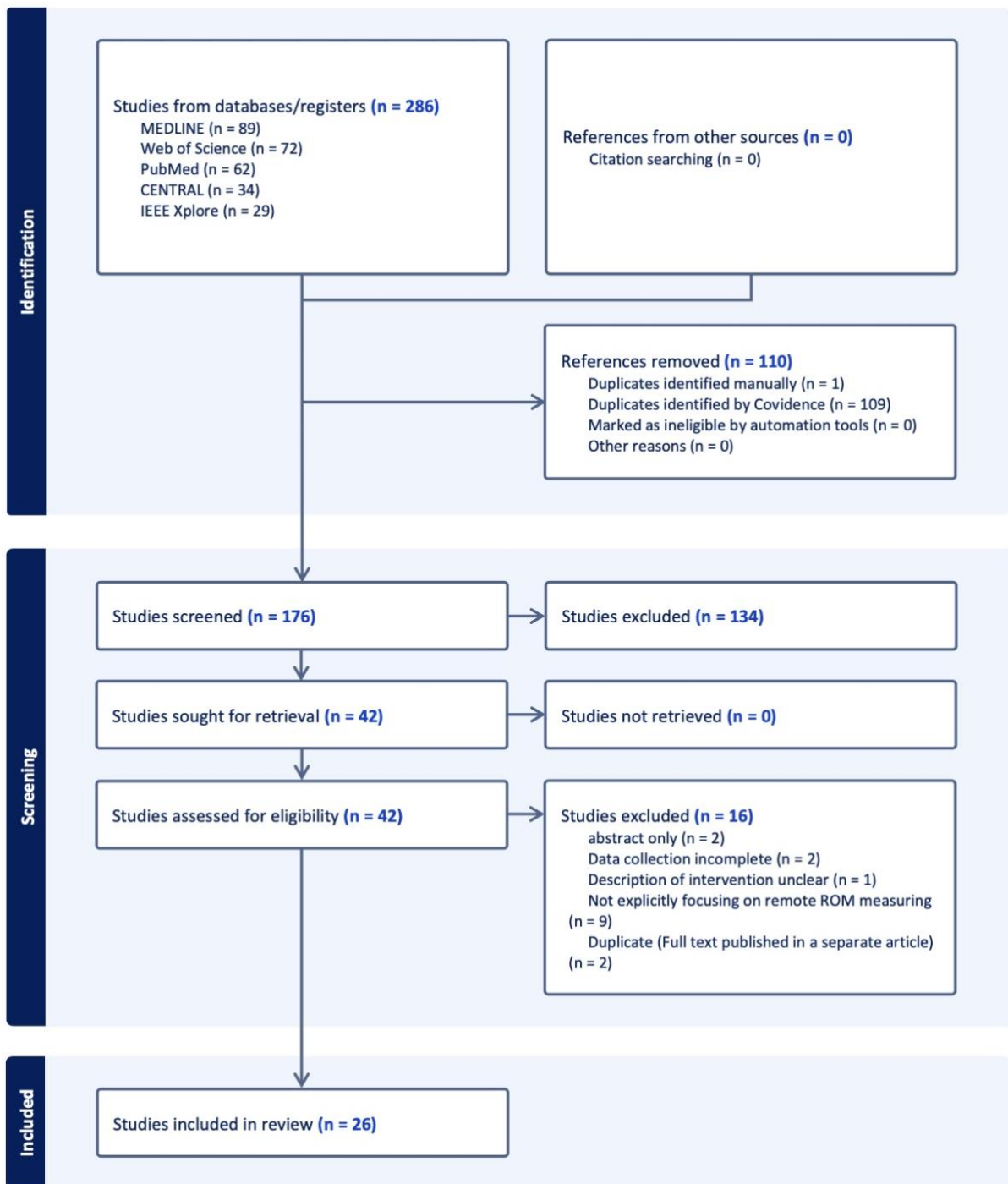


Figure 1 . PRISMA Diagram for study selection

300

301

302

303 Eligibility Criteria

304 The inclusion and exclusion criteria for this review were developed collaboratively by
 305 the authors, guided by the PICOS framework. Eligible studies focused on adults (18
 306 years or older) diagnosed with musculoskeletal shoulder conditions or healthy
 307 volunteers. Remote testing for shoulder ROM and strength using tools such as

308 wearable devices or mobile applications are included. While control comparators
309 were not required, any that were present in the studies were documented. Key
310 outcomes included accuracy and reliability outcomes reported as mean bias (in
311 degrees) or intraclass correlation coefficient (ICC). Accepted study designs included
312 randomised controlled trials (RCTs), quasi-experimental designs, prospective cohort
313 studies, comparative diagnostic accuracy studies, and case series.

314

315 Strength was considered a complementary outcome to range of motion, reflecting
316 overall movement quality and functional capacity of the shoulder. Studies assessing
317 shoulder function, including range of motion and/or muscle strength, were included
318 where these outcomes were relevant to remote assessment of movement.

319

320 Only full-text articles published in English were included. Conference abstracts,
321 editorials, and non-peer-reviewed sources were excluded due to limited
322 methodological detail. No restriction was placed on publication date, and databases
323 were initially searched from inception to May 2024.

324

325 Participants included both healthy individuals and patients with musculoskeletal
326 shoulder conditions, such as rotator cuff pathology, adhesive capsulitis, post-
327 operative rehabilitation cases, and other shoulder disorders affecting range of
328 motion.

329

330 Accuracy outcomes were defined as the agreement between remote measurement
331 methods and reference standards, typically expressed as mean bias or mean
332 difference in degrees, with associated measures of variability such as standard
333 deviation (SD) or limits of agreement (LOA). Reliability outcomes were assessed
334 using intraclass correlation coefficients (ICCs) to evaluate consistency of
335 measurements, including test-retest and inter-rater reliability where reported. The
336 choice of outcome measure was therefore dependent on whether studies assessed
337 agreement with a reference standard (accuracy) or consistency across repeated or
338 multiple measurements (reliability).

339

340 Studies were excluded if they assessed motion or strength in non-human subjects
341 (e.g., robotic arms), did not focus explicitly on remote shoulder ROM measurement

342 methods, or were unrelated to orthopaedic or musculoskeletal conditions.
343 Additionally, studies were excluded if the intervention lacked clarity, or if they were
344 not original research (such as systematic reviews, theses, or dissertations).

345

346 **Database and Search Strategy**

347 A systematic search was conducted across MEDLINE, Web of Science, PubMed,
348 CENTRAL and IEEE Xplore. The search strategy included terms related to shoulder
349 ROM, remote measurement, accuracy, and reliability. The reference lists of relevant
350 articles were manually screened for additional studies. An example database search
351 using keywords connected by Boolean operators is included as **Appendix 1**.

352 Appendix 1 provides the full search strategies used for each database, including
353 database-specific adaptations of search terms and Boolean operators, to allow
354 replication of the search process.

355

356 All databases were initially searched from inception to May 2024. To ensure the
357 currency of the review, an updated search was conducted in March 2026 using the
358 same search strategy. 33 additional studies were identified from all databases after
359 duplicate removal and assessed for eligibility. No additional studies meeting the
360 inclusion criteria were identified and there is no impact on the overall conclusion.

361

362 Trial registers (including ClinicalTrials.gov and the WHO International Clinical Trials
363 Registry Platform) were not systematically searched, as the focus of this review was
364 on published studies reporting validated accuracy and reliability outcomes. However,
365 reference list screening was undertaken to minimise the risk of missing relevant
366 studies.

367

368 **Study Selection**

369 Two independent reviewers (PC & MP) screened titles and abstracts, followed by
370 full-text review. Discrepancies were resolved through discussion or by consulting a
371 third reviewer (PE/AR).

372

373 **Data Extraction and Risk of Bias Assessment**

374 Data extraction included study characteristics, population details, measurement
375 methods, accuracy and reliability outcomes, and risk of bias assessment. The
376 QUADAS-2 was applied to all studies.⁵⁵

377

378 **Statistical Analysis**

379 For the quantitative synthesis, meta-analyses were conducted using STATA 18.0
380 (STATA Corp. 2023, TX). Mean bias was selected as the effect size to represent the
381 agreement between remote measurement methods and reference standards. The
382 mean bias is presented as remote methods minus the reference standards. Where
383 applicable, studies reporting mean difference (bias) and standard deviation (SD)
384 were directly included.

385

386 For studies that only reported 95% confidence intervals (CI) of the mean bias, the
387 standard deviation was recalculated using the formula:

$$SD = \frac{(Upper\ CI - Lower\ CI)}{3.92}$$

388

389 This approach assumes a normal distribution and aligns with established methods
390 for deriving SD from 95% CI ranges. The standard error (SE) was subsequently
391 calculated as:

$$SE = \frac{SD}{\sqrt{n}}$$

392

393 where n represents the number of shoulders measured in each study.

394

395 For studies that reported Bland-Altman limits of agreement (LOA) but did not report
396 SD, standard deviation was recalculated using the formula:

$$SD = \frac{Upper\ LOA - Lower\ LOA}{3.92}$$

397

398

399 Random-effects meta-analyses were performed using the restricted maximum
400 likelihood method to account for between-study heterogeneity. Subgroup analyses
401 were conducted based on motion type (flexion, abduction, extension, external

402 rotation [ER], and internal rotation [IR]), technology category (IMU vs. image-based),
403 participant characteristics (healthy vs. pathological shoulders), and data acquisition
404 methods (by patients vs. by assessors). Heterogeneity was assessed using the I^2
405 statistic. Differences between subgroups were assessed using a chi-squared (Q) test
406 for subgroup differences, as implemented in the random-effects meta-analysis
407 model.

408

409 Where substantial heterogeneity was identified ($I^2 > 50\%$), potential sources of
410 heterogeneity were explored through subgroup analyses based on clinically relevant
411 variables, including motion type, technology type, population characteristics, and
412 data acquisition methods. Formal meta-regression was not performed due to the
413 limited number of studies included in the quantitative synthesis. A 95% prediction
414 interval were calculated to estimate the expected range of true effects in future
415 studies, offering a more clinically meaningful interpretation of variability.

416

417 A funnel plot and Egger's regression test were used to assess publication bias.
418 Statistical significance was defined as $p < 0.05$. In accordance with PRISMA
419 guidance, included studies were also assessed for potential meta-bias through visual
420 inspection of funnel plot asymmetry and consideration of study size and reporting
421 patterns. No formal statistical correction for publication bias was applied.

422

423 Reliability outcomes were assessed using intraclass correlation coefficients (ICCs),
424 which quantify the consistency of measurements across repeated assessments or
425 between raters. ICC values were interpreted in the context of measurement
426 reliability, with higher values indicating greater agreement. Where multiple ICC
427 values were reported within a study, these were summarised in Table 1 to maintain
428 clarity and readability, with full values and confidence intervals considered during
429 data extraction.

430

431 The overall confidence in the cumulative evidence was considered in relation to
432 study quality, consistency of findings, and risk of bias. A formal GRADE assessment
433 was not undertaken due to the methodological heterogeneity and diversity of
434 outcome measures across included studies; however, key domains relevant to
435 evidence certainty were considered in the interpretation of findings.¹⁷

436

437 Where quantitative synthesis was not appropriate due to heterogeneity in study
438 design, outcome reporting, or measurement methods, a structured narrative
439 synthesis was undertaken. This approach is consistent with Cochrane guidance,
440 which recommends narrative or semi-quantitative synthesis when statistical pooling
441 is not feasible or meaningful. ^{6, 20, 42}

442

443 The narrative synthesis aimed to provide a transparent and systematic summary of
444 study findings, focusing on patterns in measurement approaches, reported accuracy
445 and reliability metrics, and methodological differences across studies. This approach
446 has been applied in previous reviews addressing heterogeneous intervention or
447 measurement studies, allowing for an analytic rather than purely descriptive
448 summary of the evidence.

449

450 **Ethics**

451 The methodology of this systematic review involves the review and synthesis of pre-
452 existing published data, thereby making this phase of the study exempt from ethical
453 approval.

454 **Results**

455 The results are presented in three stages: (1) an initial description of included
456 studies and their characteristics, (2) a narrative synthesis of studies not included in
457 the quantitative meta-analysis, and (3) quantitative synthesis of accuracy outcomes
458 using meta-analysis, followed by assessment of reliability and risk of bias.

459

460 **Study Selection and Characteristics**

461 A total of 26 studies met the inclusion criteria and were analysed in this review. Of
462 these, 17 studies focused on the accuracy of remote shoulder ROM measurement
463 tools and 10 were able to be included in the meta-analysis. Studies examined
464 reliability were narratively synthesised.

465

466 Of the 17 studies assessing accuracy, 7 were not included in the meta-analysis due
467 to methodological or reporting limitations, including lack of appropriate agreement
468 data (e.g. absence of mean bias or limits of agreement), reporting of compound or
469 non-isolated movements, or focus on outcomes not directly aligned with ROM
470 accuracy (e.g. strength or algorithm development without reference comparison).

471

472 The included studies utilised a range of remote measurement approaches, most
473 commonly smartphone-based applications and wearable inertial measurement units
474 (IMUs), with comparators including handheld goniometers, digital inclinometers,
475 optical motion capture systems, and in-person clinician assessment.

476

477 The studies varied in methodologies, sample sizes, and comparator tools, covering a
478 broad spectrum of remote measurement approaches, including smartphone
479 applications, wearable inertial measurement units (IMUs), and markerless motion
480 capture systems. The overall study characteristics are summarised in **Table 1**.

481 Table 1. Overall study characteristics summary. Abbreviations: ICC = intraclass correlation coefficient (measure of reliability for continuous data); κ = kappa coefficient (measure of agreement for
 482 categorical data); r = correlation coefficient (measure of linear association between variables); ROM = range of motion; IMU = inertial measurement unit. ICC values are reported as measures of test-retest
 483 or inter-rater reliability with individual 95% confidence interval, while κ values reflect agreement beyond chance for categorical assessments. KR-20 (Kuder–Richardson reliability) measures internal
 484 consistency reliability used for assessments with dichotomous items. Correlation coefficients (r -values) indicate the strength of association between measurement methods but do not represent agreement.

Author/ Year	Study Design	Population	Mean Age (years \pm SD)	No. of participants	Number of Shoulders Assessed	Remote ROM Method	Comparators	Shoulder Motion Metrics Studied	Validity	Reliability	Test-Retest Reliability/ Intra-rater Reliability	Inter-rater Reliability	Reliability among compared instruments
Bechard et al. (2020) ²	Measurement validity and reliability study	Healthy	Not reported	32	64	Force sensing device (FoSe)	Hand-Held dynamometer	Strength	Yes	Yes	N/A	N/A	ICC: ER: 0.74 (0.55–0.86) Left; 0.74 (0.56–0.87) Right Shoulder IR: 0.89 (0.80–0.95) Left; 0.73 (0.54–0.86) Right
Boissy et al. (2017) ³	Measurement validity study	Healthy	30.9 \pm 10.1	25	25	Digital Goniometer	Optic motion capturing system	Abduction, Flexion, ER	Yes	No			

Borresen et al. (2023) ⁴	Measurement reliability study	Pathological	Not reported	15	15	Augmented Reality Telerehabilitation System with Haptics	In-person assessment	Elevation, Depression, IR, ER, Adduction, Abduction, Protraction, Retraction	No	Yes	N/A	$\kappa =$ Elevation - 0.09 (-0.4 - 0.2) ER 0.25 (-0.2 - 0.7) Adduction - 0.11 (-0.3 - 0.0) Abduction 0.44 (-0.1 - 1.0) Protraction 0.44 (-0.1 - 1.0) Retraction 0.15 (-0.1 - 0.4)	N/A
Chan et al. (2022) ⁷	Measurement validity study	Healthy	25 (SD not reported)	19	19	Wearable IMU	Optic motion capturing system	Flexion, Extension, ER, Abduction	Yes	No			

Chen et al. (2020) ⁸	Measurement reliability study	Healthy & Pathological	53.0 ± 6.2 and 56.1 ± 13.3 (two groups)	25	25	IMU + Mobile App	Hand-Held Goniometer	Flexion, Extension, Abduction, IR, ER	No	Yes	N/A	ICC: Abduction: 0.97 (0.95–0.98) active; 0.98 (0.96–0.99) passive Flexion: 0.95 (0.92–0.97) active; 0.90 (0.83–0.94) passive Extension: 0.77 (0.64–0.87) active; 0.80 (0.68–0.89) passive ER: 0.95 (0.92–0.97) active; 0.96 (0.93–0.98) passive IR: 0.91 (0.86–0.95) active; 0.97 (0.94–0.98) passive	N/A
---------------------------------	-------------------------------	------------------------	---	----	----	------------------	----------------------	---------------------------------------	----	-----	-----	---	-----

Cui et al. (2019) ¹²	Measurement reliability study	Pathological	56.3 (SD not reported)	25	25	IMU + Virtual Reality	Hand-Held Goniometer	Flexion, Abduction, IR, ER	No	Yes	N/A	N/A	r-values: Flexion: 0.997 Abduction: 0.978 IR: 0.984 ER: 0.897
Cuesta-Vargas et al. (2016) ¹¹	Measurement validity and reliability study	Healthy & Pathological	Not reported	37	37	Mobile App	IMU	Abduction	Yes	Yes	ICC: Healthy: 0.78 (0.40 – 0.93) Pathological: 0.98 (0.94 – 0.99)	ICC: Healthy: 0.49 (0.08 – 0.82) Pathological: 0.99 (0.98 – 1.00)	N/A
Çubukçu et al. (2020) ¹⁰	Measurement validity and reliability study	Healthy	22.1± 3.1	40	80	Kinect V2	Hand-Held & Digital Goniometer	Abduction, Flexion, Extension, IR, ER	Yes	Yes	ICC (no CI reported): Abduction 0.86 Flexion: 0.85 Extension: 0.62 IR: 0.97 ER: 0.87	N/A	N/A
Gushikem et al. (2022) ¹⁶	Measurement validity and reliability study	Pathological	29.3 (SD not reported)	21	21	Tele-assessment	In-person assessment	Abduction, Strength	Yes	Yes	ICC: Abduction 0.87 (0.69 – 0.95)	N/A	N/A

Hwang et al. (2023) ²¹	Measurement validity and reliability study	Healthy & Pathological	35.2 (SD not reported)	10 Healthy, 10 Pathological	20 Surgeon Estimation	Digital Goniometer	Flexion, Abduction, ER, IR	Yes	Yes	N/A	ICC Healthy: Flexion: 0.93 (0.86–0.98); ER 0.90 (0.81–0.97); IR 0.77 (0.60–0.92). Pathological: Flexion 0.88 (0.77–0.96); ER 0.87 (0.75–0.96); IR 0.74 (0.57–0.91).	N/A
Jayaraman et al. (2020) ²²	Measurement validity study	Healthy	28 ± 4 (Phase 1) / 30 ± 3 (Phase 2)	75	75 Apple Watch	Hand-Held Goniometer	Abduction, Flexion, Extension	Yes	No			
Lam et al. (2024) ²⁷	Measurement validity and reliability study	Healthy	28.9 ± 11.8	30	30 Markerless Motion Capture System with iPad Pro	Hand-Held Goniometer	Abduction, Flexion	Yes	Yes	ICC (95% CI not reported)	N/A	N/A
Niu et al. (2024) ³⁰	Algorithmic optimisation & tracking accuracy study	Healthy	Not reported	16	16 Wearable smartwatch	Vicon Camera	Compound movements	No	No	(Static ranges): 0.17 to 0.80; ICC (Tasks ranges): 0.57 to 0.96		

Ongvisate paiboon et al. (2016) ³¹	Quasi-experimental study	Healthy	Not reported	20	20	Audio-biofeedback (ABF) system	Non-ABF system	Flexion	No	No			
Pereira et al. (2023) ³³	Measurement validity study	Healthy	Not reported	15	15	Mobile App using pose estimation	Qualisys Motion Capture System (QTM)	Compound movements	Yes	No			
Rajkumar et al. (2021) ³⁴	Measurement reliability study	Healthy	Not reported	17	34	Wearable IMU for exergames (WISE)	Kinect	Flexion-Extension, IR-ER, Abduction-Adduction (Compounded)	Yes	Yes	ICC (95% CI not reported): 0.46 to 0.90	N/A	ICC (95% CI not reported): 0.57 to 0.81
Ramkumar et al. (2018) ³⁵	Measurement validity study	Healthy	27.0 (SD not reported)	10	10	Smartphone App	Hand-Held Goniometer	Flexion, Abduction, IR, ER	Yes	No			
Roldán-Jiménez et al. (2019) ³⁸	Measurement reliability study	Healthy & Pathological	Not reported	16 Healthy, 27 Pathological	43	Smartphone in-built IMU	Wearable IMU sensors	Abduction, Flexion	No	Yes	N/A	N/A	ICC (Abduction): 0.86 to 0.97 (95% CI 0.73–0.99); ICC (flexoextension axis): 0.48 to 0.87 (95% CI 0.41–0.93); ICC (rotation axis): 0.43 to 0.62 (95% CI 0.10–0.80).

Sahu et al. (2022) ⁴¹	Measurement validity and reliability study	Healthy & Pathological	Not reported	24 Healthy, 16 Pathological	40	On-Screen App	Hand-Held Goniometer	Abduction, Flexion, IR, ER	Yes	Yes	ICC (consultant measurements): Abduction 0.97 (0.86- 0.99); Flexion 0.99 (0.97 – 0.99); IR 0.97 (0.94 – 0.99); ER 0.98 (0.94 – 1.0)	ICC: Abduction 0.86 (0.67 – 0.94); Flexion 0.94 (0.85 – 0.97); IR 0.95 (0.88 – 0.98); ER 0.97 (0.94 – 0.99).	
Seo et al. (2016) ⁴³	Measurement validity study	Healthy	28.5 (SD not reported)	10	10	Different Kinect sensor placements	3D Motion Capture	Flexion-Extension, IR-ER, Abduction-Adduction (Compounded)	Yes	No			
Shimizu et al. (2022) ⁴⁶	Measurement validity and reliability study	Healthy	Not reported	19	19	Smartphone App	Hand-Held Goniometer	Abduction, Flexion, Extension, IR, ER	Yes	Yes	ICC: Abduction 0.95 (0.91 – 0.98); Flexion 0.94 (0.88 – 0.98); Extension 0.79 (0.62 – 0.91); IR 0.91 (0.82 – 0.96); ER 0.93 (0.86 – 0.97).	N/A	N/A

Soeters et al. (2023) 47	Measurement validity and reliability study	Healthy	31.4± 11.7	30	60	Smartphone App	Hand-Held Goniometer	Abduction, Flexion, Extension, ER, IR	Yes	Yes	ICC: Abduction 0.95 (0.93 – 0.97); Flexion 0.94 (0.90 – 0.96); Extension 0.90 (0.84 – 0.93); ER 0.93 (0.88 – 0.95); IR 0.75 (0.62 – 0.84).	ICC: Abduction 0.90 (0.85 – 0.94); Flexion 0.96 (0.94 – 0.98); Extension 0.96 (0.93 – 0.97); ER 0.95 (0.92 – 0.97); IR 0.93 (0.89 – 0.95).	N/A
Tozawa et al. (2023) 50	Measurement reliability study	Healthy	Not reported	16	16	Zoom & Smartphone App	N/A	Flexion	No	Yes	ICC: Examiner 1: 0.91 (0.81 – 0.97) Examiner 2: 0.97 (0.93 – 0.99)	ICC: 0.95 (0.85 – 0.98)	N/A
Wang et al. (2020) 51	Measurement validity study	Healthy & Pathological	Not reported	30 Healthy, 1 Pathological	31	Kinect	N/A	Flexion, Abduction, Rotation	Yes	No			

Wang et al. (2022) 52	Measurement reliability study	Pathological	50.2± 16.2	32	32	Telemedicine examination	In-person assessment	Abduction, Adduction, Flexion, IR, ER, Strength	No	Yes	N/A	KR-20: Abduction 0.78 Adduction 0.88 Flexion 0.37 IR: 0.66 ER: 0.52 Strength: 0.57	KR-20: Abduction: 0.50 Adduction: 0.32 Flexion: 0.51 IR: 0.29 ER: 0.44 Strength: 0.37
Wang et al. (2023) 53	Measurement reliability study	Healthy	Not reported	25	25	Webcam based machine learning approach (Blazepose)	Optic motion capturing system (OptiTrack)	Flexion-Extension (Sagittal plane), Abduction-Adduction (Frontal plane)	No	Yes	ICC: Flexion 0.94 (0.90 – 0.96); Extension 0.93 (0.89 – 0.96); Abduction 0.97 (0.95 – 0.98); Adduction 0.92 (0.87 – 0.95).	N/A	ICC: Flexion 0.18 (-0.14 – 0.40); Extension 0.92 (0.88 – 0.94); Abduction 0.68 (0.56 – 0.77); Adduction 0.89 (0.84 – 0.92)

486 Overall study quality was mixed. Assessment of reporting bias revealed asymmetry
487 in the funnel plot and a statistically significant result from Egger's test ($p = 0.0051$),
488 suggesting the presence of publication bias. These methodological limitations should
489 be considered when interpreting the pooled meta-analytic results.

490

491 A notable characteristic of the included studies is the predominance of young,
492 healthy volunteer populations. Most studies excluded individuals with
493 musculoskeletal pathology, comorbidities, or functional impairments. As a result, the
494 generalisability of the findings to older adults or those with shoulder conditions is
495 limited.

496

497 Populations included both healthy individuals and patients with pathological shoulder
498 conditions such as rotator cuff tears, adhesive capsulitis, and post-operative
499 rehabilitation cases. ROM assessments covered flexion, abduction, extension,
500 internal rotation (IR), and external rotation (ER). The reference standards commonly
501 used were goniometry, optical motion capture, and clinician assessments.

502

503 **Studies evaluating accuracy but not included in the meta-analysis**

504 Seven studies evaluating the accuracy of remote shoulder assessment methods
505 were not included from the meta-analysis but narratively synthesised due to
506 methodological or reporting limitations but were nonetheless important for
507 understanding the broader context of remote measurement approaches.

508

509 Two of these studies, by Bechard and colleagues (2020) and Gushikem and
510 colleagues (2022), focused solely on strength assessments rather than range of
511 motion (ROM), which precluded their inclusion in a meta-analysis specifically
512 targeting ROM measurement agreement. They do, however, provide broad insights
513 in understanding important quality of shoulder motion metrics in addition to range of
514 motion alone.^{2, 16} Similarly, Niu and colleagues (2024) did not report any quantitative
515 data on accuracy or reliability. Instead, it emphasised technical advancements in
516 real-time motion tracking using ultrasound and inertial sensors, prioritising system
517 development and algorithmic accuracy without comparing against established
518 reference standards or applying statistical tools like Bland-Altman analysis. Although
519 not contributing directly to the pooled analysis, this study offers valuable insights into
520 emerging technologies.³⁰

521

522 Cuesta-Vargas and colleagues provided ROM data using both inertial sensors and
523 smartphone-based methods in healthy and pathological populations. In healthy
524 participants, mean \pm SD values were $169.07^\circ \pm 4.96^\circ$ using smartphone
525 measurements and $154.22^\circ \pm 19.27^\circ$ using inertial sensors. In pathological
526 participants, corresponding values were $93.54^\circ \pm 40.88^\circ$ and $87.86^\circ \pm 47.41^\circ$,
527 respectively. The study, however, lacked the required paired differences between
528 measurements for appropriate Bland-Altman analysis, instead reporting only mean
529 and standard deviation values for separate groups. Without direct within-subject
530 comparisons, the assumptions required to estimate limits of agreement may not be
531 valid, making the results unsuitable for meta-analysis.¹¹

532

533 Similarly, Rajkumar and colleagues (2021) reported data from a Bland-Altman
534 analysis comparing Kinect and WISE systems but presented flexion-extension as a
535 combined metric. Mean root mean square error (RMSE) values ranged from $7.72^\circ \pm$
536 2.82° to $11.29^\circ \pm 3.10^\circ$ across exercises, with lower errors observed following

537 dynamic time warping ($6.28^\circ \pm 2.57^\circ$ to $9.62^\circ \pm 3.88^\circ$). Because the meta-analysis
538 required individual movement assessments (e.g., flexion and extension separately),
539 this aggregation limited compatibility with pooled estimates. ³⁴

540

541 Pereira and colleagues (2023) assessed compound shoulder movements rather than
542 isolated ROM parameters, comparing mobile app-based pose estimation, motion
543 capture (QTM), and DTW algorithms. the mobile app-based pose estimation
544 demonstrated a mean score of 94.08 ± 1.18 , indicating high consistency in
545 movement pattern detection. In contrast, DTW method showed greater variability,
546 with a mean of 53.74 ± 25.59 . The reference QTM system demonstrated wider
547 dispersion (mean 78.8 ± 21.61), reflecting sensitivity to variations in movement
548 execution. The study found wide variability in agreement between methods,
549 especially between DTW and QTM, but its focus on multidimensional movements
550 prevented integration with studies measuring isolated joint angles. ³³

551

552 Lastly, Seo and colleagues (2016) offered critical contextual information, showing
553 that Kinect may systematically overestimate ROM due to depth perception
554 limitations. For shoulder elevation, the mean ROM error was $5^\circ \pm 1^\circ$, representing
555 the lowest error among the assessed joints. Across sensor configurations, the lowest
556 overall error was observed when the Kinect was positioned at 45° elevation in front
557 of the subject ($10^\circ \pm 1^\circ$), whereas the conventional frontal placement resulted in
558 higher error ($22^\circ \pm 2^\circ$), and the greatest error was observed at contralateral
559 positioning ($38^\circ \pm 2^\circ$). While this does not lend itself to quantitative synthesis, it offers
560 a plausible explanation for overestimation trends observed in other remote ROM
561 studies and is therefore valuable in interpreting pooled results. ⁴³

562

563 **Accuracy of Remote Shoulder ROM Measurement**

564 The meta-analysis included studies that quantitatively assessed the accuracy of
565 remote ROM measurement technologies against reference standards. Four separate
566 subgroup meta-analyses were performed:

- 567 1. By motion type
- 568 2. By technology type
- 569 3. By population (healthy vs. pathological)
- 570 4. By data acquisition methods

571

572 **Overall Accuracy Findings**

573 Across all studies, the difference of the means was 2.63° (95% CI: 1.52, 3.74),
574 indicating a small but consistent overestimation of shoulder ROM when measured
575 remotely. There was significant heterogeneity ($I^2 = 98.62\%$), reflecting variations in
576 study methodologies and sample populations. The overestimation tendency was
577 observed across all movement planes, with statistically significant overestimation
578 primarily driven by external rotation (ER), flexion, and internal rotation (IR).

579

580 **Accuracy by Motion Type (Figure 2)**

581 When analysed by shoulder movement type, the difference of the means were as
582 follows:

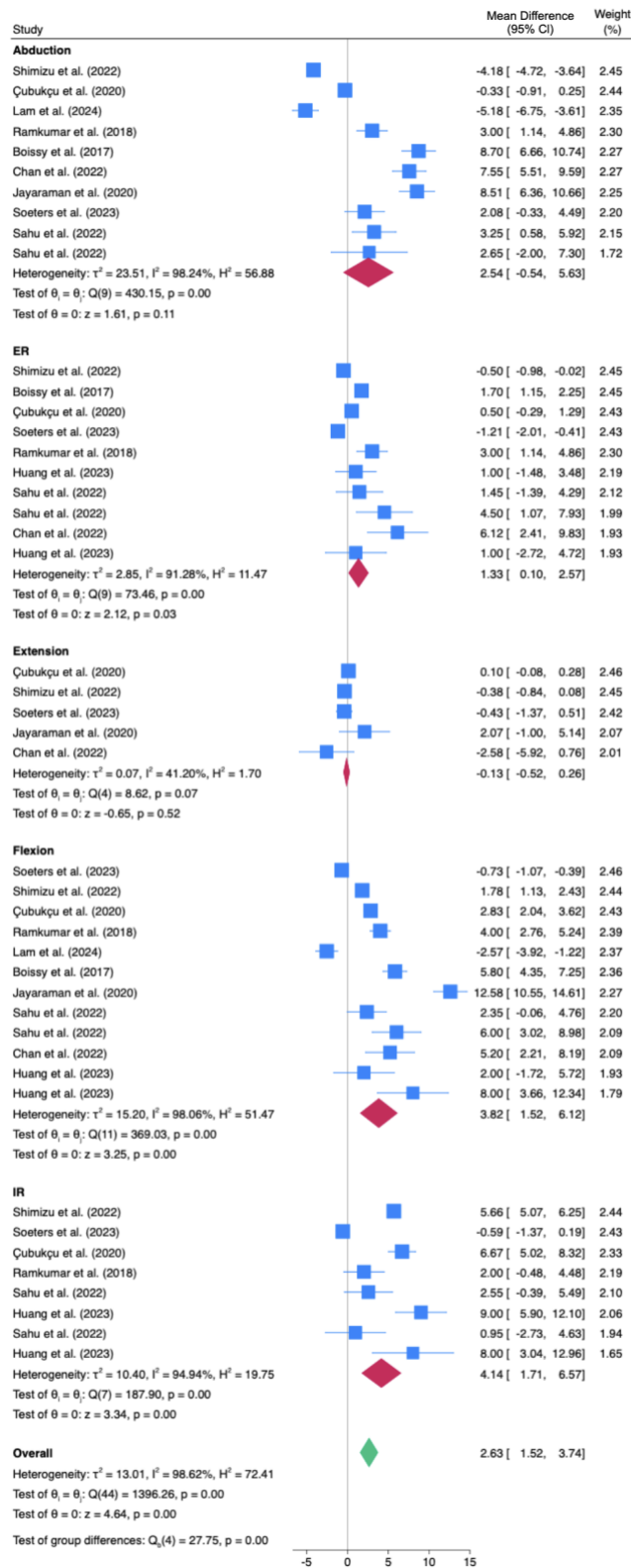
- 583 • **Abduction:** 2.54° (95% CI: -0.54, 5.63; $I^2 = 98.24\%$), which was not
584 significantly different from reference standard methods. The 95% prediction
585 interval is - 9.22° to 14.30°.
- 586 • **External Rotation (ER):** 1.33° (95% CI: 0.10, 2.57; $I^2 = 91.28\%$), showing a
587 small but statistically significant difference. The 95% prediction interval is -
588 2.83° to 5.49°.
- 589 • **Extension:** -0.13° (95% CI: -0.52, 0.26; $I^2 = 41.20\%$), which showed no
590 significant difference between remote and reference standard methods. The
591 95% prediction interval is - 1.18° to 0.92°.
- 592 • **Flexion:** 3.82° (95% CI: 1.52, 6.12; $I^2 = 98.06\%$), indicating significant
593 overestimation. The 95% prediction interval is - 5.37° to 13.01°.

594 • **Internal Rotation (IR):** 4.14° (95% CI: 1.71, 6.57; I² = 94.94%), showing the
595 largest measurement discrepancy. The 95% prediction interval is – 4.31° to
596 12.59°.

597

598 Notably, remote measurements consistently overestimated ROM across most
599 movements, except for abduction and extension, where no systematic over- or
600 underestimation was observed.

601



602
603
604
605
606
607

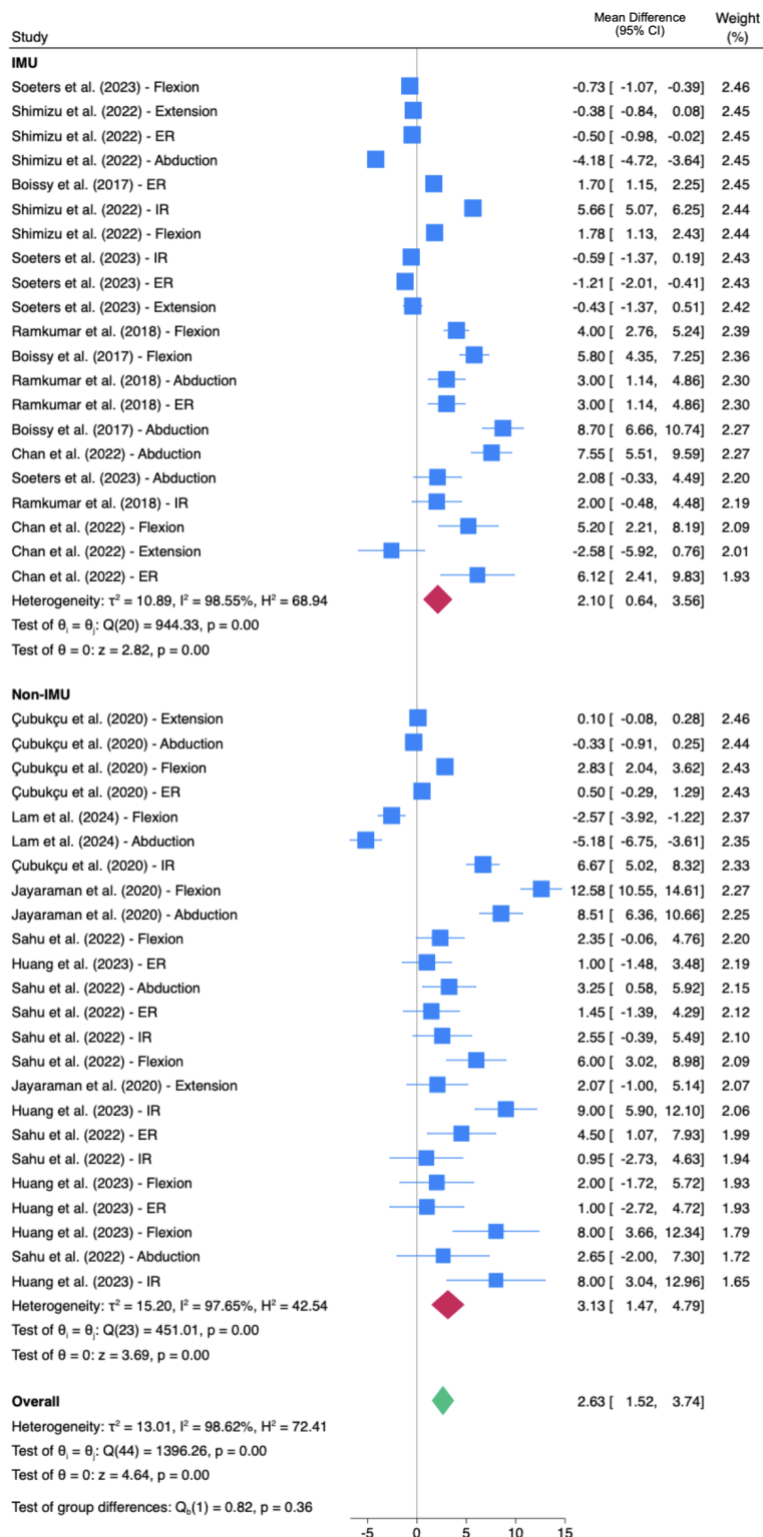
Figure 2. Forest plot of the meta-analysis of accuracy of remote shoulder assessment by motion type. Mean differences (in degrees) between remote and reference measurements are presented as point estimates (dots), with 95% confidence intervals (horizontal bars). Negative values indicate underestimation by the remote method, while positive values indicate overestimation. Statistical heterogeneity is reported using the I^2 statistic.

608 **Accuracy by technology types (Figure 3)**

- 609 • **IMU subgroup:** The difference of the means was 2.10° (95% CI: 0.64°, 3.56°;
610 $I^2 = 98.55\%$), suggesting that, on average, IMU-based methods slightly
611 overestimate shoulder ROM when compared to reference standards. This
612 result was statistically significant. The 95% prediction interval is – 5.08° to
613 10.05°.
- 614 • **Non-IMU (Image-based technologies) subgroup:** The difference of the
615 means was 3.13° (95% CI: 1.47°, 4.79°; $I^2 = 97.65\%$), also indicating an
616 overestimation with non-IMU methods, again statistically significant. The 95%
617 prediction interval is – 6.05° to 12.31°.

618

619 The test of subgroup differences yielded a p-value of 0.36, indicating that there was
620 no statistically significant difference between IMU and non-IMU methods in terms of
621 difference of the means.



622

623

624

625

626

627

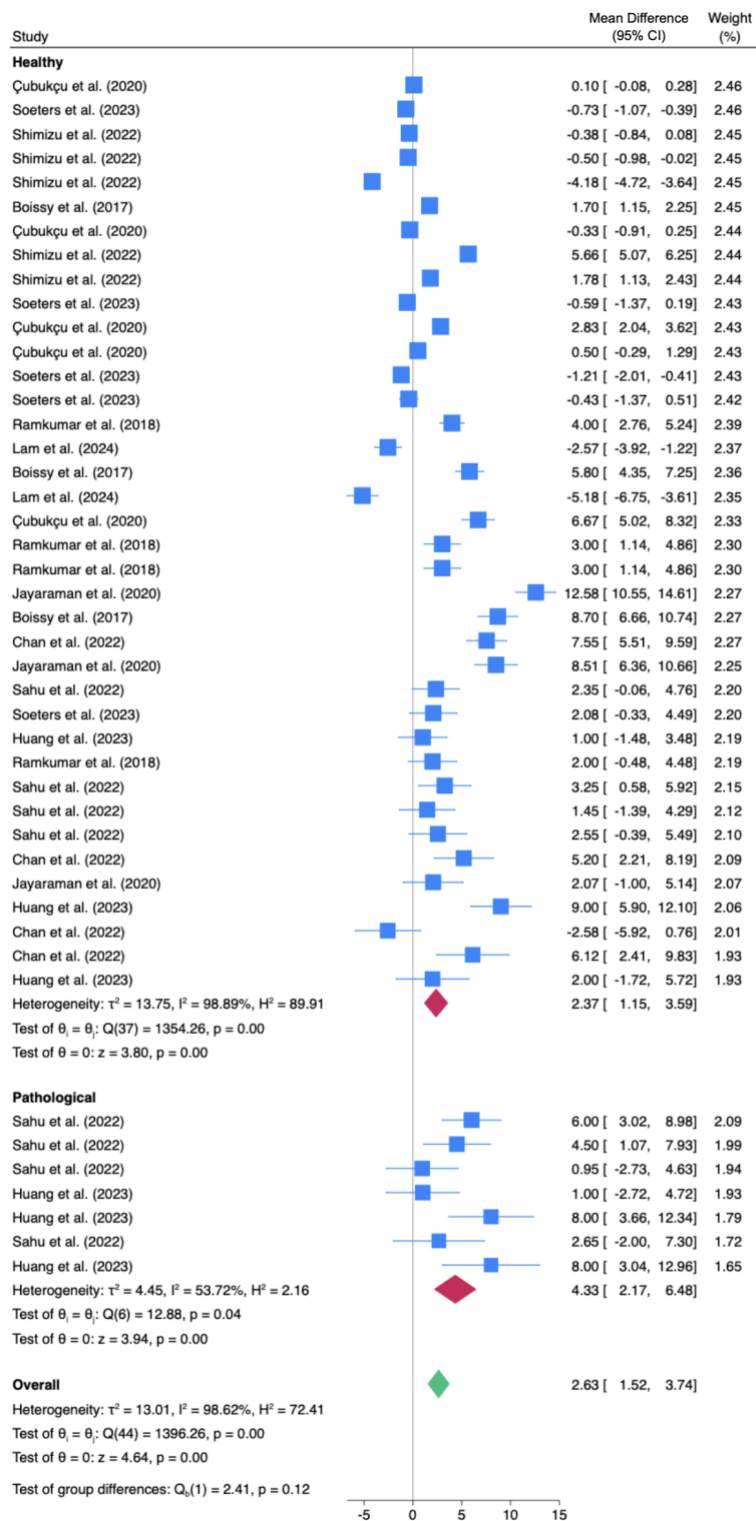
Figure 3. Forest plot of the meta-analysis of accuracy of remote shoulder assessment by technology types (IMU versus non-IMU). Mean differences (in degrees) between remote and reference measurements are presented as point estimates (dots), with 95% confidence intervals (horizontal bars). Negative values indicate underestimation by the remote method, while positive values indicate overestimation. Statistical heterogeneity is reported using the I^2 statistic.

628 **Accuracy in Healthy vs. Pathological Shoulders (Figure 4)**

- 629 • Healthy participants had a difference of the means of 2.37° (95% CI: 1.15,
630 3.59; $I^2 = 98.89\%$), suggesting overestimation. The 95% prediction interval is
631 – 5.54° to 10.28°.
- 632 • Pathological shoulders exhibited a difference of the means of 4.33° (95% CI:
633 2.17, 6.48; $I^2 = 53.72\%$), reflecting higher variability in measurement error.
634 The 95% prediction interval is – 0.53° to 9.19°.

635

636 While the test of subgroup differences was non-significant ($p = 0.12$), suggesting no
637 statistically significant difference in measurement between healthy and pathological
638 populations, the observed differences was higher in the pathological group (4.33° vs.
639 2.37°).



640

641

642

643

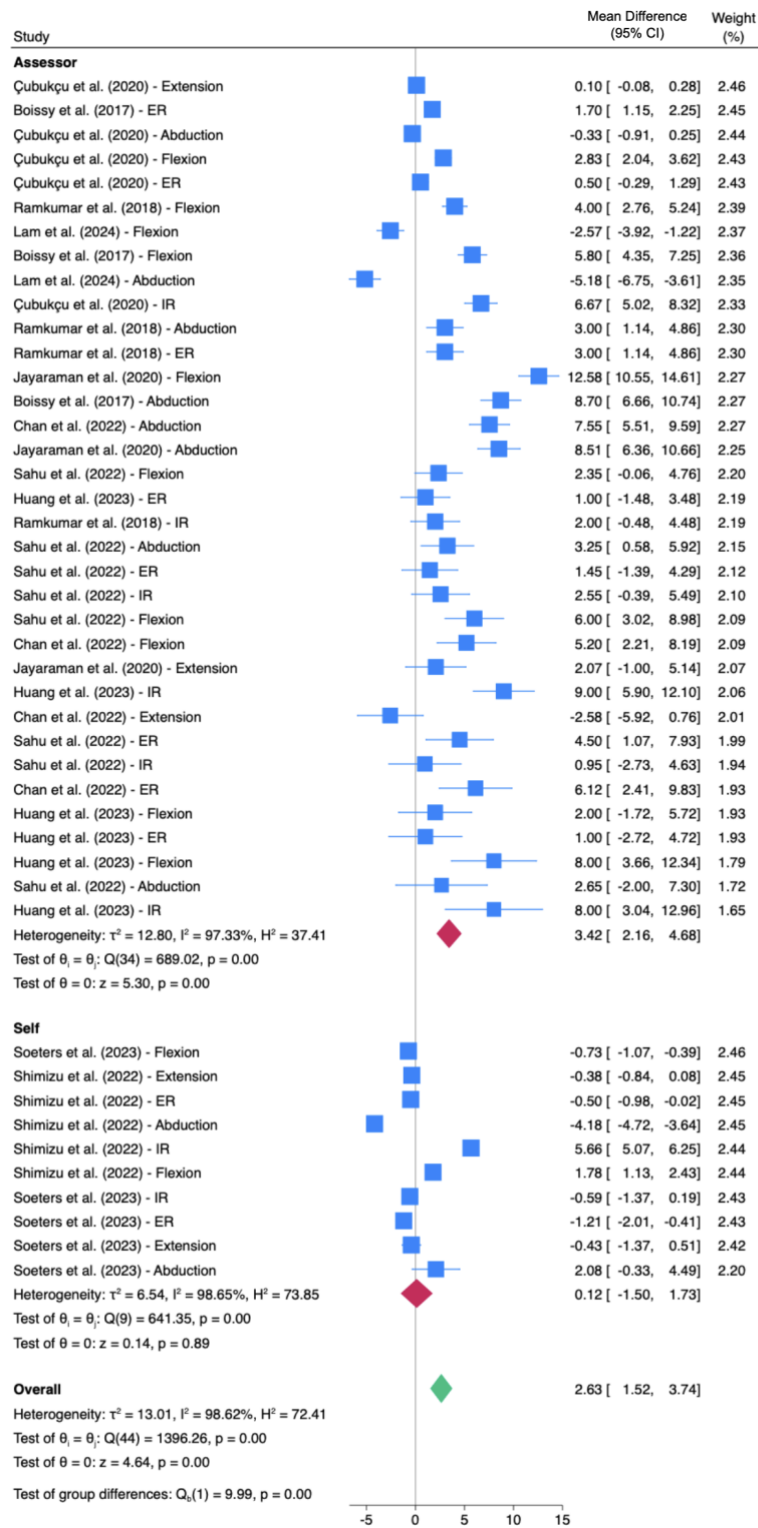
644

645

Figure 4. Forest plot of the meta-analysis of accuracy of remote shoulder assessment by population types (healthy vs. pathological). Mean differences (in degrees) between remote and reference measurements are presented as point estimates (dots), with 95% confidence intervals (horizontal bars). Negative values indicate underestimation by the remote method, while positive values indicate overestimation. Statistical heterogeneity is reported using the I^2 statistic.

646 **Accuracy by data acquisition methods (Figure 5)**

- 647
- For assessor-guided measurements, the difference of the means was 3.42° (95% CI: 2.16 to 4.68; $I^2 = 97.33\%$), suggesting a statistically significant overestimation. The 95% prediction interval is – 4.80° to 11.64°.
- 648
- For self-measurements, the difference of the means was 0.12° (95% CI: -1.50 to 1.73; $I^2 = 98.65\%$), which was not statistically significant. The 95% prediction interval is – 4.98° to 5.22°.
- 649
- 650
- 651
- 652



653

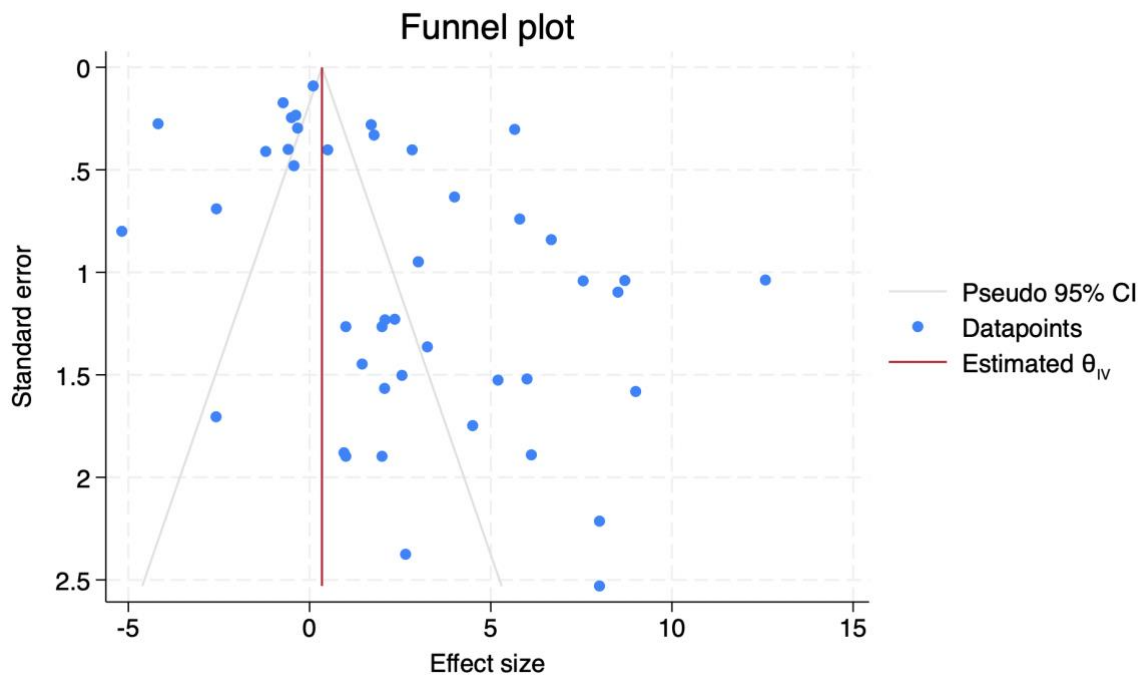
654 *Figure 5. Forest plot of the meta-analysis of accuracy of remote shoulder assessment by data acquisition methods (assessor*
 655 *vs. self). Mean differences (in degrees) between remote and reference measurements are presented as point estimates*
 656 *(dots), with 95% confidence intervals (horizontal bars). Negative values indicate underestimation by the remote method,*
 657 *while positive values indicate overestimation. Statistical heterogeneity is reported using the I^2 statistic.*

658

659 **Publication Bias**

660 The funnel plot (**Figure 6**) showed an asymmetric distribution of 45 datapoints
661 reported from 10 studies, suggesting potential small-study effects. Egger's test
662 confirmed a statistically significant bias ($p = 0.0051$), indicating the possibility of
663 publication bias.

664



665

666

Figure 6. Funnel plot of the included data

667

668 **Reliability of Remote Shoulder ROM Measurement**

669 A total of 17 studies reported reliability data relevant to remote shoulder ROM
670 measurement. Reliability was assessed using various methods, including test-retest,
671 inter-rater, and inter-instrument comparisons, across both healthy and pathological
672 populations.

673

674 **Test-Retest Reliability**

675 Reliability outcomes were reported using both intra-rater and test–retest reliability
676 metrics. Intra-rater reliability refers to the consistency of measurements taken by the
677 same assessor across repeated trials, whereas test–retest reliability reflects the
678 stability of a measurement over time under similar conditions. Although these terms
679 are sometimes used interchangeably in the literature, they represent distinct
680 constructs and have been reported accordingly in this review.

681

682 Wang and colleagues (2023) reported test–retest reliability (referred to as intra-rater
683 reliability in the study) using repeated measurements under the same conditions with
684 a BlazePose-based machine learning system. ICCs range from 0.92 (95% CI 0.87 to
685 0.95) for shoulder adduction to 0.97 (95% CI 0.95 to 0.98) for shoulder abduction.⁵³
686 Sahu and colleagues (2022) also demonstrated high test-retest reliability in patients
687 using a protractor app measured by consultant raters, with ICCs ranging from 0.97
688 (95% CI 0.86-0.99) for shoulder abduction to 0.99 (95% CI 0.97-0.99) for shoulder
689 flexion.⁴⁰ Similarly, Çubukçu and colleagues (2020) found high intra-rater reliability
690 with ICCs between 0.62 for shoulder extension and 0.97 for internal rotation using a
691 Kinect V2 system, while Shimizu and colleagues (2022) observed ICCs varied
692 depending on the range of movement from 0.79 (95% CI 0.62 to 0.91) for extension
693 to 0.95 (95% CI 0.91-0.98) for abduction with a smartphone self-measurement tool.

694 ^{10, 46}

695

696 Lam and colleagues (2024) however, highlighted greater variability in intra-rater
697 reliability using a dual-iPad LiDAR setup, with ICCs ranging from 0.17 to 0.96 across
698 static and functional tasks (95% CI not reported).²⁷ Tozawa and colleagues (2023)
699 showed excellent intra-rater reliability for shoulder flexion between two separate
700 examiners (ICC = 0.91–0.97).⁵⁰ Gushikem and colleagues (2022) reported intra-

701 rater ICCs of 0.87 (95% CI 0.69 – 0.95) for abduction strength assessments in
702 patients with brachial plexus injury. ¹⁶

703

704 **Inter-Rater Reliability**

705 Inter-rater reliability was reported across a range of settings. Soeters and colleagues
706 (2023) reported excellent ICCs between 0.90 (95% CI 0.85 – 0.94) for abduction and
707 0.96 (95% CI 0.94 – 0.98) for flexion when comparing smartphone app and
708 goniometer measures. ⁴⁸ Tozawa and colleagues (2023) showed inter-rater ICC of
709 0.95 (95% CI 0.85 – 0.98) for shoulder flexion. ⁵⁰ Chen and colleagues (2020) found
710 inter-rater ICCs from 0.77 (95% CI 0.64 – 0.87) for active shoulder extension to 0.98
711 (95% CI 0.96 – 0.99) for passive shoulder abduction. ⁸

712

713 In contrast, Borresen and colleagues (2023) reported no agreement to moderate
714 agreement reporting with kappa statistics, with $\kappa = -0.11$ (95% CI -0.3 – 0.0) for
715 shoulder adduction to 0.44 (95% CI -0.1 – 1.0) for shoulder abduction. Cuesta-
716 Vargas and colleagues reported varied ICCs for inter-rater reliability ranging from
717 0.49 (95% CI 0.08 – 0.82) for healthy volunteers to 0.99 (95% CI 0.98 – 1.00) for
718 shoulders with pathology present.

719

720 **Agreement Between Instruments**

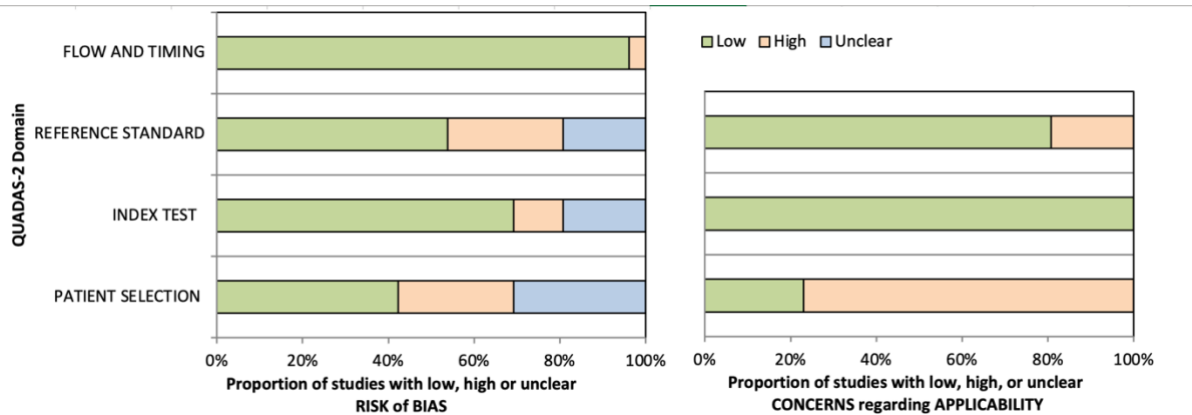
721 Studies comparing remote measurement tools against reference standards or
722 between different remote systems showed a range of agreement. Rajkumar and
723 colleagues (2021) found ICCs between 0.57 and 0.81 (95% CI not reported) for
724 various shoulder movements when comparing WISE and Kinect systems. ³⁴ Roldán-
725 Jiménez and colleagues (2019) reported excellent agreement between smartphone-
726 based and dedicated IMU sensors for abduction ICC 0.86 (95% CI 0.73-0.99), with
727 lower values for internal rotation ICC 0.43 (95% CI 0.10-0.80). ³⁸ Cui and colleagues
728 (2019) demonstrated strong correlations from $r = 0.897$ for shoulder external rotation
729 to $r = 0.997$ for shoulder flexion between wearable sensors integrated with virtual
730 reality and therapist assessments. ¹² Bechard and colleagues (2020) evaluated the
731 agreement between a novel strength device and handheld dynamometry, reporting
732 ICCs from 0.73 (95% CI 0.54-0.86) for right sided shoulder internal rotation to 0.89
733 (95% CI 0.80-0.95). ² Wang and colleagues (2022) used KR-20 to evaluate
734 dichotomous findings across telemedicine and in-person sessions, with varied

735 agreement values ranging from 0.29 for internal rotation to 0.51 (no 95% CI
736 reported) for flexion depending on the motion type and method. ⁵²

737 **Risk of Bias Assessment**

738 Risk of bias was assessed using the QUADAS-2 tool, with pre-defined signalling
739 questions tailored to the context of accuracy and reliability of instruments in remote
740 shoulder range of motion (ROM) measurements (See **Appendix 2**). A visual
741 summary of the QUADAS-2 domain assessments is shown in **Figure 7**, and detailed
742 study-level ratings are available in **Table 2**.

743



744

745

Figure 7. Visual Representation of Risk of Bias using QUADAS-2 Reporting Tool

746


Table 2. Tabular presentation for QUADAS-2 results

Study	RISK OF BIAS				APPLICABILITY CONCERNS		
	PATIENT SELECTION	INDEX TEST	REFERENCE STANDARD	FLOW AND TIMING	PATIENT SELECTION	INDEX TEST	REFERENCE STANDARD
Bechard et al. (2020)	?	😊	?	😊	😞	😊	😊
Boissy et al. (2017)	?	😞	?	😊	😞	😊	😊
Borresen et al. (2023)	😊	😊	😊	😊	😊	😊	😊
Chan et al. (2022)	😞	😞	😞	😊	😞	😊	😊
Chen et al. (2020)	😞	😞	😊	😊	😞	😊	😊
Cui et al. (2019)	😞	😊	😊	😊	😞	😊	😊
Cuesta-Vargas et al. (2016)	😊	?	😊	😊	😊	😊	😊
Çubukçu et al. (2020)	😊	😊	😞	😊	😞	😊	😊
Gushikem et al. (2022)	?	😊	😞	😊	😊	😊	😞
Hwang et al. (2023)	😞	😊	😊	😊	😞	😊	😊

Jayaraman et al. (2020)	?	?	☹️	😊	☹️	😊	😊
Lam et al. (2024)	😊	?	?	😊	☹️	😊	😊
Niu et al. (2024)	?	?	?	😊	☹️	😊	😊
Ongvisatepaiboon et al. (2016)	?	😊	☹️	?	☹️	😊	☹️
Pereira et al. (2023)	☹️	😊	☹️	😊	☹️	😊	☹️
Rajkumar et al. (2021)	?	😊	☹️	😊	☹️	😊	☹️
Ramkumar et al. (2018)	😊	😊	😊	😊	☹️	😊	😊
Roldán-Jiménez et al. (2019)	😊	😊	😊	😊	😊	😊	☹️
Sahu et al. (2022)	☹️	😊	😊	😊	😊	😊	😊
Seo et al. (2016)	😊	😊	😊	😊	☹️	😊	😊
Shimizu et al. (2022)	☹️	😊	😊	😊	☹️	😊	😊
Soeters et al. (2023)	😊	😊	😊	😊	☹️	😊	😊
Tozawa et al. (2023)	😊	😊	😊	😊	☹️	😊	😊
Wang et al. (2020)	?	?	?	😊	☹️	😊	😊

Wang et al. (2022)							
Wang et al. (2023)							

748

 Low Risk  High Risk  Unclear Risk

749 Out of 26 studies included, patient selection was the most frequent source of bias,
750 with the majority of studies (approximately 60%) rated as high or unclear risk in this
751 domain. This was primarily due to the predominance of healthy volunteer cohorts,
752 limiting applicability to clinical populations with shoulder pathology.

753

754 The index test domain demonstrated generally low risk of bias, although
755 approximately 25% of studies had unclear methods of test administration or lacked
756 sufficient detail on blinding procedures. The reference standard domain also showed
757 moderate concerns, with some studies failing to provide adequate justification or
758 validation for the gold-standard comparator. Flow and timing were mostly well
759 handled, with nearly all studies administering both index and reference tests within a
760 clinically appropriate timeframe.

761

762 Applicability concerns were similarly concentrated in patient selection, where around
763 80% of studies showed high concern due to unrepresentative populations. In
764 contrast, applicability concerns for the index test and reference standard were
765 minimal across the dataset.

766 **Discussions**

767 This systematic review and meta-analysis evaluated the accuracy and reliability of
768 remote shoulder ROM measurement methods. The results suggest that existing
769 remote measurements consistently overestimated ROM. While many remote
770 measurement technologies demonstrated acceptable reliability, their agreement with
771 reference standards was variable.

772

773 The overall mean bias of 2.63° (95% CI: 1.52, 3.74) indicates that remote ROM
774 measurements tend to slightly overestimate shoulder motion compared to reference-
775 standard assessments (**Figure 2**). This overestimation was most notable in external
776 rotation, flexion, and internal rotation, while extension showed the smallest
777 discrepancy. Due to the substantial heterogeneity across studies ($I^2 = 98.62\%$, $\tau^2 =$
778 13.01), we also calculated a 95% prediction interval to better reflect the expected
779 range of effects in future applications. The resulting interval (-5.79° to 11.05°)
780 reveals considerable variability, including the possibility of underestimation in some
781 contexts. This wide range highlights the influence of differing populations,
782 methodologies, and technologies, and calls for caution in generalising the pooled
783 estimate. An alternative interpretation of these findings is that the variability observed
784 across studies may reflect limitations in the accuracy and precision of the
785 technologies themselves, rather than solely differences in study design or
786 methodology. The presence of wide limits of agreement and large prediction
787 intervals suggests that, in some cases, remote measurement systems may produce
788 inconsistent or imprecise estimates of shoulder range of motion.

789

790 Although the mean bias was statistically significant, both the pooled estimate and the
791 prediction interval fall below the minimum clinically important difference (MCID) for
792 glenohumeral joint ROM, which ranges from 11° to 16° when assessed by a single
793 evaluator, and 14° to 24° when assessed by two. This suggests that, while
794 measurable, the observed bias may not be clinically meaningful in most cases.²⁹

795

796 Both IMU and non-IMU image-based motion tracking systems demonstrate a
797 statistically significant overestimation of ROM, the magnitude of bias remains
798 relatively small (**Figure 3**): Sensor drift and recalibration issues – Over prolonged

799 use, IMU sensors may accumulate minor errors, leading to inconsistent readings;
800 The findings of Seo and colleagues (2016) demonstrate that Kinect-based shoulder
801 ROM measurement can achieve relatively low error under optimal conditions ($\sim 5^\circ$),
802 although accuracy is highly dependent on sensor positioning.⁴³ This variability is
803 consistent with the wide prediction intervals observed in the present meta-analysis
804 and highlights the importance of methodological factors in influencing measurement
805 error. Placement inconsistencies – Small variations in sensor positioning between
806 measurement sessions or between different raters can introduce additional error;
807 Soft tissue movement artifacts – Unlike rigidly mounted motion capture markers, IMU
808 sensors placed on the skin may shift slightly due to soft tissue deformation and
809 muscle movement, leading to variations in recorded ROM.^{14, 19, 39} Importantly, these
810 biases fall below the MCID thresholds. Therefore, while measurement variability
811 exists, remote ROM assessments using either method are unlikely to introduce
812 clinically meaningful discrepancies in routine monitoring.

813

814 Although pathological shoulders demonstrated a larger mean measurement
815 difference compared to healthy participants, the test of subgroup differences was not
816 statistically significant (**Figure 4**). Therefore, any apparent increase in variability
817 should be interpreted with caution and may reflect differences in data distribution
818 rather than a true subgroup effect.

819

820 One possible explanation is that pathological shoulder conditions introduce greater
821 variability in movement execution, due to factors such as pain, restricted range,
822 compensatory movement patterns, and reduced motor control. These factors may
823 lead to less consistent positioning and movement during assessment, increasing
824 measurement variability. Alternatively, the observed variability may reflect limitations
825 in the accuracy and precision of the measurement technologies themselves when
826 applied to pathological populations. For example, altered movement patterns or
827 reduced joint range may challenge tracking algorithms or sensor alignment,
828 potentially reducing measurement accuracy. Findings from Cuesta-Vargas and
829 colleagues (2016) support the variability observed in the present meta-analysis.
830 Substantial differences in mean arm abduction angle and variability were observed
831 between smartphone and inertial sensor measurements, particularly in pathological
832 populations.¹¹

833

834 Interestingly, subgroup analysis based on who performed the ROM measurement
835 revealed that assessor-guided measurements (n = 8 studies) had a pooled mean
836 bias of 3.42° (95% CI: 2.16 to 4.68), indicating statistically significant overestimation
837 **(Figure 5)**. In contrast, self-measurements (n = 2 studies) demonstrated a much
838 smaller and non-significant bias of 0.12° (95% CI: -1.50 to 1.73). While these
839 findings suggest that self-measurement—when properly guided—may yield
840 comparable accuracy with tighter agreement, the limited number of self-assessment
841 studies precludes definitive conclusions. Nonetheless, this highlights the potential of
842 patient-led measurements in remote settings, warranting further investigation in
843 future research.

844

845 17 studies assessed the reliability of remote shoulder ROM measurements,
846 encompassing test-retest, inter-rater, and inter-instrument comparisons across
847 healthy and pathological populations. The wide variation in inter-rater ICC values
848 across studies has important clinical implications. High ICC values (e.g. >0.80)
849 indicate good agreement between assessors and suggest that remote measurement
850 tools may be reliable in controlled or optimised conditions. However, lower ICC
851 values (e.g. <0.60) indicate reduced agreement and raise concerns regarding
852 consistency between users. In a clinical context, this variability suggests that
853 measurements obtained using remote technologies may differ depending on the
854 assessor, which could impact clinical decision-making, monitoring of progress, and
855 treatment planning. This inconsistency may limit the reliability of these technologies
856 in routine practice, particularly in settings where multiple clinicians or patients are
857 involved in measurement, such as tele-rehabilitation or self-assessment contexts.

858

859 The observed variability in ICC values may be attributed to several factors, including
860 differences in assessor experience, variability in patient movement execution, and
861 sensitivity of the technology to positioning or environmental conditions. These factors
862 may disproportionately affect remote assessment tools compared to traditional in-
863 person measurement methods.

864

865 It is important to distinguish between intraclass correlation coefficients (ICC) and
866 kappa statistics, as they assess different types of agreement. ICC is used for

867 continuous data and evaluates the consistency of quantitative measurements
868 between raters or across time. In contrast, kappa statistics assess agreement for
869 categorical data, such as classification of movement quality or presence/absence of
870 impairment. As these metrics evaluate different constructs, differences in reported
871 values between ICC and kappa are expected and do not represent conflicting
872 results. Consequently, studies reporting kappa values may reflect agreement in
873 clinical classification, whereas ICC values provide insight into measurement
874 precision. Both are important but address different aspects of reliability.

875

876 From a clinical perspective, the variability in reliability metrics suggests that while
877 remote technologies show promise, their consistency is not yet sufficient to fully
878 replace standard measurement methods in all contexts. Clinicians should therefore
879 interpret remote ROM measurements with caution, particularly when monitoring
880 small changes over time or when multiple assessors are involved.

881

882 The QUADAS-2 assessment highlights several important limitations in the current
883 evidence base. First, the predominance of healthy participants introduces a key
884 source of bias and limits the applicability of findings to clinical populations. This
885 interpretation is further supported by the subgroup analysis, which demonstrated a
886 greater mean measurement difference in pathological shoulders (4.33°) compared to
887 healthy participants (2.37°), although this difference did not reach statistical
888 significance. While this finding should be interpreted with caution, it suggests a trend
889 towards increased measurement discrepancy in pathological populations. Although
890 remote measurement technologies may perform well under ideal conditions, their
891 accuracy and reliability in patients with restricted or compensatory movement
892 patterns remain uncertain. Second, inconsistencies in reporting or justifying the
893 reference standards used—such as goniometers, optical motion capture, or clinician
894 ratings—may have impacted the observed variability in agreement. In some studies,
895 the reference standard itself was subject to known inter-rater error, further
896 complicating interpretation. Additionally, few studies reported whether the
897 interpretation of remote and reference measures was blinded, introducing a potential
898 for bias in outcome assessment. Finally, while the flow and timing domain was
899 largely well addressed, the lack of standardisation in the interval between
900 assessments and test procedures may still contribute to unexplained variability.

901

902 This review builds on prior work examining digital ROM measurement methods but
903 uniquely focuses on remote ROM assessment. In contrast to a previous systematic
904 review by Shepherd and colleagues that found no significant differences between
905 digital measurement tools and traditional goniometer tools, our study identifies
906 potential sources of measurement error and highlights specific movement planes
907 and pathological population that are more prone to inaccuracies, thereby highlighting
908 the need for further evaluation and optimisation required in remote assessment
909 contexts.⁴⁵

910

911 While this review provides a comprehensive analysis, several limitations must be
912 acknowledged: Heterogeneity across studies in sample populations, measurement
913 techniques, and reference standards contributed to high variability ($I^2 > 90\%$ in meta-
914 analysis). The use of prediction intervals in this review partially addresses this
915 limitation by providing an estimate of the expected range of effects in future settings.
916 Several additional methodological approaches could be used to further explore this
917 heterogeneity. Meta-regression analysis may allow investigation of whether study-
918 level factors, such as population characteristics (e.g. healthy vs pathological), type of
919 technology (e.g. IMU vs camera-based systems), or measurement protocol,
920 contribute to variability in effect estimates.

921

922 Additionally, sensitivity analyses could be conducted to assess the robustness of
923 findings by excluding studies with high risk of bias, small sample sizes, or non-
924 standardised measurement approaches. This may help determine whether specific
925 studies disproportionately influence pooled estimates.

926

927 However, the feasibility of such analyses in the present review was limited by the
928 relatively small number of studies within each subgroup and variability in reported
929 outcome measures. As such, these approaches may not have yielded reliable or
930 interpretable results. Future research with more standardised reporting and larger
931 datasets would enable more advanced quantitative exploration of heterogeneity.

932

933 Limited representation of certain movement planes (e.g., extension) means that
934 findings may not be fully generalisable across all ROM measurements. Potential

935 publication bias was detected, indicating that negative or low-accuracy findings may
936 be underreported. The demographic profile of the included studies presents an
937 important limitation. Most of the evidence stems from research involving young,
938 healthy individuals, which may not reflect the performance of remote ROM tools in
939 real-world clinical populations—particularly older adults or patients with pathological
940 shoulders, such as rotator cuff tears or adhesive capsulitis. These groups may
941 present with limited mobility, compensatory movement patterns, or pain, all of which
942 can affect measurement accuracy and usability.

943

944 The findings of this review suggest that remote technologies may have a role in
945 supporting shoulder assessment, particularly in tele-rehabilitation and follow-up
946 settings, although their variability means they should be used as adjuncts rather than
947 replacements for standard clinical measures. The observed inconsistency in
948 accuracy and reliability highlights the importance of clinician training, particularly in
949 areas such as sensor positioning, patient instruction, and standardisation of
950 measurement protocols to optimise performance. While these technologies may
951 improve access to care and patient engagement, limitations in detecting small but
952 clinically meaningful changes may affect their ability to guide treatment decisions
953 and monitor subtle progress. From a health economic perspective, remote
954 assessment has the potential to reduce healthcare utilisation and associated costs,
955 although this must be balanced against the risk of reduced measurement precision
956 and the potential need for confirmatory in-person assessment. Additionally, future
957 research should address digital inequality, as access to and familiarity with digital
958 technologies may vary by age, socioeconomic status, and health literacy. Ensuring
959 that remote assessment tools are validated across diverse populations will be
960 essential to support effective implementation.

961

962 **Conclusion**

963 Remote technologies for assessing shoulder range of motion show generally small
964 differences compared to reference standards, but performance is inconsistent, with
965 substantial heterogeneity and variable reliability across studies. Although average
966 error is often below reported minimal clinically important differences (MCID), the
967 wide variability observed suggests these tools may not reliably detect small but
968 clinically meaningful changes. As such, they may have value for tracking overall
969 trends in shoulder function, but their use in clinical decision-making—particularly in
970 pathological populations—should be approached with caution.

971 **References**

972 1. Barzegar Khanghah A, Fernie G, Roshan Fekr A. Joint angle estimation during
973 shoulder abduction exercise using contactless technology. *BioMedical Engineering*
974 *OnLine* 2024;23:11. doi:10.1186/s12938-024-01203-5

975 2. Bechard L, Bell K, Lynch A. Preliminary validation of a mobile force Sensing device for
976 clinical and telerehabilitation. *J Biomech* 2020;110:109973.
977 doi:10.1016/j.jbiomech.2020.109973

978 3. Boissy P, Diop-Fallou S, Lebel K, Bernier M, Balg F, Tousignant-Laflamme Y. Trueness
979 and Minimal Detectable Change of Smartphone Inclinometer Measurements of
980 Shoulder Range of Motion. *Telemed J E Health* 2017;23:503-6.
981 doi:10.1089/tmj.2016.0205

982 4. Borresen A, Chakka K, Wu R, Lin CK, Wolfe C, Prabhakaran B, Annaswamy TM.
983 Comparison of in-person and synchronous remote musculoskeletal exam using
984 augmented reality and haptics: A pilot study. *Pm r* 2023;15:891-8.
985 doi:10.1002/pmrj.12883

986 5. Bradley KE, Cook C, Reinke EK, Vinson EN, Mather RC, III, Riboh J, et al. Comparison
987 of the accuracy of telehealth examination versus clinical examination in the
988 detection of shoulder pathology. *Journal of Shoulder and Elbow Surgery*
989 2021;30:1042-52. doi:10.1016/j.jse.2020.08.016

990 6. Campbell KJ, Louie PK, Bohl DD, Edmiston T, Mikhail C, Li J, et al. A novel, automated
991 text-messaging system is effective in patients undergoing total joint arthroplasty. *J*
992 *Bone Joint Surg Am* 2019;101:145-51. doi:10.2106/JBJS.17.01505

993 7. Chan LYT, Chua CS, Chou SM, Seah RYB, Huang Y, Luo Y, et al. Assessment of
994 shoulder range of motion using a commercially available wearable sensor-a
995 validation study. *MHealth* 2022;8:30-. doi:10.21037/mhealth-22-7

996 8. Chen YP, Lin CY, Tsai MJ, Chuang TY, Lee OK. Wearable Motion Sensor Device to
997 Facilitate Rehabilitation in Patients With Shoulder Adhesive Capsulitis: Pilot Study to
998 Assess Feasibility. *J Med Internet Res* 2020;22:e17032. doi:10.2196/17032

999 9. Constant CR, Murley AH. A clinical method of functional assessment of the shoulder.
1000 *Clin Orthop Relat Res* 1987:160-4.

1001 10. Çubukçu B, Yüzgeç U, Zileli R, Zileli A. Reliability and validity analyzes of Kinect V2
1002 based measurement system for shoulder motions. *Med Eng Phys* 2020;76:20-31.
1003 doi:10.1016/j.medengphy.2019.10.017

1004 11. Cuesta-Vargas AI, Roldan-Jimenez C. Validity and reliability of arm abduction angle
1005 measured on smartphone: a cross-sectional study. *BMC Musculoskelet Disord*
1006 2016;17:93. doi:10.1186/s12891-016-0957-3

1007 12. Cui J, Yeh SC, Lee SH. Wearable Sensors Integrated with Virtual Reality: A Self-Guided
1008 Healthcare System Measuring Shoulder Joint Mobility for Frozen Shoulder. *J Healthc*
1009 *Eng* 2019;2019:7681237. doi:10.1155/2019/7681237

1010 13. Faber M, Andersen MH, Sevel C, Thorborg K, Bandholm T, Rathleff M. The majority
1011 are not performing home-exercises correctly two weeks after their initial instruction-
1012 an assessor-blinded study. *PeerJ* 2015;3:e1102. doi:10.7717/peerj.1102

1013 14. Falbriard M, Meyer F, Mariani B, Millet GP, Aminian K. Drift-Free Foot Orientation
1014 Estimation in Running Using Wearable IMU. *Front Bioeng Biotechnol* 2020;8:65.
1015 doi:10.3389/fbioe.2020.00065

- 1016 15. Gismervik SØ, Drogset JO, Granviken F, Rø M, Leivseth G. Physical examination tests
1017 of the shoulder: a systematic review and meta-analysis of diagnostic test
1018 performance. *BMC Musculoskeletal Disorders* 2017;18:41. doi:10.1186/s12891-017-
1019 1400-0
- 1020 16. Gushikem A, Gomes Costa RR, Lima Cabral AL, Lopes Bomtempo LF, de Mendonça
1021 Cardoso M. Validity of range of motion, muscle strength, sensitivity, and Tinel sign
1022 tele-assessment in adults with traumatic brachial plexus injury. *Acta Neurochir*
1023 (Wien) 2022;164:1317-28. doi:10.1007/s00701-022-05164-3
- 1024 17. Guyatt GH, Oxman AD, Vist GE, Kunz R, Falck-Ytter Y, Alonso-Coello P, Schünemann
1025 HJ. GRADE: an emerging consensus on rating quality of evidence and strength of
1026 recommendations. *Bmj* 2008;336:924-6. doi:10.1136/bmj.39489.470347.AD
- 1027 18. Hayes K, Walton JR, Szomor ZL, Murrell GAC. Reliability of five methods for assessing
1028 shoulder range of motion. *Australian Journal of Physiotherapy* 2001;47:289-94.
1029 doi:[https://doi.org/10.1016/S0004-9514\(14\)60274-9](https://doi.org/10.1016/S0004-9514(14)60274-9)
- 1030 19. Höglund G, Grip H, Öhberg F. The importance of inertial measurement unit
1031 placement in assessing upper limb motion. *Medical Engineering & Physics* 2021;92:1-
1032 9. doi:<https://doi.org/10.1016/j.medengphy.2021.03.010>
- 1033 20. Howlett N, Trivedi D, Troop NA, Chater AM. Are physical activity interventions for
1034 healthy inactive adults effective in promoting behavior change and maintenance,
1035 and which behavior change techniques are effective? A systematic review and meta-
1036 analysis. *Transl Behav Med* 2019;9:147-57. doi:10.1093/tbm/iby010
- 1037 21. Hwang S, Ardebol J, Ghayyad K, Pak T, Bonadiman JA, Denard PJ, Menendez ME.
1038 Remote visual estimation of shoulder range of motion has generally high
1039 interobserver reliability but limited accuracy. *JSES Int* 2023;7:2528-33.
1040 doi:10.1016/j.jseint.2023.07.002
- 1041 22. Jayaraman S, Joshy J, Priya PK, Poduval M, Thangavel AB. SmartWatch as a
1042 Kinaesthetic System for Shoulder Function Assessment. *2020 IEEE International*
1043 *Conference on Pervasive Computing and Communications Workshops (PerCom*
1044 *Workshops)*. 2020:1-6.
- 1045 23. Johnson BM, Crawford Z, Harley J, Bonamer JP, Grawe BM. A comparison of
1046 quantitative and qualitative shoulder range of motion with patient-reported
1047 outcomes. *J Shoulder Elbow Surg* 2023;32:1364-9. doi:10.1016/j.jse.2023.02.133
- 1048 24. Kane LT, Thakar O, Jamgochian G, Lazarus MD, Abboud JA, Namdari S, Horneff JG.
1049 The role of telehealth as a platform for postoperative visits following rotator cuff
1050 repair: a prospective, randomized controlled trial. *Journal of Shoulder and Elbow*
1051 *Surgery* 2020;29:775-83. doi:10.1016/j.jse.2019.12.004
- 1052 25. Knapp PW, Keller RA, Mabee KA, Pillai R, Frisch NB. Quantifying Patient Engagement
1053 in Total Joint Arthroplasty Using Digital Application-Based Technology. *J Arthroplasty*
1054 2021;36:3108-17. doi:10.1016/j.arth.2021.04.022
- 1055 26. Kolber MJ, Hanney WJ. The reliability and concurrent validity of shoulder mobility
1056 measurements using a digital inclinometer and goniometer: a technical report. *Int J*
1057 *Sports Phys Ther* 2012;7:306-13.
- 1058 27. Lam WWT, Fong KNK. Validity and Reliability of Upper Limb Kinematic Assessment
1059 Using a Markerless Motion Capture (MMC) System: A Pilot Study. *Arch Phys Med*
1060 *Rehabil* 2024;105:673-81.e2. doi:10.1016/j.apmr.2023.10.018
- 1061 28. Mitchell K, Gutierrez SB, Sutton S, Morton S, Morgenthaler A. Reliability and validity
1062 of goniometric iPhone applications for the assessment of active shoulder external

- 1063 rotation. *Physiother Theory Pract* 2014;30:521-5.
1064 doi:10.3109/09593985.2014.900593
- 1065 29. Muir SW, Corea CL, Beaupre L. Evaluating change in clinical status: reliability and
1066 measures of agreement for the assessment of glenohumeral range of motion. *N Am J*
1067 *Sports Phys Ther* 2010;5:98-110.
- 1068 30. Niu X, Zou K, Shen D, Drew S, Wu S, Guo G, Chen R. UltraMotion: High-Precision
1069 Ultrasonic Arm Tracking for Real-World Exercises. *IEEE Transactions on Mobile*
1070 *Computing* 2024;23:1846-62. doi:10.1109/TMC.2023.3241077
- 1071 31. Ongvisatepaiboon K VV, Chignell M, Mekhora K, Chan JH. Smartphone-based audio-
1072 biofeedback system for shoulder joint tele-rehabilitation. *J Med Imaging Health*
1073 *Inform* 2016;6:1127—34. doi:<https://doi.org/10.1166/jmihi.2016.1810>.
- 1074 32. Padua R, de Girolamo L, Grassi A, Cucchi D. Choosing patient-reported outcome
1075 measures for shoulder pathology. *EFORT Open Rev* 2021;6:779-87.
1076 doi:10.1302/2058-5241.6.200109
- 1077 33. Pereira B, Cunha B, Viana P, Lopes M, Melo ASC, Sousa ASP. A Machine Learning App
1078 for Monitoring Physical Therapy at Home. *Sensors (Basel)* 2023;24.
1079 doi:10.3390/s24010158
- 1080 34. Rajkumar A, Vulpi F, Bethi SR, Raghavan P, Kapila V. Usability study of wearable
1081 inertial sensors for exergames (WISE) for movement assessment and exercise.
1082 *Mhealth* 2021;7:4. doi:10.21037/mhealth-19-199
- 1083 35. Ramkumar PN, Haeberle HS, Ramanathan D, Cantrell WA, Navarro SM, Mont MA, et
1084 al. Remote patient monitoring using mobile health for total knee arthroplasty:
1085 Validation of a wearable and machine learning-based surveillance platform. *J*
1086 *Arthroplasty* 2019;34:2253-9. doi:10.1016/j.arth.2019.05.021
- 1087 36. Rangan A, Handoll H, Brealey S, Jefferson L, Keding A, Martin BC, et al. Surgical vs
1088 nonsurgical treatment of adults with displaced fractures of the proximal humerus:
1089 the PROFHER randomized clinical trial. *JAMA* 2015;313:1037-47.
1090 doi:10.1001/jama.2015.1629
- 1091 37. Rigoni M, Gill S, Babazadeh S, Elsewaisy O, Gillies H, Nguyen N, et al. Assessment of
1092 Shoulder Range of Motion Using a Wireless Inertial Motion Capture Device-A
1093 Validation Study. *Sensors (Basel)* 2019;19. doi:10.3390/s19081781
- 1094 38. Roldán-Jiménez C, Martín-Martín J, Cuesta-Vargas AI. Reliability of a Smartphone
1095 Compared With an Inertial Sensor to Measure Shoulder Mobility: Cross-Sectional
1096 Study. *JMIR Mhealth Uhealth* 2019;7:e13640. doi:10.2196/13640
- 1097 39. Rong R, Kuo C. Dynamic Soft Tissue Artifacts during Impulsive Loads: Measurement
1098 Errors Vary With Wearable Inertial Measurement Unit Sensor Design. *IEEE Trans*
1099 *Biomed Eng* 2024;71:3275-82. doi:10.1109/tbme.2024.3416378
- 1100 40. Sahu D, Shah D, Joshi M, Shaikh S, Gaikwad P, Shyam A. Validation of an on-screen
1101 application-based measurement of shoulder range of motion over telehealth
1102 medium. *J Shoulder Elbow Surg* 2022;31:201-8. doi:10.1016/j.jse.2021.06.017
- 1103 41. Sahu D, Shah D, Joshi M, Shaikh S, Gaikwad P, Shyam A. Validation of an on-screen
1104 application-based measurement of shoulder range of motion over telehealth
1105 medium. *J Shoulder Elbow Surg* 2022;31:201-8. doi:10.1016/j.jse.2021.06.017
- 1106 42. Sallis JF, Prochaska JJ, Taylor WC. A review of correlates of physical activity of
1107 children and adolescents. *Med Sci Sports Exerc* 2000;32:963-75.
1108 doi:10.1097/00005768-200005000-00014

- 1109 43. Seo NJ, Fathi MF, Hur P, Crocher V. Modifying Kinect placement to improve upper
1110 limb joint angle measurement accuracy. *J Hand Ther* 2016;29:465-73.
1111 doi:10.1016/j.jht.2016.06.010
- 1112 44. Shah NV, Gold R, Dar QA, Diebo BG, Paulino CB, Naziri Q. Smart Technology and
1113 Orthopaedic Surgery: Current Concepts Regarding the Impact of Smartphones and
1114 Wearable Technology on Our Patients and Practice. *Curr Rev Musculoskelet Med*
1115 2021;14:378-91. doi:10.1007/s12178-021-09723-6
- 1116 45. Shepherd J, Hansjee S, Divall P, Raval P, Singh HP. How do digital range of motion
1117 measurement devices 'measure-up' to traditional goniometry in assessing shoulder
1118 range of motion? A systematic review and meta-analysis. *Shoulder Elbow*
1119 2024;16:363-81. doi:10.1177/17585732231195554
- 1120 46. Shimizu H, Saito T, Kouno C, Shimoura K, Kawabe R, Shinohara Y, et al. Validity and
1121 reliability of a smartphone application for self-measurement of active shoulder range
1122 of motion in a standing position among healthy adults. *JSES Int* 2022;6:655-9.
1123 doi:10.1016/j.jseint.2022.04.005
- 1124 47. Soeters R, Damodar D, Borman N, Jacobson K, Shi J, Pillai R, Mehran N. Accuracy of a
1125 smartphone software application compared with a handheld goniometer for
1126 measuring shoulder range of motion in asymptomatic adults. *Orthop J Sports Med*
1127 2023;11:23259671231187297-. doi:10.1177/23259671231187297
- 1128 48. Soeters R, Damodar D, Borman N, Jacobson K, Shi J, Pillai R, Mehran N. Accuracy of a
1129 Smartphone Software Application Compared With a Handheld Goniometer for
1130 Measuring Shoulder Range of Motion in Asymptomatic Adults. *Orthop J Sports Med*
1131 2023;11:23259671231187297. doi:10.1177/23259671231187297
- 1132 49. Sullivan GM. A primer on the validity of assessment instruments. *J Grad Med Educ*
1133 2011;3:119-20. doi:10.4300/jgme-d-11-00075.1
- 1134 50. Tozawa R, Ishii N, Onuma R, Kawasaki T. The reliability and validity of joint range of
1135 motion measurement using zoom and a smartphone application. *J Phys Ther Sci*
1136 2023;35:538-41. doi:10.1589/jpts.35.538
- 1137 51. Wang CH, Hwang YS, Chen YL, Chen CC, Tsai KY. Implementation of interactive games
1138 to a shoulder rehabilitation and evaluation system. *Technol Health Care*
1139 2020;28:431-7. doi:10.3233/thc-202173
- 1140 52. Wang G, Fiedler AK, Warth RJ, Bailey L, Shupe PG, Gregory JM. Reliability and
1141 accuracy of telemedicine-based shoulder examinations. *J Shoulder Elbow Surg*
1142 2022;31:e369-e75. doi:10.1016/j.jse.2022.04.005
- 1143 53. Wang XM, Smith DT, Zhu Q. A webcam-based machine learning approach for three-
1144 dimensional range of motion evaluation. *PLoS One* 2023;18:e0293178.
1145 doi:10.1371/journal.pone.0293178
- 1146 54. Werner BC, Holzgrefe RE, Griffin JW, Lyons ML, Cosgrove CT, Hart JM, Brockmeier SF.
1147 Validation of an innovative method of shoulder range-of-motion measurement using
1148 a smartphone clinometer application. *J Shoulder Elbow Surg* 2014;23:e275-82.
1149 doi:10.1016/j.jse.2014.02.030
- 1150 55. Whiting PF, Rutjes AW, Westwood ME, Mallett S, Deeks JJ, Reitsma JB, et al.
1151 QUADAS-2: a revised tool for the quality assessment of diagnostic accuracy studies.
1152 *Ann Intern Med* 2011;155:529-36. doi:10.7326/0003-4819-155-8-201110180-00009
1153
1154

1155 **Appendix 1: Full database search strategies**

1156 The following search strategies were used for each database. Search terms were
1157 adapted as required for database-specific indexing systems (e.g. MeSH terms in
1158 MEDLINE) and syntax. The full search strategy for MEDLINE is presented below,
1159 with equivalent adaptations applied across the other databases (PubMed, Web of
1160 Science, CENTRAL, and IEEE Xplore) using the same combination of keywords and
1161 Boolean operators. This approach allows replication of the search process across all
1162 databases.

1163

1164 The MEDLINE (Ovid) search strategy was implemented as follows:

Search: Ovid Medline Inception-2024

Searched 10th May 2024

Topics	Descriptor, keywords and truncations
--------	--------------------------------------

Remote

1. Remote
2. Remote Sensing Technology
3. Remote consultation
4. Teleworking
5. Telerehabilitation
6. Telemedicine
7. Online Systems
8. OR/1-7

Shoulder

9. Shoulder
10. Shoulder pain
11. Shoulder Impingement Syndrome
12. Shoulder Joint
13. Shoulder Fracture*
14. Shoulder Dislocation
15. Shoulder Arthroplasty
16. Shoulder Replacement
17. Shoulder Prosthesis
18. Shoulder Injur*
19. Brachial Plexus Neuritis
20. Rotator Cuff Tear Arthropathy

	21. Shoulder Patholog*
	22. OR/9-21
Range of motion & Strength	23. Range of motion
	24. Range-of-motion
	25. ROM
	26. Joint flexibility
	27. Kinematic*
	28. Biomechanic*
	29. Strength
	30. Muscle Strength
	31. Dynamomet*
	32. OR/23-31
Accuracy & Reliability	33. Accuracy
	34. Precision
	35. Reliability
	36. Reproducibility
	37. Validity
	38. OR/33-37
Boolean operators	s8. AND s22. AND s32. AND s38.

1165

1166

1167 Final MEDLINE search string:

1168 (Remote OR "Remote Sensing Technology" OR "Remote consultation" OR Teleworking
1169 OR Telerehabilitation OR Telemedicine OR "Online Systems") AND (Shoulder OR
1170 "Shoulder pain" OR "Shoulder Impingement Syndrome" OR "Shoulder Joint" OR Shoulder
1171 Fracture OR Shoulder Dislocation OR Shoulder Arthroplasty OR Shoulder Replacement
1172 OR Shoulder Prosthesis OR Shoulder Injur OR "Brachial Plexus Neuritis" OR "Rotator
1173 Cuff Tear Arthropathy" OR Shoulder Patholog*) AND ("Range of motion" OR Range-of-
1174 motion OR ROM OR "Joint flexibility" OR Kinematic* OR Biomechanic* OR Strength OR
1175 "Muscle Strength" OR Dynamomet*) AND (Accuracy OR Precision OR Reliability OR
1176 Reproducibility OR Validity)**

1177

1178 For other databases (PubMed, Web of Science, CENTRAL, and IEEE Xplore), equivalent
1179 search strategies were used with appropriate adjustments to database-specific syntax and
1180 indexing terms. No restrictions were placed on publication date, and all searches were
1181 conducted up to May 2024 (subsequently updated to March 2026).

1182 **Appendix 2: Signalling questions used for the QUADAS-2**
1183 **Tool**

1184

1185 **1. Patient Selection**

1186 **Risk of Bias Signalling Questions:**

- 1187 • Was a consecutive or random sample of participants enrolled?
1188 • Was a case-control design avoided?
1189 • Did the study avoid inappropriate exclusions (e.g. excluding certain ROM
1190 impairments without justification)?

1191 **Applicability Concerns:**

- 1192 • Do the included patients reflect the population in which remote ROM
1193 assessments would be used (e.g. clinical, post-op, or community settings)?
1194

1195 **2. Index Test: Remote ROM Assessment Method**

1196 **Risk of Bias Signalling Questions:**

- 1197 • Were the index test results interpreted without knowledge of the reference
1198 standard results?
1199 • Was the remote ROM assessment conducted and interpreted according to a
1200 pre-specified protocol or clearly described procedure?

1201 **Applicability Concerns:**

- 1202 • Does the remote assessment method reflect the intended use in practice (e.g.
1203 telehealth, unsupervised home setting)?
1204

1205 **3. Reference Standard**

1206 **Risk of Bias Signalling Questions:**

- 1207 • Is the reference standard (e.g. universal goniometer or 3D motion capture)
1208 likely to correctly measure the target ROM?
1209 • Were the reference standard results interpreted without knowledge of the
1210 index test results?

1211 **Applicability Concerns:**

- 1212 • Is the reference standard appropriate and reflective of best clinical practice for
1213 shoulder ROM evaluation?analysis

1214 **4. Flow and Timing**

1215 **Risk of Bias Signalling Questions:**

- 1216 • Was there an appropriate time interval between the remote and reference
1217 measurements?
- 1218 • Did all patients receive the same reference standard?
- 1219 • Were all patients included in the analysis.