

A Comprehensive Survey of Action Quality Assessment: Method and Benchmark

Kanglei Zhou^{a,b}, Ruizhi Cai^b, Liyuan Wang^a, Hubert P. H. Shum^c, Xiaohui Liang^{b,d,*}

^a*Department of Psychological and Cognitive Sciences, Tsinghua University, No. 30 Shuangqing Road, Haidian District, Beijing 100084, China*

^b*State Key Laboratory of Virtual Reality Technology and Systems, Beihang University, No. 37 Xueyuan Road, Haidian District, Beijing 100191, China*

^c*Department of Computer Science, Durham University, Stockton Rd, Durham DH1 3LE, United Kingdom*

^d*Zhongguancun Laboratory, Beijing, China*

Abstract

Action Quality Assessment (AQA) aims to automatically evaluate how well human actions are performed and has been widely applied in sports analysis, skill assessment, and healthcare. However, AQA studies are often developed under heterogeneous datasets and evaluation settings, making systematic comparison across methods difficult. To address these challenges, we present a comprehensive survey of recent advances in AQA. In particular, we propose a modality-driven hierarchical taxonomy that organizes existing methods into video-based, skeleton-based, and multi-modal approaches, and analyze the methodological evolution of representative models. We further establish a unified benchmark for representative video-based AQA methods by integrating diverse datasets and standardized evaluation protocols, enabling consistent comparison in terms of both accuracy and computational efficiency. Finally, we analyze emerging research trends, identify key challenges in current AQA research, and outline future directions ranging from near-term methodological advances to longer-term opportunities enabled by emerging AI paradigms. The project webpage is available at <https://ZhouKanglei.github.io/AQA-Survey>.

Keywords: Action Quality Assessment, Sports Scoring, Skill Assessment, Exercise Assessment

*Corresponding author.

Email addresses: zhoukanglei@tsinghua.edu.cn (Kanglei Zhou),
craaaaazy@buaa.edu.cn (Ruizhi Cai), wly19@tsinghua.org.cn (Liyuan Wang),
hubert.shum@durham.ac.uk (Hubert P. H. Shum), liang_xiaohui@buaa.edu.cn (Xiaohui Liang)

Preprint submitted to Pattern Recognition

1 **1. Introduction**

2 Action Quality Assessment (AQA) [1, 2] aims to automatically assess how well an
 3 action is performed. Unlike action recognition, which focuses on identifying different
 4 action categories, AQA evaluates performance differences within the same action cat-
 5 egory. Consequently, accurately assessing action quality requires sensitivity to subtle
 6 variations in motion dynamics and domain-specific evaluation criteria, making it signif-
 7 icantly more challenging than traditional recognition tasks. By providing an objective
 8 alternative to subjective judgments, AQA has been widely applied in domains such as
 9 sports analysis [3, 4], skill assessment [5, 6], and healthcare [7, 8]. More broadly, AQA
 10 can support intelligent perception systems by enabling quantitative evaluation of action
 11 performance, facilitating applications such as skill coaching and human–robot collabo-
 12 ration in embodied AI environments [9, 10]. Since AQA spans different domains with
 13 varying terminology, this survey uses *AQA* as the primary term, encompassing related
 14 concepts such as action scoring, skill assessment, and performance evaluation.

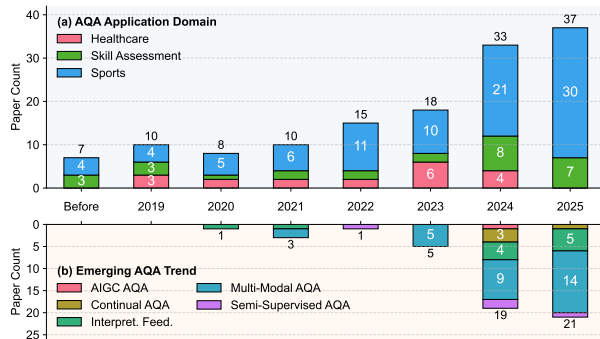


Figure 1: Annual statistics of methodological AQA papers in major CV and ML venues included in this survey. Papers are categorized by (a) application domains and (b) research directions, which have not been systematically summarized in prior AQA surveys.

15 Recent advancements in deep learning [11, 12] have driven substantial progress
 16 in AQA (see Fig. 1(a)), enabling increasingly powerful methodologies and expanding
 17 application scenarios (see Fig. 1(b)). Although several surveys [1, 2, 13, 14] have sum-
 18 marized AQA research (see Tab. 1), they remain limited in scope and methodological
 19 analysis. Early surveys primarily adopt domain-driven taxonomies [13, 14], which
 20 fail to capture shared methodological foundations across tasks, while trend-oriented
 21 reviews [2] summarize research developments but do not explicitly unify overlapping
 22 techniques across paradigms. In addition, none of the existing surveys provides a uni-
 23 fied benchmark for systematic comparison of AQA methods. Furthermore, recent sur-

Table 1: Comparison of existing AQA surveys. 🏊, 🏥, and 🎓 indicate sports, medical care, and skill assessment icons. † and ‡ denote two-level taxonomy standards. 📺, 📷, 📱, 📺, 📺, 📺, 📺, and 📺 represent video, skeleton, text, audio, gaze, and flow icons. 📁 and 📊 denote the dataset and metric icons.

Survey	Year	Domains	Taxonomy	Modality	Benchmark	Core Contributions
Lei et al. [13]	2019	🏊🏥🎓	† Feature type, ‡ Domain	📺📷	None	An early overview of motion detection and data preprocessing techniques relevant to AQA.
Wang et al. [14]	2021	🏊🏥	† Domain, ‡ Feature type	📺	None	A domain-oriented survey emphasizing categorization by application fields and publishing venues.
Liu et al. [1]	2024	🏊🏥🎓	† Input mode	📺📷	None	A descriptive review of representative AQA studies, summarizing methods and application scenarios.
Yin et al. [2]	2025	🏊🏥🎓	† Research trends	📺📷📱📺📺	None	A trend-oriented survey, highlighting research challenges with a high-level discussion of evaluation practices.
Ours	2024	🏊🏥🎓	† Input mode, ‡ Mode specialty	📺📷📱📺📺	6 📁 with 7 📊 Unified setting	A modality-driven hierarchical taxonomy, a unified benchmark, specific AQA tasks, challenges, and prospects.

1 veys [1, 2] mainly summarize existing studies but provide limited synthesis of emerging
 2 research trends and little discussion of future directions. These limitations reveal a fun-
 3 damental gap in the current literature: existing surveys mainly summarize AQA studies
 4 from domain or trend perspectives and lack a unified methodological framework that
 5 organizes modeling paradigms and supports consistent evaluation across methods.

6 As a result, AQA research remains fragmented across datasets, application do-
 7 mains, and modeling paradigms (see Fig. 1), making it difficult to obtain a systematic
 8 understanding of the field. Such fragmentation manifests at multiple levels: heteroge-
 9 neous datasets and evaluation protocols hinder reproducibility; methods are often val-
 10 idated in narrowly defined settings, limiting generalization; and the diversity of appli-
 11 cation domains obscures shared modeling principles across tasks. This fragmentation
 12 highlights the need for a more systematic and unified perspective on AQA research.

13 To bridge this gap, we present a systematic review of recent advances in AQA.
 14 **First**, we argue that the type of input data, rather than the application domain, is the
 15 primary determinant of AQA model design. Based on this premise, we propose a
 16 modality-driven hierarchical taxonomy. Specifically, we categorize approaches into
 17 video-based, skeleton-based, and multi-modal methods, further distinguishing them
 18 by modality-specific attributes. **Second**, we establish a unified AQA benchmark, to
 19 the best of our knowledge, integrating six widely used datasets and seven evaluation
 20 metrics. By standardizing experimental settings, we facilitate consistent comparisons
 21 across diverse methods, evaluating both accuracy and the often-overlooked compu-
 22 tation overhead. **Third**, beyond summarizing existing work, we analyze emerging
 23 trends and outline both near-term and longer-term directions for AQA. Near-term di-
 24 rections focus on improving reliability and interpretability, while longer-term opportu-
 25 nities arise from emerging paradigms such as generative models and embodied AI.

1 This survey targets researchers developing AQA systems in domains such as sports
 2 analytics and healthcare, as well as related areas including human–computer interac-
 3 tion, where objective action evaluation supports decision-making. The remainder of
 4 this paper is organized as follows. Sec. 2 describes the review methodology, Sec. 3
 5 introduces the foundational framework, Sec. 4 analyzes representative AQA methods,
 6 Sec. 5 reviews task-specific applications, Sec. 6 summarizes representative datasets and
 7 presents the benchmark, Sec. 7 discusses emerging trends, remaining challenges, and
 8 future research directions, and Sec. 8 concludes the whole paper.

9 2. Review Methodology

10 To ensure consistent and fair comparison across diverse AQA application domains,
 11 we follow the PRISMA framework (see Fig. 2) to improve transparency, reproducibil-
 12 ity, and reliability in selecting representative studies for this survey.

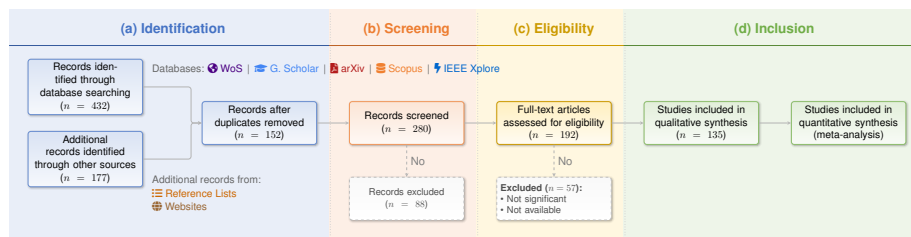


Figure 2: PRISMA flow diagram: (a) identification, (b) screening, (c) eligibility, and (d) inclusion.

13 **Identification Strategy.** We searched IEEE Xplore, Web of Science, Scopus,
 14 Google Scholar, and arXiv using keywords such as *action quality assessment*, *skill*
 15 *assessment*, and *action scoring*. Publications from 2014 to present were considered,
 16 spanning from early handcrafted to recent deep learning-based approaches. We further
 17 conducted backward citation analysis on representative papers to improve coverage.

18 **Inclusion and Exclusion Criteria.** Studies were included if they: (1) addressed
 19 AQA in terms of score, grade, or rank prediction; (2) proposed learning-based methods,
 20 benchmark datasets, or evaluation frameworks directly related to AQA; and (3) were
 21 peer-reviewed papers published in leading journals or conferences in computer vision
 22 and pattern recognition (e.g., TPAMI, IJCV, PR, CVPR, ICCV, ECCV), or widely rec-
 23 ognized arXiv preprints with distinct technical contributions. To ensure quality, we

1 consider four criteria: (1) methodological clarity, (2) experimental rigor, (3) repro-
 2 ducibility, and (4) benchmark relevance. Benchmark relevance and methodological
 3 clarity are treated as essential inclusion criteria, and studies failing to meet them are
 4 excluded. Experimental rigor and reproducibility are used for AQA and prioritiza-
 5 tion when multiple studies address similar problems. We further exclude non-English,
 6 tutorial, abstract-only, duplicate studies, and works without explicit quality prediction.

7 **Screening and Selection Process.** After deduplication, titles and abstracts were
 8 first screened to remove irrelevant records, followed by full-text assessment based on
 9 the above criteria. Papers with insufficient technical details or unclear experimental
 10 validation were excluded during full-text review. Disagreements were resolved through
 11 discussion among the authors. This resulted in 135 methodological AQA papers in-
 12 cluded in our final survey (see Fig. 2).

13 **Data Extraction and Categorization.** From each study, we extracted input modal-
 14 ities, model architectures, supervision types, datasets, and application domains. These
 15 attributes support our taxonomy and analysis across AQA methods.

16 **Reproducibility and Resources.** To promote transparency, we maintain a public
 17 webpage with curated papers and metadata, which will be continuously updated.

18 3. Fundamentals of Common AQA Setup

19 Although AQA spans diverse domains, most methods share a common framework.
 20 To enable comprehensive analysis, this section formulates a unified framework for
 21 multi-modal inputs and introduces typical scenarios and evaluation metrics.

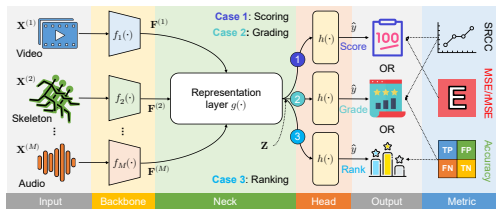


Figure 3: Three-stage AQA framework: (1) a backbone for modality-aware feature extraction; (2) a neck for representation learning that derives compact embeddings; and (3) a regression or classification head that outputs continuous scores (Case 1), discrete grades (Case 2), or ranks (Case 3).

22 3.1. Problem Formulation

23 Fig. 3 illustrates the unified framework for AQA, focusing on evaluating the qual-
 24 ity of actions through diverse input modalities, including video, skeleton data, and

Table 2: Multidimensional comparison of different AQA scenarios.

Scenarios	Outputs	Precision of Feedback	Complexity of Annotation	Application Suitability
Continuous Scoring	$\hat{y} \in [a, b]$	★★★: Delivers highly exact feedback under certain criteria	★★★: Requires precise annotations often from domain experts	Sports (e.g., gymnastics [15], diving [16], figure skating [17])
Discrete Grading	$\hat{y} \in \{1, 2, \dots\}$	★★☆: Offers output that is easier to interpret but less granular	★★☆: Demands annotations with moderate granularity	Skill assessment (e.g., surgery [18], physical exercise [7])
Discrete Ranking	$\hat{y} \in \{-1, 0, 1\}$	★☆☆: Provides general feedback with limited detail	★☆☆: Requires minimal effort	Comparative tasks (e.g., cooking [19], basketball [20], teaching [21])

1 sensor inputs. Let $\mathcal{X} = \{\mathbf{X}^{(1)}, \mathbf{X}^{(2)}, \dots, \mathbf{X}^{(M)}\}$ represent a set of $M \geq 1$ input modal-
 2 ities, where each $\mathbf{X}^{(m)} = [\mathbf{x}_1^{(m)}; \mathbf{x}_2^{(m)}; \dots; \mathbf{x}_T^{(m)}]$ corresponds to a sequence of observa-
 3 tions for modality m over time $t \in \{0, 1, \dots, T\}$. The AQA process in a neural network
 4 architecture can be divided into three primary components: feature extraction (back-
 5 bone), representation learning layer (neck), and score prediction (head). Each input
 6 modality $\mathbf{X}^{(m)}$ is processed through a feature extractor $f_m(\cdot)$, generating modality-
 7 specific features $\mathbf{F}^{(m)} = f_m(\mathbf{X}^{(m)})$. Then the extracted features are passed through
 8 the neck $g(\cdot)$, which fuses and transforms the features into an overall representation
 9 $\mathbf{Z} = g(\mathbf{F}^{(1)}, \mathbf{F}^{(2)}, \dots, \mathbf{F}^{(M)})$. Finally, the learned representation \mathbf{Z} is passed through a
 10 quality prediction head $h(\cdot)$ to estimate the predicted quality score

$$\hat{y} = h(g(f_1(\mathbf{X}^{(1)}), f_2(\mathbf{X}^{(2)}), \dots, f_M(\mathbf{X}^{(M)}))). \quad (1)$$

11 Our formulation retains generality and can handle both unimodal and multi-modal
 12 AQA methods. For unimodal methods ($M = 1$), the focus is on a single input type,
 13 such as video (see Sec. 4.2) or skeletal data (see Sec. 4.3), each of which provides
 14 full human-centric cues crucial for AQA. For multi-modal methods ($M > 1$), auxiliary
 15 modalities aid in comprehensive analysis (see Sec. 4.4).

16 3.2. Typical Scenarios

17 In AQA, the form of the output varies depending on the requirements of different
 18 applications. Broadly, AQA methods typically produce continuous or discrete outputs.

19 **3.2.1. Continuous Score.** Continuous scores are commonly used in scenarios that re-
 20 quire precise and detailed evaluation, such as in Olympic sports AQA [3, 4, 15–17].
 21 The score \hat{y} in Eq. (1) typically ranges within a specific interval, $[a, b]$ (e.g., $[0, 100]$),
 22 reflecting the exact quality of the action. The primary objective is to predict a contin-
 23 uous numerical value that accurately quantifies action quality, providing precise feed-
 24 back essential for competitive evaluation.

1 **3.2.2. Discrete Grade or Rank.** Discrete outputs, such as grade and rank, simplify
2 AQA by mapping performances to levels or relative orders, which is helpful when pre-
3 cise scores are unnecessary. **Grading** assigns each performance to a discrete level (e.g.,
4 “excellent”, “good”, “fair”, “poor”) [22, 23], providing interpretable and coarse evalu-
5 ation suited for applications like skill assessment. **Ranking** compares performances to
6 establish a relative order, useful for pairwise or group comparisons [19, 21, 24] where
7 absolute scores are less critical. Given paired samples (i, j) , the model outputs:

$$\hat{y}_{(i,j)} = \begin{cases} 1, & i \text{ outperforms } j \\ -1, & i \text{ underperforms } j \\ 0, & \text{no preference} \end{cases} \quad (2)$$

8 **3.2.3. Implications for Modeling.** As summarized in Tab. 2, different AQA scenarios
9 impose different supervision and modeling needs. *Continuous scoring* treats AQA as
10 fine-grained regression for high-stakes settings but requires precise and robust labels;
11 *discrete grading* favors efficiency and interpretability when coarse feedback suffices;
12 *ranking* relies on relative comparisons, reducing annotation cost and emphasizing or-
13 dinal consistency. These choices reflect trade-offs between precision, cost, and robust-
14 ness, motivating the joint use of correlation- and precision-based metrics in Sec. 3.3.

15 3.3. Evaluation Metrics

16 We mainly discuss metrics for evaluating prediction performance in AQA, while
17 computational efficiency metrics are considered later in the unified benchmark.

18 **3.3.1. Correlation Metrics.** The correlation metric in AQA is the Spearman Rank
19 Correlation Coefficient (SRCC). This metric assesses the strength and direction of the
20 relationship between the predicted scores \hat{y}_i ($i \in \{1, 2, \dots, N\}$) and the true scores y_i
21 based on their ranks \hat{r}_i and r_i . SRCC can be defined as:

$$\text{SRCC} = \frac{\sum_{i=1}^N (r_i - \bar{r})(\hat{r}_i - \bar{\hat{r}})}{\sqrt{\sum_{i=1}^N (r_i - \bar{r})^2} \sqrt{\sum_{i=1}^N (\hat{r}_i - \bar{\hat{r}})^2}}, \quad (3)$$

1 where \bar{r} and $\bar{\hat{r}}$ are the mean ranks of true and predicted scores, respectively. SRCC
 2 is particularly useful in scenarios where the precise numerical score is less important
 3 than the relative ranking of actions. This makes it a robust metric for evaluating ordinal
 4 data, as it emphasizes the consistency of rankings.

5 **3.3.2. Precision Metrics.** Two common types of precision metrics assess the accuracy
 6 of predicted scores, tailored to either continuous or discrete outputs (see Sec. 3.2).

7 **Score error** quantifies the difference between the predicted and actual scores, al-
 8 lowing for an assessment of the model’s accuracy in predicting continuous quality
 9 scores. Common measures include Mean Absolute Error (MAE) and Mean Squared
 10 Error (MSE), which provide insights into the magnitude of prediction errors. Recently,
 11 relative Mean Squared Error (rMSE) [25, 26] has gained popularity as it avoids the
 12 impact of differing score scales across actions in different categories. rMSE is:

$$\text{rMSE} = \frac{1}{N} \sum_{i=1}^N \left(\frac{y_i - \hat{y}_i}{y_{\max} - y_{\min}} \right)^2 \times 100, \quad (4)$$

13 where y_{\max} and y_{\min} denote the maximum and minimum scores, respectively. Note
 14 that our definition differs slightly from the relative ℓ_2 in [25]. Normalization scales the
 15 score to a decimal, and squaring reduces it further. To address this, we multiply by 100
 16 to rescale it for a precise comparison.

17 **Accuracy** measures the proportion of correctly predicted grades or ranked pairs
 18 relative to the total predictions made for applications where actions are categorized
 19 into discrete grades or ranks, which is:

$$\text{Accuracy} = \frac{\sum_{i=1}^N \mathbb{1}(\hat{y}_i = y_i)}{N} \times 100\%, \quad (5)$$

20 where $\mathbb{1}(\cdot)$ denotes the indicator function, which equals 1 if $\hat{y}_i = y_i$ and 0 otherwise.

21 **3.3.3. Task-Specific Metrics.** Beyond general AQA, task-specific metrics are often
 22 required to reflect application objectives. For AQA with temporal segmentation [27–
 23 29], Intersection over Union (IoU) evaluates segment-level localization quality, which
 24 is crucial for fine-grained feedback but less indicative of overall scoring accuracy. For
 25 continual AQA [30, 31], average forgetting measures knowledge retention over time,

1 highlighting the trade-off between stability and adaptability to new actions.

2 **3.3.4. Metric Trade-offs across Scenarios.** Different metrics capture complementary

3 aspects and are suited to different scenarios. For ranking-based AQA, correlation met-

4 rics (e.g., SRCC) are preferred as they emphasize relative ordering but may tolerate

5 large absolute errors. For scoring or grading with continuous or discrete outputs, pre-

6 cision metrics (e.g., rMSE, accuracy) are essential to measure numerical or categorical

7 correctness, though they may not fully capture ranking consistency. Task-specific met-

8 rics further tailor evaluation to specialized goals but reduce cross-task comparability.

9 Thus, many AQA settings jointly adopt correlation- and precision-based metrics.

10 4. A Modality-Driven Taxonomy of AQA Methods

11 A clear taxonomy is essential to address fragmentation in AQA methods. We or-

12 ganize AQA approaches by input modality rather than domains, as it directly deter-

13 mines network design and modeling strategy (see Sec. 4.1). Accordingly, methods are

14 grouped into video-based, skeleton-based, and multi-modal categories (see Secs. 4.2

15 to 4.4), and we discuss how modality-specific trade-offs between data acquisition cost,

16 representation capability, and robustness guide AQA system design (see Sec. 4.5).

17 4.1. Modality-Driven Taxonomy: Rationale and Overview

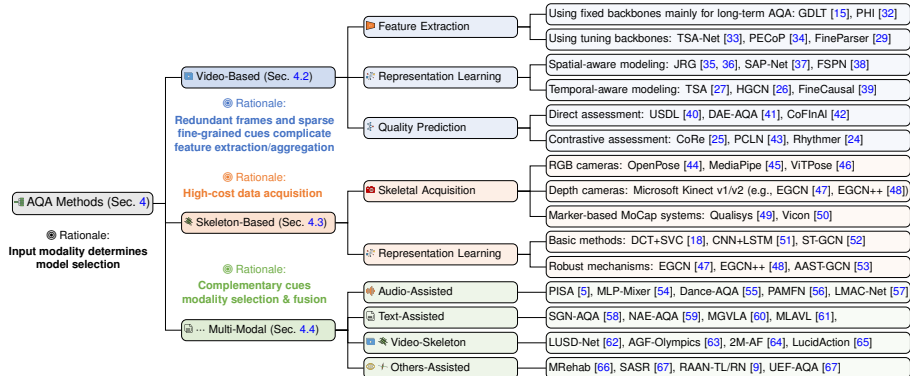


Figure 4: Our hierarchical taxonomy organized by input modalities and key method design rationales.

18 As illustrated in Fig. 4, we organize AQA methods using a modality-driven taxon-

19 omy because the input modality largely determines the available motion cues, data rep-

1 representation, and modeling pipeline. At the **first** level, methods are divided into video-
2 based, skeleton-based, and multi-modal branches. Video-based and skeleton-based
3 AQA form two primary paradigms, as they capture motion through visual appearance
4 and structured joint trajectories, respectively. At the **second** level, each branch is fur-
5 ther refined according to its core challenges. Video-based AQA focuses on extracting
6 subtle cues from long and redundant visual sequences, skeleton-based AQA empha-
7 sizes robust modeling under noisy pose data and costly acquisition, while multi-modal
8 AQA centers on effectively fusing complementary modalities. In this survey, we distin-
9 guish paradigm-level advances from incremental improvements. The former introduces
10 new modeling formulations (e.g., modality shifts or supervision paradigms), while the
11 latter mainly refines existing pipelines through architectural modifications.

12 4.2. Video-Based Approach

13 Video-based AQA methods exploit rich RGB visual information and the wide avail-
14 ability of video data to capture fine-grained motion and contextual cues, making them
15 the dominant paradigm in gymnastics [68, 69], diving [16, 70], etc. However, assessing
16 action quality from raw videos faces intrinsic challenges, including long and redundant
17 temporal sequences versus limited computational resources [4, 15], sparse and subtle
18 quality cues easily overwhelmed by background variations [71, 72], coarse features
19 transferred from action recognition with domain shift [32, 42], and limited labeled data
20 due to high annotation cost [30, 73]. These challenges arise at different stages of the
21 processing pipeline and motivate distinct design choices. Accordingly, we organize
22 video-based AQA methods by *how they address these issues across three stages*: fea-
23 ture extraction, video-level representation learning, and quality prediction.

24 **4.2.1. Feature Extraction.** As illustrated in Fig. 5, early AQA methods [6, 74] re-
25 lied on handcrafted features, which were adequate for small-scale datasets but limited
26 in capturing rich spatiotemporal cues. With the emergence of deep learning, more
27 expressive representations became possible. However, because AQA datasets remain
28 relatively small due to costly annotations, training deep models from scratch is of-
29 ten impractical for long and highly redundant video data. Consequently, most works
30 adopt backbone networks [75–79] pre-trained on large-scale image or action recogni-

1 tion datasets [75, 76, 80, 81]. Yet these models are not optimized for the fine-grained
 2 requirements of AQA, and the high computational cost of long videos further limits
 3 extensive fine-tuning. As a result, recent approaches either adapt the backbone to the
 4 AQA domain or refine representations at later stages. In the following, we review
 5 commonly used backbones and representative adaptation strategies.

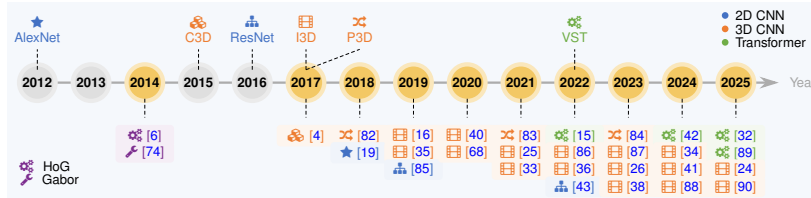


Figure 5: Timeline of backbone architectures (top) and representative video-based AQA papers (bottom).

6 **Representative Backbones in AQA.** As summarized in Tab. 3 and Fig. 5, com-
 7 monly used backbones in video-based AQA can be broadly grouped into three cate-
 8 gories: 2D CNNs, 3D CNNs, and transformer-based models, reflecting different induc-
 9 tive biases for modeling action dynamics. Early studies typically employed lightweight
 10 2D CNNs, such as AlexNet [75] (e.g., [19]) and ResNet [91] (e.g., [43, 85]), whose
 11 strong spatial inductive bias enables efficient appearance modeling but lacks explicit
 12 temporal reasoning. To better capture motion dynamics, subsequent works increas-
 13 ingly adopted 3D CNNs, including C3D [76] (e.g., [4, 70]), P3D [77] (e.g., [82]), and
 14 especially I3D [78] (e.g., [26, 29, 35, 86]), which introduce a local spatiotemporal in-
 15 ductive bias through 3D convolutions and have become the dominant choice for short-
 16 to medium-length AQA tasks. More recently, transformer-based backbones such as
 17 VST [79], adopted in [15, 42, 58, 92], relax these locality assumptions and rely on
 18 global attention to model long-range temporal dependencies, making them particularly
 19 suitable for long-term AQA, albeit at a substantially higher computational cost.

20 **Efficiency, Transferability, and Practical Choices.** The different inductive biases
 21 of 2D CNNs, 3D CNNs, and transformers lead to distinct trade-offs in efficiency and
 22 temporal modeling. Compared with 2D CNNs, 3D CNNs incur substantially higher
 23 computational and memory costs, while transformer-based models further increase this
 24 burden due to global attention. To make training tractable, many studies adopt temporal
 25 down-sampling strategies, such as key clip selection [23] or uniform clip partitioning,

Table 3: Representative backbone architectures for video-based AQA and their characteristics.

Backbone	Type	Typical Works	Pre-Trained Data	Spatial-Temporal Modeling	Efficiency	Generalization
AlexNet [75]	2D CNN	[19]	ImageNet [75]	★☆☆☆☆: Strong spatial appearance; no explicit temporal modeling.	★★★★★ Lightweight and computationally efficient	★★★☆☆ Limited cross-domain generalization
ResNet [91]	2D CNN	[43, 85]	ImageNet [75]	★★★★☆ Strong spatial features; temporal modeling requires extra modules	★★★★☆ Efficient with deeper representations	★★★★☆ Good transferability with fine-tuning
C3D [76]	3D CNN	[4, 70]	Sports1M [76]	★★★★☆ Joint spatiotemporal modeling via 3D convolutions	★★★★☆ High computation and memory cost	★★★☆☆ Moderate transferability
P3D [77]	3D CNN	[82-84]	Kinetics-400 [78]	★★★★☆ Factorized spatial and temporal convolutions	★★★★☆ More efficient than full 3D CNNs	★★★☆☆ Moderate transferability
I3D [78]	3D CNN	[26, 29, 35, 86]	Kinetics-400 [78]	★★★★★ Strong spatiotemporal modeling via inflated 2D filters	★★★★☆ Computationally heavy due to 3D operations	★★★★☆ Strong transferability with fine-tuning
VST [79]	Transformer	[15, 42, 58, 92]	Kinetics-600 [81]	★★★★★ Strong long-range temporal modeling with global attention	★★★★☆ Expensive due to self-attention	★★★★☆ Good generalization in complex scenarios

1 although this may discard subtle yet critical cues and weaken long-range temporal
 2 coherence. Meanwhile, domain shift remains a fundamental issue. Backbones pre-
 3 trained for action recognition often encode coarse semantics that are insufficient for the
 4 fine-grained requirements of AQA [32, 34, 42]. In practice, a clear trade-off emerges:
 5 2D CNNs favor efficiency, 3D CNNs balance spatiotemporal modeling and cost for
 6 moderate-length actions, while transformer-based backbones are more suitable for long
 7 and complex sequences that require long-range temporal reasoning.

8 **Backbone Tuning Strategies.** To mitigate the domain gap between pre-trained and
 9 AQA tasks, full fine-tuning is the most straightforward approach. However, it is prone
 10 to overfitting on small-scale AQA datasets [31]. Consequently, several works explore
 11 constrained backbone adaptation strategies that improve task relevance while limiting
 12 parameter updates. For example, TSA-Net [33] injects human-centric masks into I3D
 13 to emphasize critical regions with minimal overhead (see Fig. 7(a)), and this idea is
 14 further extended by FineParser [29] and Uni-FineParser [90]. In contrast, PECoP [34]
 15 introduces lightweight 3D components together with self-supervised objectives to learn
 16 in-domain spatiotemporal features, and then freezes the adapted backbone to reduce
 17 overfitting risk in downstream AQA tasks. These strategies offer a practical compro-
 18 mise between transferability and efficiency. For *long-term AQA* [15, 42, 58, 92], where
 19 minute-level videos are processed, directly fine-tuning the backbone is often impracti-
 20 cal due to prohibitive computational cost. Such scenarios are therefore more commonly
 21 addressed by feature adaptation at later stages (see Sec. 4.2.2).

22 **4.2.2. Video-Based Representation Learning.** Although powerful backbones provide
 23 strong clip-level features, feature extraction (see Sec. 4.2.1) in AQA still faces two
 24 fundamental limitations: (1) to control computational cost, long videos are usually

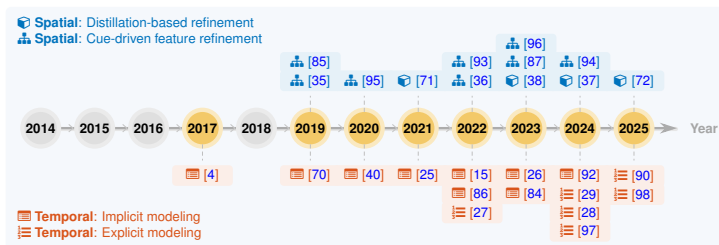


Figure 6: Timeline of representative video-based AQA methods in representation learning, grouped into spatial-aware (above the axis) and temporal-aware (below the axis) strategies over time.

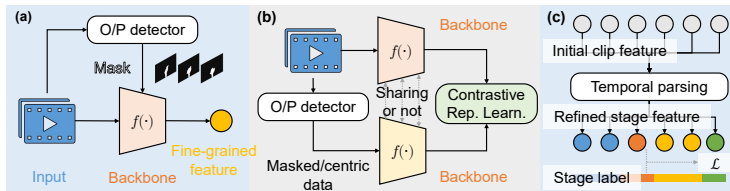


Figure 7: Three representative fine-grained reasoning paradigms in AQA. (a) Spatial reasoning injects object- or pose-centric masks into backbones to emphasize critical regions. (b) Siamese reasoning contrasts raw and masked inputs to highlight object-centric cues, sometimes using unshared branches for richer representations. (c) Temporal reasoning models procedural structure by parsing subactions, either implicitly or explicitly.

1 processed as short clips, which breaks temporal continuity and weakens long-range
 2 dependency modeling; and (2) features transferred from action recognition are often
 3 coarse and domain-shifted, and fine-tuning them is impractical for *long-term AQA*
 4 with minute-level videos due to prohibitive cost. Video-level representation learning
 5 is therefore introduced as an intermediate stage between feature extraction and quality
 6 prediction, aiming to refine clip-level features into more quality-aware representations
 7 and aggregate them into a coherent video-level embedding that preserves global tempo-
 8 ral context. By explicitly modeling finer spatial and temporal cues and integrating in-
 9 formation across clips, this stage mitigates both coarse semantics and fragmented tempo-
 10 ral modeling. Most existing methods address both aspects to some extent, but differ
 11 in their primary emphasis. As summarized in Fig. 6, representative works are reviewed
 12 below and categorized into *spatial-aware* approaches [37, 38, 71, 93], which enhance
 13 human- and object-centric cues, and *temporal-aware* approaches [70, 86, 99, 100],
 14 which focus on procedural structure and long-range dependencies.

15 **Spatial-Aware Modeling.** *Spatial-aware* modeling aims to emphasize human- and
 16 object-centric regions while suppressing background distractions, since subtle quality
 17 differences in AQA are often reflected in localized body parts, tools, or interactions
 18 rather than global appearance. Instead of explicitly modifying the backbone with masks

1 (e.g., Fig. 7(a)), most works implicitly refine spatial representations during represen-
 2 tation learning to improve flexibility (see Fig. 7(b)). These works mainly fall into two
 3 representative directions.

4 One direction is **distillation-based refinement**, which guides the model to focus on
 5 actor-centric regions through strategies such as human-masked supervision and scene-
 6 adversarial objectives [71], teacher–student consistency with pseudo actor-centric la-
 7 bels [37], background swapping augmentation [72], and causality-aware foreground–
 8 background modeling [39]. These strategies suppress background bias and transfer
 9 fine-grained spatial knowledge to the student network, enabling more discriminative
 10 representations for accurate AQA.

11 Another line of work follows **cue-driven feature refinement**, where spatial rep-
 12 resentations are enriched by integrating auxiliary cues, such as frame-wise spatial at-
 13 tention to highlight critical regions [85], semantic component modeling in surgical
 14 scenes [93], body-part graph reasoning based on pose estimation [35, 36], and inter-
 15 action relationship [87, 95] or group structure modeling [96] for multi-agent actions.
 16 These designs introduce richer semantics into spatial modeling, leading to more inter-
 17 pretable representations of critical regions, while often relying on the quality of aux-
 18 iliary cues. Cue-driven designs enrich spatial semantics but depend on external cue
 19 quality and add complexity, which may hurt robustness. In contrast, distillation-based
 20 refinement injects spatial priors implicitly, offering a lightweight and flexible alter-
 21 native. This reveals a key trade-off between explicit semantic structure and implicit,
 22 data-driven refinement, motivating designs that balance expressiveness and efficiency.

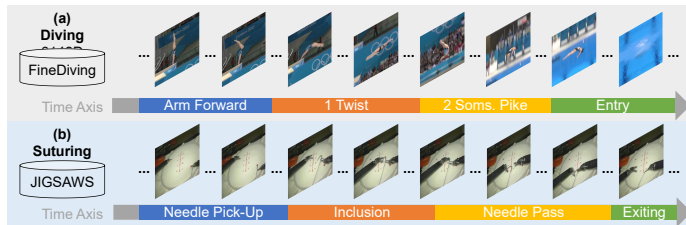


Figure 8: The procedural nature of actions in fine temporal modeling for AQA. (a) illustrates a diving example from the FineDiving dataset [27], while (b) depicts a suturing example from the JIGSAWS dataset [101].

23 **Temporal-Aware Modeling.** *Temporal-aware* modeling aims to recover the pro-
 24 cedural structure of actions (see Fig. 8) from fragmented clip features and to capture

1 long-range dependencies that are critical for AQA. In many AQA settings, an action
2 naturally consists of a sequence of subactions (e.g., approach, takeoff, flight, and entry
3 in Fig. 8(a)), and errors in any stage may substantially affect the final score [26, 27, 29,
4 83]. However, due to efficiency constraints, videos are usually processed as disjoint
5 clips, and naive aggregation (e.g., average pooling [4]) often leads to vague temporal
6 modeling and loss of contextual cues. To address this, existing works mainly fall into
7 two representative directions.

8 One line of work follows **implicit temporal modeling**, which learns temporal de-
9 pendencies directly from data using sequence models. Early studies employed LSTMs
10 or TCNs to aggregate clip features [18, 70, 102], while more recent methods adopt
11 Transformers to better capture long-range interactions [15, 86, 99, 100]. Typical de-
12 signs include grade-aware decoding via cross-attention [15], temporal parsing with
13 learnable queries [86], hierarchical graph reasoning over clips [26, 103], multi-stage
14 temporal parsing for subaction decomposition [83, 104, 105], causality-aware model-
15 ing [39, 39, 89], and event pattern analysis [22, 106]. These approaches are flexible
16 and do not require extra annotations, enabling broad applicability across tasks, but
17 often suffer from limited interpretability, higher computational cost, and the risk of
18 shortcut learning or overfitting on small datasets [92, 107].

19 In contrast, **explicit temporal modeling** injects human knowledge of action proce-
20 dures through temporal annotations. Representative works such as TSA and FineParser [27–
21 29, 90] leverage stage-level supervision to guide feature aggregation and score pre-
22 diction, leading to more interpretable and fine-grained assessment of each subaction.
23 Weakly supervised variants further reduce annotation cost by aligning procedures across
24 videos sharing similar structures [97, 98]. While these methods offer stronger inter-
25 pretability and stage-level feedback, they rely heavily on costly temporal annotations,
26 are less scalable, and often generalize poorly to unseen actions. Overall, temporal-
27 aware modeling reveals a clear trade-off: implicit approaches favor flexibility and scal-
28 ability but often lack transparency, whereas explicit designs provide structured and
29 interpretable reasoning at the cost of increased annotation demands and limited gener-
30 alization. These observations suggest that hybrid strategies, which combine data-driven
31 temporal learning with lightweight procedural priors, may offer a promising direction

1 for improving the robustness and practicality of AQA.

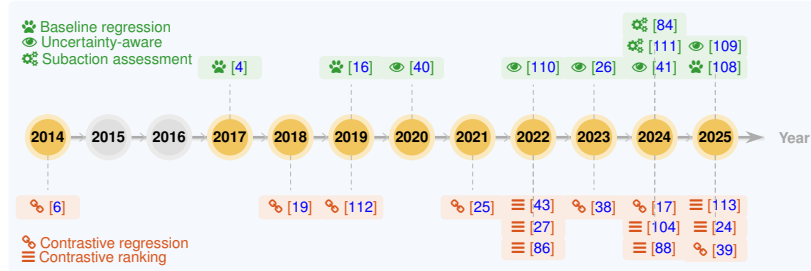


Figure 9: Timeline of representative quality prediction methods in AQA, grouped into direct assessment (top) and contrastive assessment (bottom) paradigms.

2 **4.2.3. Quality Prediction.** Quality prediction is the final stage of AQA, which maps
 3 video-level representations to assessment outputs and ultimately determines how action
 4 quality is quantified. From a modeling perspective, this stage differs mainly in whether
 5 quality is inferred from an individual action in isolation or from its relation to other
 6 actions. As visualized in Fig. 9, existing methods can therefore be broadly categorized
 7 into *direct* and *contrastive* assessments.

8 **Direct Assessment.** *Direct* assessment methods [16, 26, 30, 41, 42, 108, 114] for-
 9 mulate AQA as regression or classification, predicting absolute scores or grades by
 10 minimizing errors with respect to ground truth. They are simple to implement and
 11 computationally efficient, but often show limited sensitivity and generalization when
 12 quality differences are subtle or action contexts vary. To alleviate these issues, two
 13 representative strategies have been explored. **Uncertainty modeling** [26, 40, 41, 62]
 14 explicitly accounts for ambiguity in human judgments by predicting score distributions
 15 rather than deterministic values. USDL and MUSDL [40] soften labels with Gaussian
 16 distributions for single- and multi-judge settings, while DAE [41] employs VAE-based
 17 modeling to regress scores from latent distributions. These approaches improve robust-
 18 ness to noisy annotations, but often rely on assumptions about label variance and may
 19 generalize poorly with limited data. **Subaction assessment** [84, 94, 111, 115] decom-
 20 poses an action into meaningful stages and predicts sub-scores that are later aggregated
 21 into a final score. This enables fine-grained diagnosis and improves generalization by
 22 leveraging the internal structure of actions.

23 **Contrastive Assessment.** *Contrastive assessment* methods [19, 25, 27, 38, 86,

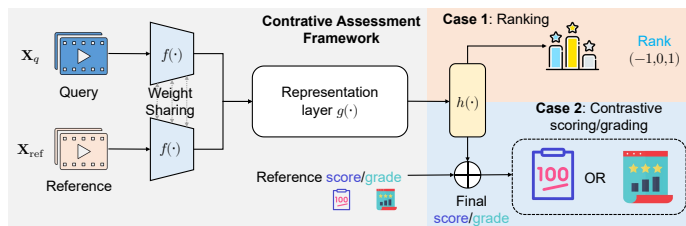


Figure 10: Pairwise contrastive framework for AQA. The model compares a target action with a reference action to learn relative quality cues. Case 1 uses relative ranking (see Eq. (2)) without exact labels, whereas Case 2 performs contrastive scoring or grading with labeled references.

1 [116](#)] learn action quality by modeling *relative differences* between actions rather than
 2 relying solely on absolute scores. Typically implemented with Siamese architectures
 3 (see Fig. 10), these methods compare a target video against reference examples, mak-
 4 ing them particularly sensitive to fine-grained performance variations. Originally intro-
 5 duced in ranking-based formulations [6, 19, 112], contrastive learning has since been
 6 extended to both ranking and scoring settings. **Ranking-oriented methods** [6, 19–
 7 21, 24] avoid absolute annotations and instead predict relative orderings between ac-
 8 tions. This formulation substantially reduces annotation cost and is especially attractive
 9 when expert scores are expensive, subjective, or unreliable. Such approaches have been
 10 successfully applied across domains, including surgical skill assessment [22, 117] and
 11 teaching quality evaluation [21, 108]. **Contrastive regression methods** incorporate
 12 pairwise contrastive learning into score prediction [25, 27–29, 37, 38, 86, 116]. Rep-
 13 resentative designs include group-aware contrastive regression [25], hierarchical con-
 14 trastive parsing of subactions [86], and multi-view or replay-based contrastive frame-
 15 works [17, 88]. By explicitly enforcing relative consistency alongside regression ob-
 16 jectives, these methods significantly improve sensitivity to subtle quality differences.

17 **Direct and Contrastive Trade-offs.** As summarized in Tab. 4, direct and con-
 18 trastive assessments differ in several ways. Direct methods efficiently regress absolute
 19 scores but rely on precise annotations and are less sensitive to subtle differences. In
 20 contrast, contrastive methods learn from relative comparisons, improving fine-grained
 21 discrimination and reducing annotation burden, albeit at higher computational cost.
 22 In practice, direct assessment suits scenarios requiring absolute scoring and real-time
 23 inference, while contrastive formulations favor fine-grained differentiation under lim-
 24 ited labels. This trade-off has motivated hybrid strategies that use relative comparisons

Table 4: Comparison of direct and contrastive assessment methods.

Assessment	Typical Works	Easy to Implement	Sensitivity to Variations	Required Data	Generalization	Computation
Direct	[40, 41, 70, 94]	★★★★★	★★★☆☆	★★★☆☆	★★★☆☆	★★★★★
		Simple	Hard to capture details	Less for train	Moderate	Efficient
Contrastive	[19, 25, 27, 37, 86, 112]	★★★☆☆	★★★★★	★★★★★	★★★★★	★★★★★
		Complex pairwise design	Robust for subtle differences	More for pairs	Better for relation focus	High for reference compute

Table 5: Statistics of popular skeleton acquisition techniques in AQA.

Data Acquisition	Input	Output	Typical Works
OpenPose [44]	Image	2D Joints	[120, 122]
MediaPipe [45]	Image	2D Joints	[123, 124]
ViTPose [46]	Image	2D Joints	[125]
CPM [126]	Image	2D Joints	[127]
MMPose [128]	Image	3D Joints	[129]
LTHP [130]	Multi-view Images	3D Joints	[131]
MotioAGF [132]	2D Joints	3D Joints	[125]
PosePrior [133]	2D Joints	3D Joints	[127]
VideoPose3D [134]	Video	3D Joints	[53, 135, 136]
MubyNet [137]	Image	3D Pose & Shape	[138]
Depth Sensors (e.g., Kinect)	Depth	3D Joints	[7, 47–51, 139–143]
Marker-Based (MoCap)	Human Body	3D joints	[49, 50]

1 during training and absolute score prediction at inference, combining efficient global
 2 calibration with fine-grained supervision [108, 118].

3 4.3. Skeleton-Based Approach

4 Skeleton-based AQA methods [49, 119–121] assess action quality by modeling hu-
 5 man joint trajectories in 2D or 3D space, explicitly focusing on body kinematics rather
 6 than visual appearance. Compared with video-based AQA, skeleton-based approaches
 7 inherently suppress background variations, making them well-suited for human-centric
 8 scenarios with limited object interaction, such as rehabilitation. However, this abstrac-
 9 tion introduces several practical challenges [33, 74]. First, reliable pose acquisition is
 10 difficult in many real-world settings due to occlusions, viewpoint changes, and lim-
 11 ited sensing conditions. Second, skeletal representations are inherently sparse and may
 12 miss subtle motion cues required for fine-grained quality assessment. Third, skeleton-
 13 only representations lack environmental context, which is often important for evaluat-
 14 ing action quality when interactions with objects or scene elements are involved. Since
 15 quality prediction largely overlaps with video-based AQA, we focus this section on
 16 skeletal acquisition and representation learning, which most clearly reflect the unique
 17 challenges of skeleton-based AQA. Given the relatively limited number of skeleton-
 18 based studies, this section emphasizes their distinctive modeling challenges rather than
 19 providing an exhaustive methodological breakdown.

1 **4.3.1. Skeletal Acquisition Methods.** High-quality skeletal data significantly affects
2 model performance. Acquisition methods include RGB cameras, depth sensors, and
3 optical motion capture systems, each with distinct trade-offs (Tab. 5). **RGB camera-**
4 **based methods** [53, 122, 124, 125, 138] estimate 2D/3D poses from video. They are
5 low-cost and accessible, but face challenges such as occlusion, depth ambiguity, and
6 pose inaccuracies. Improvements like 2D-to-3D lifting [125], multi-view reconstruc-
7 tion [138], and context integration have expanded their use in sports, healthcare, and
8 fitness. **Depth sensors** (e.g., Kinect) [50, 51, 142] directly capture 3D joints with high
9 accuracy, making them suitable for medical applications like physiotherapy [47]. How-
10 ever, they are more costly and sensitive to the environment. Combining RGB and depth
11 (RGB-D) leverages both modalities for enhanced pose estimation, particularly in reha-
12 bilitation. **Optical motion capture (MoCap) systems** (e.g., Vicon, Qualisys) [49, 50]
13 provide the highest precision for skeletal tracking, but require complex, expensive se-
14 tups and are mainly used in controlled labs. RGB cameras are affordable and flexible,
15 but less precise. Depth sensors improve 3D accuracy but have environmental and cost
16 constraints. MoCap offers unmatched accuracy at the expense of flexibility and cost.
17 The choice depends on the balance between accuracy and practicality.

18 **4.3.2. Skeletal Representation Learning Methods.** Unlike video-based representa-
19 tion learning, which focuses on pixel-based features from raw images or frames, skele-
20 tal representation learning directly works with abstract body joint coordinates and their
21 dynamic interactions over time. The complexity of joint relationships makes skeletal
22 representation learning distinct, as it emphasizes the spatial and temporal dependencies
23 between body parts. This section first reviews typical feature representation methods
24 and then introduces robust mechanisms for enhanced AQA.

25 **Basic Learning Methods.** Skeletal representation learning in AQA has evolved
26 from handcrafted descriptors to deep neural models. Early approaches relied on man-
27 ually designed features combined with traditional classifiers, such as DCT-based de-
28 scriptors or self-similarity features with SVMs [122], but these methods struggled
29 to capture complex kinematic patterns and generalized poorly. With the advent of
30 deep learning, CNNs, LSTMs, GCNs, and Transformers have been explored, among

1 which graph convolutional networks (GCNs) have become dominant due to their natu-
 2 ral compatibility with skeletal data [50, 51, 119]. By representing joints as nodes and
 3 spatial–temporal connections as edges, GCNs provide a principled way to model struc-
 4 tured body-part relations [139, 140]. The introduction of ST-GCN further enabled joint
 5 spatial–temporal modeling by integrating temporal convolution [52], laying the foun-
 6 dation for modern skeleton-based AQA. Subsequent works build upon this paradigm
 7 with enhanced strategies to improve assessment robustness [47, 48, 141, 142].

8 **Robust Learning Mechanisms.** Beyond backbone architectures, recent advances
 9 focus on improving robustness and discriminability under noisy, sparse, and variable-
 10 length skeletal data. We summarize four representative mechanisms. *Hierarchical*
 11 *body-part modeling* exploits structured priors by decomposing the skeleton into body
 12 parts or motion hierarchies [49, 51, 120, 140]. By modeling joints, limbs, and body
 13 parts at different granularities, these methods enhance sensitivity to localized motion
 14 quality while preserving global coordination. *Multi-input and ensemble modeling* im-
 15 proves robustness by integrating complementary skeletal cues. EGCN and its exten-
 16 sions [47, 48] combine positional and orientational features through data- and model-
 17 level fusion, enabling richer motion representations. Related works further incorpo-
 18 rate multi-task learning, jointly optimizing quality assessment with auxiliary objectives
 19 such as abnormality detection or reference comparison [50, 139]. *Contrastive and re-*
 20 *lational learning* introduces relative reasoning into skeletal representation learning. By
 21 comparing test actions with reference or standard motions, contrastive objectives en-
 22 hance discrimination between subtle quality differences [141], complementing abso-
 23 lute regression-based learning. *Self-supervised learning* addresses annotation scarcity
 24 and domain variability by encouraging consistent motion representations without ex-
 25 plicit labels [142, 144]. These approaches improve generalization by learning intrinsic
 26 motion patterns before or alongside supervised quality prediction.

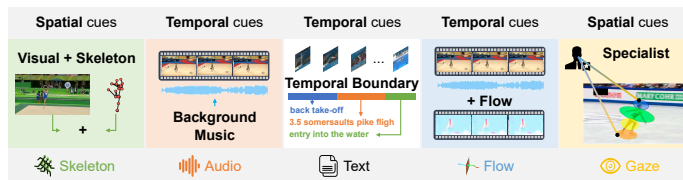


Figure 11: Illustration of typical multi-modal AQA methods.

1 4.4. Multi-Modal Approach

2 Multi-modal AQA methods improve robustness and assessment accuracy by inte-
3 grating complementary information from multiple data sources, alleviating the intrinsic
4 limitations of unimodal approaches. For example, RGB-based methods are sensitive to
5 lighting and viewpoint changes, while pose or sensor signals may be noisy or incom-
6 plete. By fusing modalities such as audio, text, skeletons, optical flow, and wearable
7 sensors, multi-modal AQA provides a more comprehensive understanding of action
8 quality (see Fig. 11). Such approaches are particularly effective in domains requiring
9 fine-grained and reliable assessment, including sports analysis [54–56, 58–60, 145] and
10 healthcare [66, 67, 146–148]. We organize existing methods by auxiliary modality and
11 discuss how each modality contributes to AQA.

12 **Audio-Assisted Methods.** Audio provides complementary temporal and rhythmic
13 cues that are difficult to infer from visual motion alone, particularly in activities in
14 which action quality is closely tied to timing, pace, or synchronization with sound.
15 Early audio-assisted AQA studies therefore focused on rhythm-dominated domains
16 such as figure skating [54], dance [55], piano performance [5], and long-term assess-
17 ment [56, 57]. Skating-Mixer [54] jointly models audio and visual streams with mem-
18 ory recurrent units to capture long-term temporal dependencies in skating routines,
19 while contrastive audio–visual learning was introduced to align dance movements with
20 musical structure and beat dynamics [55]. Subsequent works extended audio–visual
21 fusion to richer multi-modal settings. PAMFN [56] progressively aggregates RGB, op-
22 tical flow, and audio features through modality-specific and shared branches to enhance
23 fine-grained quality estimation. To handle long-term sequences, LMAC-Net [57] em-
24 ploys attention-driven fusion to align RGB, optical flow, and audio, improving tempo-
25 ral coherence over extended actions. More recent methods further incorporate semantic
26 guidance: MLAVL [61] injects language prompts to align audio–visual grading with
27 textual criteria, while MCMoE [149] uses mixture-of-experts to complete missing au-
28 dio or visual modalities. These developments demonstrate that audio cues consistently
29 enhance rhythm awareness, temporal alignment, and robustness in AQA.

30 **Text-Assisted Methods.** Text-assisted AQA incorporates criteria instructions and

1 language feedback into visual assessment, enabling more context-aware parsing and
2 interpretable evaluation. Recent works increasingly explore vision-language interac-
3 tion and prompt-based assessment, reflecting a broader trend of integrating language
4 supervision into visual reasoning. One line of work uses textual information to en-
5 hance visual representation learning and improve assessment performance. VLAKL
6 [60] embeds domain-specific terminology into multi-level representations to address
7 coarse action segmentation, while SGN [58] adopts a teacher-student framework to
8 transfer semantic knowledge from text to visual features. Another line of work focuses
9 on generating language-based feedback to improve interpretability. NAE-AQA [59]
10 reformulates score regression as video-text matching using prompt-guided transform-
11 ers, producing textual descriptions rather than numerical scores to improve readability.
12 ExpertAF [150] generates expert-like free-form feedback from RGB video and 3D
13 pose conditioned on text prompts. CROSSTRAINER [151] mines transferable skill
14 attributes from video-language pairs. EFA [152] introduces chain-of-thought reason-
15 ing for video-language form assessment, while QGVL [153] leverages quality-aware
16 prompts to guide scoring. These studies highlight the growing role of language in
17 injecting high-level semantics and interpretability into AQA systems.

18 **Video-Skeleton Methods.** Video-skeleton methods combine visual appearance
19 with explicit human kinematics, capturing both contextual cues and precise motion
20 dynamics. Such fusion is particularly beneficial when appearance alone is ambiguous
21 or when fine-grained body control is critical. Early two-stream designs integrate RGB
22 and pose features to improve robustness [154]. Subsequent works focus on long-term
23 and complex actions. Zahan et al. [63] introduce the AGF-Olympics dataset and a
24 discriminative attention module for long sports routines. 2M-AF [64] fuses RGB and
25 skeleton streams via preference-based feature aggregation. LucidAction [65] provides
26 multi-view RGB with 2D and 3D poses and adapts multiple AQA architectures. Re-
27 cent advances further refine pose-guided fusion. MELO [118] introduces scene-aware
28 contrastive fusion for multi-person routines. PGT [155] injects 2D and 3D pose priors
29 into transformers. HP-MCoRe [105] adopts hierarchical pose-guided contrastive re-
30 gression. DanceFix [156] demonstrates the effectiveness of RGB-pose fusion for ski-
31 ing and group dance assessment. These methods show that video-skeleton integration

1 substantially improves fine-grained action understanding.

2 **Other Modality-Assisted Methods.** Beyond audio, text, and skeletons, several
3 works incorporate additional sensing modalities to enhance AQA. Wearable and envi-
4 ronmental sensors provide complementary kinematic and physiological signals. MMW-
5 AQA [146] combines video, IMU, and GPS for windsurfing assessment. Gaze infor-
6 mation is fused with kinematics to analyze figure skating jumps [147]. Optical flow is
7 integrated with RGB using cross-attention for fine-grained evaluation [67]. UJ-AQA-
8 CricketVision [157] trains a dual-stream autoencoder over RGB and keypoints. FLEX
9 [148] further expands modality coverage by providing RGB, 3D pose, sEMG, and
10 physiological signals for fitness benchmarking. These studies illustrate the potential of
11 sensor-rich settings for precise and robust AQA.

12 4.5. Modality Selection and Trade-offs

13 As summarized in Tab. 6, the choice of modality in AQA reflects a fundamen-
14 tal trade-off between information richness, annotation cost, robustness, and computa-
15 tional efficiency. Video-based methods provide the most comprehensive visual cues
16 and remain the dominant choice for complex actions, but they are sensitive to back-
17 ground variations and incur high computation and annotation cost. Skeleton-based
18 approaches abstract actions into joint trajectories, offering efficiency and interpretabil-
19 ity for human-centric scenarios, yet they are vulnerable to pose estimation errors and
20 often lack contextual cues needed for fine-grained assessment. Multi-modal methods
21 further enhance robustness by integrating complementary signals, but their reliance
22 on synchronized data acquisition, modality alignment, and full-modality availability at
23 inference significantly limits scalability and real-world deployment.

Table 6: Comparison of video-based, skeleton-based, and multi-modal AQA approaches from the perspective of supervision, robustness, efficiency, and application scenarios.

Modality	Information Richness	Annotation Cost	Robustness	Computational Cost	Interpretability	Suitable Scenarios
Video-based (RGB)	High	High	Medium	High	Low	Sports judging, in-the-wild videos
Skeleton-based	Medium	Medium	Medium-Low	Low	High	Rehabilitation, technique analysis
Multi-modal	Very High	Very High	High	Very High	Medium-High	Medical assessment, expert analysis

24 In practice, modality selection should be guided by application constraints rather
25 than performance alone. Skeleton-based methods are preferable when precise kine-
26 matic analysis and efficiency are prioritized, video-based approaches suit scenarios
27 where visual context is essential, and multi-modal designs are most appropriate when

1 annotation cost and inference complexity are acceptable in exchange for maximum ro-
2 bustness. An emerging direction is to design modality-flexible frameworks that can
3 selectively activate auxiliary signals, balancing performance and practicality.

4 **5. Task-Specific Applications**

5 Despite strong benchmark performance, AQA still faces real-world challenges in
6 label scarcity, continual adaptation, and interpretability. To address these challenges,
7 we identify several underexplored research directions for practical AQA deployment.

8 **Semi-Supervised AQA.** *Labeled data scarcity* remains a fundamental bottleneck
9 for AQA, as high-quality annotations often require expert knowledge and are costly to
10 obtain. Semi-supervised AQA methods [37, 73, 158, 159] aim to alleviate this issue
11 by exploiting large amounts of unlabeled data together with limited labeled samples,
12 using techniques such as self-training, consistency regularization, and pseudo-labeling.
13 Early work such as S⁴AQA [73] explored masked segment recovery to learn robust rep-
14 resentations from unlabeled videos, while TRS-AQA [158] adopted a teacher–student
15 paradigm to generate reliable pseudo labels. These studies indicate that semi-supervised
16 learning is a promising direction for scaling AQA to new domains, though ensuring
17 pseudo-label reliability remains an open challenge.

18 **Continual AQA.** *Non-stationary variations* in action execution, recording con-
19 ditions, and user populations pose significant challenges to the robustness of AQA
20 systems. Continual AQA [30, 114] investigates how models can incrementally adapt
21 to new data distributions without catastrophic forgetting, drawing on advances in con-
22 tinual learning [11]. Representative methods include PECoP [34], which mitigates do-
23 main shift during incremental training, methods that explicitly address the introduction
24 of new action categories [114], and MAGR [30, 31], which leverages manifold pro-
25 jection and graph regularization to stabilize learning while preserving privacy. While
26 these methods show encouraging progress, balancing adaptability, memory efficiency,
27 and long-term stability remains an open research problem for practical AQA systems.

28 **Interpretable AQA.** *Limited interpretability* is also a major barrier to the real-
29 world adoption of AQA, as numerical scores alone provide little actionable guidance

Table 7: Overview of popular AQA datasets, including modality, domain, number of classes, sample size, average frames, annotations, and access URLs. Datasets span sports, skill assessment, and healthcare.

Dataset	Year	Modality	Domain	# Class	# Samples	# Average Frames	Annotations	URL
MIT Olympic [74]	2014	Video 2D Skeleton	Sports	2	309	Dive: 150 Figure Skate: 4200	Score	↗
UNLV Olympic [4]	2017	Video	Sports	3	717	Dive: 150 Figure Skate: 4200 Vault: 75	Score	↗
AQA-7 [70]	2019	Video	Sports	7	1189	Dive: 97 105 156 Vault: 87 Big Air: 122 132 Trampoline: 634	Score	↗
MTL-AQA [16]	2019	Video	Sports	16	1412	96	Score	↗
Fis-V [3]	2019	Video	Sports	1	500	4300	TES PCS	↗
RG [68]	2020	Video	Sports	1	250	2375	Difficulty Execution	↗
FineDiving [27]	2022	Video	Sports	52	3000	105	Total Step Score	↗
TaiChi-24 [154]	2022	RGB-D 3D Skeleton	Sports	24	1408	72-695	Score	-
FS1000 [54]	2023	Video Audio	Sports	7	1604	5000	TES PCS Detailed PCS	↗
FineFS [62]	2023	Video 2D/3D Skeleton	Sports	4	1167	5000	Detailed Score Subaction Class Segmentation	↗
LOGO [96]	2023	Video	Sports	12	200	5100	Action Class Formation	↗
LucidAction [65]	2024	Video 2D/3D Skeleton	Sports	259	6702	103	Score Penalty	-
FineDiving-HM [29]	2024	Video Mask	Sports	52	3000	105	Mask + Score	↗
Trampoline-AQA [69]	2025	Video Optical Flow	Sports	1	206	+902 (30.05 s @ 30fps)	Score Label Competition Label Meta Label	↗
UI-AQA-CricketVision [157]	2025	Video 2D Skeleton	Sports	1	8500	28	Phase Scores	↗
FineDiving-Pose [105]	2025	Video 2D/3D Skeleton	Sports	1	3,000	96	Pose + Score	↗
BDJ [109]	2025	Video	Sports	1	920	400	Action Specification Score Execution Level Score Total Score	↗
FLEX [148]	2025	Video 3D Pose sEMG Physiological	Sports	20	7500	234	Skill Score Knowledge Graph	↗
CoT-AFA [152]	2025	Video Text	Sports	141	3,392	108	Form + CoT Labels	↗
JIGSAWS [101]	2014	Video	Skill Assessment	3	103	-	Surgeries Class Rating	↗
EPIC-Skills [19]	2018	Video	Skill Assessment	7	216	-	Relative Rank	↗
BEST [112]	2019	Video	Skill Assessment	5	500	6400	Relative Rank	↗
PISA [5]	2021	Video Audio	Skill Assessment	1	992	160	Skill Level Difficulty	↗
TAQR [21]	2024	Video	Skill Assessment	4	300	488	Relative Rank	↗
EgoExo-Fitness [161]	2024	Video (Ego/Exo)	Skill Assessment	12	6,131	600 10-30 s @ 30fps	Score Boundaries Comment	↗
EgoExoLearn [9]	2024	Video	Skill Assessment	8	3304	250	Relative Rank	↗
BASKET [117]	2025	Video	Skill Assessment	21	66,000	16,200 (~9 min @ 30fps)	Score	↗
TAQA [162]	2026	Video Text	Skill Assessment	4	3738	348	Score Text Prompt	↗
UI-PRMD [163]	2018	3D Skeleton Joint Pos. & Ori.	Healthcare	10	1326	-	Binary Class	↗
KIMORE [7]	2019	Video, 3D Skeleton Joint Pos. & Ori.	Healthcare	5	1560	-	Score	↗
EHE [50]	2021	3D Skeleton Joint Pos. & Ori.	Healthcare	6	869	-	Binary Class	↗
FineRehab [143]	2024	Video, 3D Skeleton Joint Pos. & Ori.	Healthcare	16	4215	-	Score	↗
GAIA [164]	2024	Video	AIGC	510	9180	70	Subject Completeness Interaction	↗
Human-AGVQA [165]	2025	Video Text	AIGC	44	6000	32	Quality Ratings	↗

1 to users. Interpretable AQA methods aim to augment predictions with understand-
2 able feedback. Narrative feedback approaches [59, 94, 138, 150, 160] usually gen-
3 erate descriptive explanations, error diagnoses, or coaching-style comments through
4 vision-language foundation models. Systems such as AIFit [138], NSAQA [94], NAE-
5 AQA [59], ExpertAF [150], and TechCoach [160] illustrate different levels of inter-
6 pretability, ranging from keypoint-level analysis to free-form textual feedback. Com-
7plementary visual feedback mechanisms [8, 138, 150], including motion overlays and
8 side-by-side comparisons, further help users understand performance differences. De-
9 spite this progress, generating informative explanations remains an open challenge.

1 **6. Dataset and Benchmark**

2 We review representative datasets and observe that video-based AQA is the most
3 common setting. Thus, we use it as a case study to build a comprehensive benchmark.

4 *6.1. Datasets*

5 As shown in Tab. 7, existing AQA datasets exhibit three evolving trends.

6 First, there is a growing emphasis on **modality diversity and fine-grained anno-**
7 **tation**, moving beyond single-stream RGB videos toward richer supervision signals.
8 Early datasets such as MIT Olympic [74], UNLV Olympic [4], and AQA-7 [70] mainly
9 provide video-level scores, whereas more recent benchmarks introduce multi-modal in-
10 puts and structured labels, including step-level annotations in FineDiving [27], detailed
11 subaction scores and segmentation in FineFS [62], and pose-enhanced supervision in
12 FineDiving-Pose [105]. Healthcare-oriented datasets such as FineRehab [143] and KI-
13 MORE [7] further incorporate 3D skeletal signals to support precise AQA.

14 Second, datasets are rapidly expanding in **scale and temporal coverage**, enabling
15 more robust training and evaluation of data-hungry models. While earlier benchmarks
16 are limited to a few hundred samples, large-scale datasets such as FS1000 [54], BAS-
17 KET [117], and GAIA [164] provide thousands to tens of thousands of videos, often
18 with long temporal durations. Recent long-sequence datasets, including LOGO [96]
19 and Trampoline-AQA [69], explicitly target minute-level actions, posing new chal-
20 lenges for temporal modeling and efficiency.

21 Third, AQA datasets are expanding into **broader application domains**, reflect-
22 ing the practical diversification of AQA research. AQA datasets have steadily ex-
23 panded into more diverse application domains: early extensions focused on skill assess-
24 ment and healthcare (e.g., EPIC-Skills [19], UI-PRMD [163], EHE [50]), while more
25 recent releases further broaden the scope to teaching assessment (e.g., TAQR [21],
26 TAQA [162]) and to emerging or fine-grained settings such as AI-generated action
27 evaluation (GAIA [164], Human-AGVQA [165]) and traditional activities with spe-
28 cialized criteria (e.g., BDJ [109]).

1 6.2. A Unified AQA Benchmark

2 Existing AQA studies are often evaluated on disparate datasets with inconsistent
3 protocols, hindering fair comparisons and obscuring real progress across methods. To
4 mitigate this issue, we construct a unified benchmark centered on *video-based AQA*,
5 which remains the most mature, widely adopted, and reproducible research setting
6 in the literature. Skeleton-based and multi-modal approaches are not included at this
7 stage, primarily due to the limited availability of standardized datasets and open-source
8 implementations that support controlled benchmarking.

9 **6.2.1. Experimental Setting.** Our benchmark covers six widely used datasets, seven
10 state-of-the-art baseline methods, and seven evaluation metrics.

11 **Benchmark Datasets.** We select datasets according to three criteria: (1) public
12 availability with clearly defined train/test splits, (2) coverage of both short-term and
13 long-term actions, and (3) compatibility with standard video-based pipelines without
14 requiring additional sensors or annotations. Accordingly, we include short-term bench-
15 marks such as MTL-AQA [16], AQA-7 [70], and FineDiving [27], as well as long-term
16 datasets including Fis-V [3], RG [68], and LOGO [96].

17 **Benchmark Methods.** We include representative baselines selected based on four
18 criteria: (1) publicly available implementations to ensure reproducibility; (2) competi-
19 tive performance reported on standard AQA benchmarks; (3) methodological diversity,
20 covering major paradigms in Sec. 4; and (4) compatibility with unified evaluation set-
21 tings (see Sec. 3.3) for fair comparison. Specifically, USDL [40] and DAE [41] rep-
22 resent uncertainty-aware direct regression, CoRe [25] and T²CR [88] represent con-
23 trastive regression and multi-view reasoning, GDLT [15] and CoFInAI [42] represent
24 transformer-based methods, and HGCN [26] represents structured graph reasoning.

25 **Evaluation Metrics.** We report SRCC for ranking, MSE, and rMSE for accuracy,
26 as well as computational metrics such as training time, model size, inference speed,
27 and computational complexity, providing a balanced view of accuracy and efficiency.

28 **Implementation Details.** All benchmark results are **reproduced** using publicly
29 available implementations under a unified evaluation protocol. Unless otherwise spec-
30 ified, experiments are conducted on a single RTX 3090 GPU with an I3D backbone

Table 8: Results of our video-based AQA benchmark. The training time unit is in hours. The average SRCC is calculated using Fisher-z transformation to ensure comparability across datasets. The best results are presented in **bold**, while the second-best results are underlined.

Method	Publisher	MTL-AQA [16]				AQA-7 ^{AV} [70]				FineDiving [27]				RG ^{AV} [68]				Fis-V ^{AV} [3]				LOGO [96]			
		SRCC	MSE	rMSE	Time	SRCC	MSE	rMSE	Time	SRCC	MSE	rMSE	Time	SRCC	MSE	rMSE	Time	SRCC	MSE	rMSE	Time	SRCC	MSE	rMSE	Time
MUSDL [40]	CVPR'20	0.9350	39.7753	0.3642	11.32	0.8202	268.1272	3.0767	1.57	0.8812	53.6996	0.4927	14.08	0.4897	12.0564	4.3557	0.01	0.3514	478.8537	41.1493	0.07	0.7044	368.2800	3.6828	0.02
CoRe [25]	ICCV'21	0.9519	32.9234	0.3015	46.81	0.7838	256.4000	2.8041	4.44	0.9413	24.4000	0.2440	26.62	0.7348	6.0166	2.2178	0.12	0.7255	183.9800	1.8398	0.20	0.5852	45.3446	4.7783	0.10
GDLT [15]	CVPR'22	0.9395	43.5769	0.3990	10.84	0.8057	238.9922	2.7120	1.73	0.9342	27.9100	0.2791	8.84	0.7884	6.2177	2.2767	0.04	0.7203	190.2050	1.9021	0.04	0.7453	401.5500	4.0155	0.04
HGCN [26]	TCSVT'23	0.9522	30.7432	0.2815	14.79	0.8451	235.1635	2.7642	1.53	0.9383	24.7700	0.2477	9.14	0.7121	6.8748	2.5240	0.03	0.7270	227.8600	2.2786	0.03	0.8022	369.0600	3.6906	0.04
DAE [41]	NCAA'24	0.9497	31.3355	0.2869	11.11	0.7916	307.4355	3.4763	1.99	0.9350	26.5100	0.2651	8.58	0.7412	7.1299	2.7019	0.03	0.7447	181.1600	1.8116	0.03	0.6701	356.6300	3.5663	0.02
T ² CR [88]	INFS'24	0.9529	29.8640	0.2735	49.62	0.7910	303.2644	3.4754	4.98	0.9424	23.5100	0.2351	30.56	0.6581	7.0013	2.5215	0.28	0.6988	236.1867	2.3619	1.09	0.5199	43.3161	4.5645	0.26
CoFinAI [42]	ICAI'24	0.9461	37.7907	0.3461	14.80	0.8195	249.9134	2.7769	1.57	0.9317	36.4681	0.2887	8.64	0.7534	10.8178	4.0537	0.08	0.6974	401.9300	4.0193	0.05	0.5972	38.6218	4.0698	0.05
Diving ^{QA-V} [70] Gym Vault ^{QA-V} [70] BigSki ^{QA-V} [70] BigSnow ^{QA-V} [70] Sync. 10m ^{QA-V} [70]																									
Method	Publisher	SRCC	MSE	rMSE	Time	SRCC	MSE	rMSE	Time	SRCC	MSE	rMSE	Time	SRCC	MSE	rMSE	Time	SRCC	MSE	rMSE	Time	SRCC	MSE	rMSE	Time
MUSDL [40]	CVPR'20	0.8238	129.5963	1.2960	3.21	0.7300	133.8101	2.1927	1.44	0.5369	729.6780	8.4625	1.44	0.7109	366.3390	3.6634	1.64	0.9205	126.1327	1.5508	0.85	0.8416	124.1512	1.3939	0.86
CoRe [25]	ICCV'21	0.8432	85.5388	0.8554	11.53	0.7520	124.9170	2.0470	3.41	0.6545	365.8256	4.2427	3.47	0.5328	516.1882	5.1619	4.17	0.9061	108.2117	1.3305	2.02	0.8452	337.7187	3.5469	2.03
GDLT [15]	CVPR'22	0.8425	82.4375	0.8244	3.77	0.7665	159.1891	2.6086	1.50	0.6540	397.9806	4.6157	1.56	0.5805	526.6038	5.2661	1.79	0.9078	64.1058	0.7882	0.85	0.9049	207.0555	2.1746	0.86
HGCN [26]	TCSVT'23	0.8871	101.4813	1.0148	3.16	0.7725	220.5605	3.6143	1.36	0.6701	483.7507	5.6104	1.40	0.6487	369.5956	3.6960	1.59	0.9174	96.7002	1.1890	0.86	0.9070	149.0894	1.4713	1.85
DAE [41]	NCAA'24	0.8468	99.1924	0.9919	3.27	0.7526	158.1046	2.5908	1.43	0.5439	558.6434	6.4789	1.55	0.4824	589.9491	5.8995	3.23	0.9329	160.5100	1.9735	1.27	0.8883	278.3217	2.9231	1.19
T ² CR [88]	INFS'24	0.8334	96.7862	0.9679	9.13	0.7585	187.0455	3.0651	4.93	0.5887	482.7082	5.5983	5.03	0.4790	546.7196	5.4672	5.55	0.9087	242.4116	2.9805	2.58	0.9109	265.3465	2.7868	2.71
CoFinAI [42]	ICAI'24	0.8652	150.5001	1.5050	3.17	0.7685	126.4775	2.0726	1.43	0.6119	394.2921	4.5722	1.45	0.5990	539.6546	5.3965	1.64	0.9073	46.2191	0.5683	0.85	0.9334	243.4674	2.5570	0.87
Ball ^{HT} [68] Clubs ^{HT} [68] Hoop ^{HT} [68] Ribbon ^{HT} [68] TES ^{HT} [3] PCS ^{HT-V} [3]																									
Method	Publisher	SRCC	MSE	rMSE	Time	SRCC	MSE	rMSE	Time	SRCC	MSE	rMSE	Time	SRCC	MSE	rMSE	Time	SRCC	MSE	rMSE	Time	SRCC	MSE	rMSE	Time
MUSDL [40]	CVPR'20	0.4897	433.5650	4.3557	0.01	0.4746	439.7667	4.3977	0.01	0.4911	452.2267	4.5223	0.01	0.6128	466.8100	4.6681	0.01	0.0000	7823.4400	78.2344	0.02	0.6506	407.5000	4.0753	0.02
CoRe [25]	ICCV'21	0.7510	160.1400	1.6014	0.06	0.7592	232.6600	2.3246	0.06	0.7554	211.4200	2.1148	0.12	0.7439	246.2800	2.4628	0.06	0.6578	186.3000	1.8630	0.20	0.7812	181.6600	1.8166	0.20
GDLT [15]	CVPR'22	0.7599	197.7200	1.9772	0.02	0.7365	263.7000	2.6370	0.02	0.7616	179.1000	1.7910	0.04	0.8007	202.0100	2.0201	0.02	0.6304	237.0700	2.3707	0.04	0.7911	143.3400	1.4334	0.02
HGCN [26]	TCSVT'23	0.7007	197.3500	1.9735	0.02	0.7585	246.6900	2.4669	0.02	0.7701	176.4000	1.7640	0.03	0.7603	215.9700	2.1597	0.02	0.6830	273.5100	2.7351	0.03	0.7657	182.2100	1.8221	0.03
DAE [41]	NCAA'24	0.7425	143.1000	1.4310	0.02	0.7518	402.1000	4.0210	0.02	0.7431	183.4900	1.8349	0.03	0.7399	388.5100	3.8851	0.02	0.6748	198.8100	1.9881	0.03	0.8014	163.5100	1.6351	0.03
T ² CR [88]	INFS'24	0.6483	185.0900	1.8509	0.28	0.7335	237.6300	2.3763	0.28	0.6810	284.4400	2.8444	0.28	0.7269	231.5000	2.3150	0.28	0.6896	184.9400	1.8494	1.14	0.7536	178.8300	1.7883	1.05
CoFinAI [42]	ICAI'24	0.7616	150.9593	1.5016	0.14	0.7443	6.4109	3.0076	0.12	0.7482	24.4233	9.1924	0.11	0.7593	6.4777	2.4232	0.06	0.6539	37.2062	3.1223	0.06	0.8153	17.4047	2.7803	0.12

1 pre-trained on Kinetics-400, trained for 100 epochs using the Adam optimizer (learn-
2 ing rate 1×10^{-4} , weight decay 1×10^{-4}). For datasets without released configurations,
3 we follow official settings from prior works; for RG and Fis-V, official VST features are
4 directly adopted. To ensure fairness and reproducibility, no per-dataset or per-method
5 hyperparameter tuning is performed. Results are averaged over three runs (seeds 1022,
6 1023, and 1024) on the official splits, following common AQA benchmarking practice.
7 Task-specific tuning may further improve individual methods.

8 **6.2.2. Assessment Performance.** Tab. 8 compares seven AQA methods across six
9 benchmark datasets. Results are based on standardized implementations with official
10 code for fair comparison, and the reported standard deviations indicate generally stable
11 performance across runs. Transformer-based methods, such as CoFinAI [42], perform
12 competitively on several long-sequence datasets (e.g., RG, LOGO), while GDLT [15]
13 achieves strong SRCC performance on RG and LOGO. However, MUSDL/MUSDL [40]
14 shows larger variance on challenging datasets such as Fis-V, suggesting sensitivity to
15 noise and complex temporal dependencies. Contrastive regression methods (CoRe [25],
16 T²CR [88]) perform competitively on large-scale datasets (MTL-AQA, FineDiving),

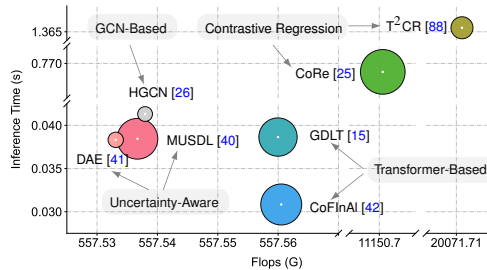


Figure 12: Computation comparison with selected baselines under the identical device settings on the MTL-AQA dataset. The x-axis represents the FLOPs, the y-axis indicates the inference time, and the bubble size corresponds to the number of parameters.

1 while their performance varies more across smaller datasets. Uncertainty-aware DAE [41]
 2 shows strong results on Fis-V, indicating the benefit of modeling prediction uncertainty.
 3 GCN-based HGCN [26] also performs competitively on AQA-7, suggesting the effec-
 4 tiveness of graph-based local and global feature modeling for robust performance.

5 **6.2.3. Computational Performance.** Tab. 8 and Fig. 12 compare training time, FLOPs,
 6 parameters, and inference speed across methods (statistics from MTL-AQA). Con-
 7 trastive regression methods (CoRe [25], T²CR [88]) are the most computationally in-
 8 tensive, with training times of 46 and 49 hours, respectively, and the longest infer-
 9 ence times (up to 1.4s per sample), mainly due to exemplar computation and ensemble
 10 strategies. T²CR’s dual time-scale design further increases computation. In contrast,
 11 direct regression and transformer-based methods (e.g., CoFiInAI [42]) are more effi-
 12 cient, with CoFiInAI showing the fastest inference (0.03s). HGCN [26] and DAE [41]
 13 have the lowest parameter counts (around 12.5M), and DAE has the lowest FLOPs
 14 (557.53G), whereas T²CR’s complexity is over 30 times higher.

15 **6.2.4. Experimental Observations.** As shown in Tab. 8 and Fig. 12, our benchmark
 16 provides several key insights into the strengths and limitations of existing AQA meth-
 17 ods. First, no single method consistently achieves the best performance across all
 18 datasets; methods that perform well on one dataset often fail to maintain the same
 19 advantage on others. Second, performance is strongly influenced by dataset charac-
 20 teristics, particularly differences in temporal length and domain structure. Third, higher
 21 predictive accuracy is often accompanied by increased computational cost, indicating
 22 a clear trade-off between performance and efficiency. These findings highlight funda-
 23 mental limitations of existing AQA approaches and emphasize the need for models that
 24 are both computationally efficient and robust across diverse datasets.

7. Trends, Challenges, and Future Directions

As shown in Fig. 13, AQA is evolving beyond benchmark-driven accuracy toward more robust and adaptive systems for real-world deployment. This section synthesizes key methodological trends, remaining challenges, and future directions in AQA.

7.1. Current Trends

Recent AQA research is shifting from simple score regression toward structure-aware modeling. A prominent trend is the adoption of **fine-grained reasoning** [27, 29, 39] to support model-level interpretability, where internal representations are aligned with meaningful temporal segments or sub-actions. In parallel, increasing attention is paid to **output-level interpretability** [150, 160], moving beyond scalar scores toward diagnostic and actionable feedback. Driven by limited annotations, AQA is increasingly integrated with **semi/self-supervised** learning [37, 158] to better exploit sparse labels, **contrastive** objectives [25, 29] to enhance assessment robustness, and **multi-modal** data [60, 61] to leverage complementary cues. At the architectural level, transformers [15, 86] and **pre-trained models** [78, 79] are becoming common backbones, enabling long-range temporal modeling and improved data efficiency. Together, these trends indicate a transition toward more interpretable and generalizable AQA systems.

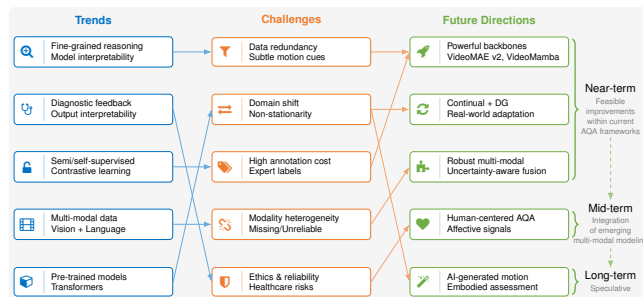


Figure 13: Evolutionary roadmap of AQA. The diagram connects current methodological **trends** (left) with emerging **challenges** (middle), and highlights potential **future directions** (right).

7.2. Under-explored Challenges

Despite recent advances, several challenges in AQA remain under-explored.

Data Redundancy versus Subtle Cues. High-dimensional video inputs contain considerable redundancy, while action quality depends on subtle motion variations. Balancing information suppression and cue amplification remains a core challenge.

1 **Domain Shift and Non-stationarity.** Gaps between pre-trained domains and AQA
2 tasks, together with real-world non-stationary variations, significantly degrade model
3 performance. Effectively adapting to such shifts remains largely unsolved.

4 **High Annotation Cost.** Fine-grained quality annotations, such as stage-level or
5 temporal labels, require expert knowledge and are expensive to obtain. Learning reli-
6 able quality representations from coarse or sparse supervision remains challenging.

7 **Modality Heterogeneity.** Although multi-modal AQA improves performance,
8 real-world data often involve missing or unreliable modalities. Designing robust mod-
9 els that gracefully degrade under incomplete multi-modal inputs remains under-explored.

10 **Ethics and Reliability.** In sensitive domains such as healthcare and rehabilitation,
11 AQA systems must ensure reliability, fairness, and interpretability. Addressing ethical
12 risks and enabling responsible deployment remain open challenges.

13 *7.3. Future Directions*

14 Building on the above challenges, we outline several promising directions for ad-
15 vancing AQA, organized roughly from near-term extensions of existing pipelines to
16 longer-term opportunities enabled by emerging AI paradigms (see Fig. 13).

17 **Integrating Powerful Foundation Models.** Our analysis in Sec. 4 reveals that
18 progress in AQA closely follows backbone evolution (see Fig. 5), and recent works [32,
19 42] show that stronger backbones consistently improve performance. Meanwhile, our
20 benchmark analysis shows that current AQA datasets are limited by label scarcity (see
21 Sec. 6), which constrains supervised learning and limits generalization. This moti-
22 vates the adoption of foundation models that leverage large-scale pre-training to learn
23 transferable representations, such as VideoMAE V2 [166] and VideoMamba [167] for
24 spatiotemporal representation, MotionGPT [168] for skeletal modeling, and Video-
25 LLaMA [169] for video-language understanding. While these models provide richer
26 representations and cross-modal reasoning capabilities, effective adaptation to AQA-
27 specific requirements remains an open challenge.

28 **Continual Learning and Domain Generalization.** Existing continual AQA stud-
29 ies [31, 114, 170] still largely rely on labeled data for model updates. Future research
30 should jointly address continual adaptation and generalization to unseen domains, par-

1 ticularly in ego-view settings [161], where viewpoints and motion patterns evolve sig-
2 nificantly over time, enabling robust AQA under non-stationary conditions. Moreover,
3 these methods rely on rehearsal strategies that store previously seen data, which may
4 raise privacy concerns and limit their applicability in real-world scenarios.

5 **Robust Multi-modal Learning.** Our analysis in Sec. 4 reveals that most multi-
6 modal AQA methods assume complete and reliable modality inputs, limiting their ap-
7 plicability in real-world scenarios where modalities may be missing or noisy. This mo-
8 tivates the design of AQA models that are robust to incomplete or unreliable modalities
9 [149, 170]. Future work should focus on flexible architectures with uncertainty-aware
10 fusion and graceful degradation under partial multi-modal inputs.

11 **Human-centered and Trustworthy AQA.** As discussed in Sec. 5 and Sec. 6, AQA
12 is often applied in sensitive domains such as healthcare [8], where reliability and inter-
13 pretability are critical. However, current methods remain vulnerable to adversarial per-
14 turbations and lack sufficient interpretability for practical deployment. This motivates
15 the development of trustworthy AQA systems that emphasize robustness, interpretabil-
16 ity, and ethical considerations. Future work should also incorporate affective signals,
17 such as emotion and engagement, to enable more human-centered assessment.

18 **AI-generated and Embodied Action Assessment.** The rapid progress of action
19 generation foundation models [168, 171] has enabled realistic motion synthesis, mak-
20 ing reliable evaluation increasingly important. In this context, AQA can play a critical
21 role in assessing the realism, correctness, and execution quality of generated motions
22 for applications such as e-commerce avatars and AI hosts [164]. Beyond motion syn-
23 thesis, similar evaluation challenges arise in embodied learning and robotics, where
24 agents must perform actions in physical or simulated environments [172, 173]. Re-
25 lated problems have also been studied in surgical skill assessment, where methods are
26 used to evaluate how humans operate robotic or medical systems [87, 174, 175]. Future
27 research should explore how AQA methods can be integrated with generative and em-
28 bodied models to provide automatic and fine-grained evaluation, potentially forming a
29 feedback loop that further improves motion generation and embodied action learning.

1 8. Conclusion

2 This survey reviews recent advances in AQA, proposes a modality-driven hierar-
3 chical taxonomy, and introduces a unified benchmark for evaluating both accuracy and
4 computational efficiency. By organizing fragmented literature and inconsistent com-
5 parisons, we clarify the methodological evolution across modalities and applications.
6 We also summarize key trends, open challenges, and future research directions. This
7 survey has several limitations. Some emerging application-specific datasets are not
8 discussed in depth due to space constraints. In addition, the benchmark focuses on
9 video-based AQA, as skeleton-based and multi-modal methods still lack standardized
10 datasets and reproducible pipelines. Thus, the evaluation emphasizes fair comparison
11 under official settings rather than extensive task-specific tuning. We hope this survey
12 serves as a useful reference and promotes more robust and unified AQA research.

13 **Acknowledgements.** This work was supported in part by NSFC under Grant 62272019.

14 References

- 15 [1] J. Liu, H. Wang, K. Stawarz et al., Vision-based human action quality assessment: A
16 systematic review, *ESWA* 263 (2025).
17 [2] H. Yin, P. Parmar, D. Xu et al., A decade of action quality assessment: Largest systematic
18 survey of trends, challenges, and future directions, *IJCV* (2026).
19 [3] C. Xu, Y. Fu, B. Zhang et al., Learning to score figure skating sport videos, *TCSVT*
20 30 (12) (2019) 4578–4590.
21 [4] P. Parmar, B. Tran Morris, Learning to score olympic events, in: *CVPRW*, 2017, pp.
22 20–28.
23 [5] P. Parmar, J. Reddy, B. Morris, Piano skills assessment, in: *MMSP*, 2021, pp. 1–5.
24 [6] Q. Zhang, B. Li, Relative hidden markov models for video-based evaluation of motion
25 skills in surgical training, *TPAMI* 37 (6) (2014) 1206–1218.
26 [7] M. Capecci, M. G. Ceravolo, F. Ferracuti et al., The kimore dataset: Kinematic assessment
27 of movement and clinical scores for remote monitoring of physical rehabilitation, *TNSRE*
28 27 (7) (2019) 1436–1448.
29 [8] K. Zhou, R. Cai, Y. Ma et al., A video-based augmented reality system for human-in-
30 the-loop muscle strength assessment of juvenile dermatomyositis, *TVCG* 29 (5) (2023)
31 2456–2466.
32 [9] Y. Huang, G. Chen, J. Xu et al., Egoexolearn: A dataset for bridging asynchronous ego-
33 and exo-centric view of procedural activities in real world, in: *CVPR*, 2024, pp. 22072–
34 22086.
35 [10] C. H. Wu, K. Ashutosh, K. Grauman, Skillsight: Efficient first-person skill assessment
36 with gaze, arXiv preprint arXiv:2511.19629 (2025).
37 [11] L. Wang, X. Zhang, H. Su et al., A comprehensive survey of continual learning: theory,
38 method and application, *TPAMI* 46 (8) (2024) 5362–5383.

- 1 [12] Z. Chang, G. A. Koulteris, H. J. Chang et al., On the design fundamentals of diffusion
2 models: A survey, *PR* (2025) 111934.
- 3 [13] Q. Lei, J.-X. Du, H.-B. Zhang et al., A survey of vision-based human action evaluation
4 methods, *Sensors* 19 (19) (2019) 4129.
- 5 [14] S. Wang, D. Yang, P. Zhai et al., A survey of video-based action quality assessment, in:
6 *INSAI, IEEE*, 2021, pp. 1–9.
- 7 [15] A. Xu, L.-A. Zeng, W.-S. Zheng, Likert scoring with grade decoupling for long-term
8 action assessment, in: *CVPR*, 2022, pp. 3222–3231.
- 9 [16] P. Parmar, B. T. Morris, What and how well you performed? a multitask learning approach
10 to action quality assessment, in: *CVPR*, 2019, pp. 304–313.
- 11 [17] Y. Liu, X. Cheng, T. Ikenaga, A figure skating jumping dataset for replay-guided action
12 quality assessment, in: *ACM MM*, 2023, pp. 2437–2445.
- 13 [18] T. Wang, Y. Wang, M. Li, Towards accurate and interpretable surgical skill assessment: A
14 video-based method incorporating recognized surgical gestures and skill levels, in: *MIC-*
15 *CAI*, 2020, pp. 668–678.
- 16 [19] H. Doughty, D. Damen, W. Mayol-Cuevas, Who’s better? who’s best? pairwise deep
17 ranking for skill determination, in: *CVPR*, 2018, pp. 6057–6066.
- 18 [20] G. Bertasius, H. Soo Park, S. X. Yu et al., Am i a baller? basketball performance assess-
19 ment from first-person videos, in: *ICCV*, 2017, pp. 2177–2185.
- 20 [21] M. Fang, X. Du, Q. Liu et al., Which is the better teacher action? a new ranking model
21 and dataset, in: *ICASSP*, 2024, pp. 7695–7699.
- 22 [22] X. Ding, X. Xu, X. Li, Sedskill: Surgical events driven method for skill assessment from
23 thoracoscopic surgical videos, in: *MICCAI*, Vol. 14228, 2023, pp. 35–45.
- 24 [23] Y. Li, X. Chai, X. Chen, Scoringnet: Learning key fragment for action quality assessment
25 with ranking loss in skilled sports, in: *ACCV*, Vol. 11366, 2019, pp. 149–164.
- 26 [24] Z. Luo, Y. Xiao, F. Yang et al., Rhythmer: Ranking-based skill assessment with rhythm-
27 aware transformer, *TCSVT* 35 (1) (2025) 259–272.
- 28 [25] X. Yu, Y. Rao, W. Zhao et al., Group-aware contrastive regression for action quality as-
29 sessment, in: *ICCV*, 2021, pp. 7899–7908.
- 30 [26] K. Zhou, Y. Ma, H. P. H. Shum et al., Hierarchical graph convolutional networks for action
31 quality assessment, *TCSVT* 33 (12) (2023) 7749–7763.
- 32 [27] J. Xu, Y. Rao, X. Yu et al., Finediving: A fine-grained dataset for procedure-aware action
33 quality assessment, in: *CVPR*, 2022, pp. 2939–2948.
- 34 [28] J. Xu, Y. Rao, J. Zhou et al., Procedure-aware action quality assessment: Datasets and
35 performance evaluation, *IJCV* 132 (12) (2024) 6069–6090.
- 36 [29] J. Xu, S. Yin, G. Zhao et al., Fineparser: A fine-grained spatio-temporal action parser for
37 human-centric action quality assessment, in: *CVPR*, 2024, pp. 14628–14637.
- 38 [30] K. Zhou, L. Wang, X. Zhang et al., Magr: Manifold-aligned graph regularization for
39 continual action quality assessment, in: *ECCV*, Vol. 15069, 2024, pp. 375–392.
- 40 [31] K. Zhou, Q. Pan, X. Zhang et al., Continual action quality assessment via adaptive
41 manifold-aligned graph regularization, *arXiv preprint arXiv:2510.06842* (2025).
- 42 [32] K. Zhou, H. P. H. Shum, F. W. B. Li et al., Phi: Bridging domain shift in long-term action
43 quality assessment via progressive hierarchical instruction, *TIP* 34 (2025) 3718–3732.
- 44 [33] S. Wang, D. Yang, P. Zhai et al., Tsa-net: Tube self-attention network for action quality
45 assessment, in: *ACM MM*, 2021, pp. 4902–4910.
- 46 [34] A. Dadashzadeh, S. Duan, A. Whone et al., Pecop: Parameter efficient continual pretrain-
47 ing for action quality assessment, in: *WACV*, 2024, pp. 42–52.
- 48 [35] J.-H. Pan, J. Gao, W.-S. Zheng, Action assessment by joint relation graphs, in: *ICCV*,
49 2019, pp. 6340–6349.
- 50 [36] J. Pan, J. Gao, W.-S. Zheng, Adaptive action assessment, *TPAMI* 44 (12) (2022) 8779–

- 1 8795.
- 2 [37] K. Gedamu, Y. Ji, Y. Yang et al., Self-supervised sub-action parsing network for semi-
- 3 supervised action quality assessment, *TIP* 33 (2024) 6057–6070.
- 4 [38] K. Gedamu, Y. Ji, Y. Yang et al., Fine-grained spatio-temporal parsing network for action
- 5 quality assessment, *TIP* 32 (2023) 6386–6400.
- 6 [39] R. Han, K. Zhou, A. Atapour-Abarghouei et al., Finecausal: A causal-based framework
- 7 for interpretable fine-grained action quality assessment, in: *CVPRW*, 2025, pp. 6008–
- 8 6017.
- 9 [40] Y. Tang, Z. Ni, J. Zhou et al., Uncertainty-aware score distribution learning for action
- 10 quality assessment, in: *CVPR*, 2020, pp. 9836–9845.
- 11 [41] B. Zhang, J. Chen, Y. Xu et al., Auto-encoding score distribution regression for action
- 12 quality assessment, *Neural Comput. Appl.* 36 (2) (2024) 929–942.
- 13 [42] K. Zhou, J. Li, R. Cai et al., Cofinal: Enhancing action quality assessment with coarse-to-
- 14 fine instruction alignment, in: *IJCAI*, 2024, pp. 1771–1779.
- 15 [43] M. Li, H.-B. Zhang, Q. Lei et al., Pairwise contrastive learning network for action quality
- 16 assessment, in: *ECCV*, Vol. 13664, 2022, pp. 457–473.
- 17 [44] Z. Cao, T. Simon, S.-E. Wei et al., Realtime multi-person 2d pose estimation using part
- 18 affinity fields, in: *CVPR*, 2017, pp. 7291–7299.
- 19 [45] C. Lugaresi, J. Tang, H. Nash et al., Mediapipe: A framework for building perception
- 20 pipelines, *arXiv preprint arXiv:1906.08172* (2019).
- 21 [46] Y. Xu, J. Zhang, Q. Zhang et al., Vitpose: Simple vision transformer baselines for human
- 22 pose estimation, *NeurIPS* 35 (2022) 38571–38584.
- 23 [47] B. Yu, Y. Liu, X. Zhang et al., Egcen: An ensemble-based learning framework for exploring
- 24 effective skeleton-based rehabilitation exercise assessment, in: *IJCAI*, 2022, pp. 3681–
- 25 3687.
- 26 [48] X. B. Bruce, Y. Liu, K. C. C. Chan et al., Egcen++: A new fusion strategy for ensemble
- 27 learning in skeleton-based rehabilitation exercise assessment, *TPAMI* 46 (9) (2024) 6471–
- 28 6485.
- 29 [49] X. Wang, J. Wang, H. Hu, Skeleton-based action quality assessment via partially con-
- 30 nected lstm with triplet losses, in: *PRCV*, Vol. 13536, 2022, pp. 220–232.
- 31 [50] X. Bruce, Y. Liu, K. C. Chan et al., Skeleton-based human action evaluation using graph
- 32 convolutional network for monitoring alzheimer’s progression, *PR* 119 (2021) 108095.
- 33 [51] Y. Liao, A. Vakanski, M. Xian, A deep learning framework for assessing physical reha-
- 34 bilitation exercises, *TNSRE* 28 (2) (2020) 468–477.
- 35 [52] S. Yan, Y. Xiong, D. Lin, Spatial temporal graph convolutional networks for skeleton-
- 36 based action recognition, in: *AAAI*, Vol. 32, 2018.
- 37 [53] C. Zhou, J. Zeng, L. Qiu et al., An attention-based adaptive spatial-temporal graph con-
- 38 volutional network for long-video ergonomic risk assessment, *EAAI* 131 (2024) 107780.
- 39 [54] J. Xia, M. Zhuge, T. Geng et al., Skating-mixer: Long-term sport audio-visual modeling
- 40 with mlps, in: *AAAI*, Vol. 37, 2023, pp. 2901–2909.
- 41 [55] Y. Zhong, F. Zhang, Y. Demiris, Contrastive self-supervised learning for automated multi-
- 42 modal dance performance assessment, in: *ICASSP*, 2023, pp. 1–5.
- 43 [56] L.-A. Zeng, W.-S. Zheng, Multimodal action quality assessment, *TIP* 33 (2024) 1600–
- 44 1613.
- 45 [57] X. Wang, P.-J. Li, Y.-Y. Shen, Attention-driven multimodal alignment for long-term action
- 46 quality assessment, *Appl. Soft Comput.* 183 (2025).
- 47 [58] Z. Du, D. He, X. Wang et al., Learning semantics-guided representations for scoring figure
- 48 skating, *TMM* 26 (2023) 4987–4997.
- 49 [59] S. Zhang, S. Bai, G. Chen et al., Narrative action evaluation with prompt-guided multi-
- 50 modal interaction, in: *CVPR*, 2024, pp. 18430–18439.

- 1 [60] H. Xu, X. Ke, Y. Li et al., Vision-language action knowledge learning for semantic-aware
2 action quality assessment, in: ECCV, Vol. 15100, 2024, pp. 423–440.
- 3 [61] H. Xu, X. Ke, H. Wu et al., Language-guided audio-visual learning for long-term sports
4 assessment, in: CVPR, 2025, pp. 23967–23977.
- 5 [62] Y. Ji, L. Ye, H. Huang et al., Localization-assisted uncertainty score disentanglement
6 network for action quality assessment, in: ACM MM, 2023, pp. 8590–8597.
- 7 [63] S. Zahan, G. M. Hassan, A. Mian, Learning sparse temporal video mapping for action
8 quality assessment in floor gymnastics, *IEEE Transactions on Instrumentation and Mea-
9 surement* (2024).
- 10 [64] Y. Ding, S. Zhang, S. Liu et al., 2m-af: A strong multi-modality framework for human
11 action quality assessment with self-supervised representation learning, in: ACM MM,
12 2024, pp. 1564–1572.
- 13 [65] L. Dong, W. Wang, Y. Qiao et al., Lucidaction: A hierarchical and multi-model dataset
14 for comprehensive action quality assessment, in: NeurIPS, 2024.
- 15 [66] M. A. A. H. Khan, H. Shahriar, Mrehab: Mutlimodal data acquisition and modeling
16 framework for assessing stroke and cardiac rehabilitation exercises, in: COMPSAC,
17 IEEE, 2022, pp. 452–453.
- 18 [67] D.-W. Kim, J. E. Park, M.-J. Kim et al., Automatic assessment of upper extremity function
19 and mobile application for self-administered stroke rehabilitation, *TNSRE* 32 (2024) 652–
20 661.
- 21 [68] L.-A. Zeng, F.-T. Hong, W.-S. Zheng et al., Hybrid dynamic-static context-aware attention
22 network for action assessment in long videos, in: ACM MM, 2020, pp. 2526–2534.
- 23 [69] F. Lin, J. Huang, Z. Chen et al., Enhancing long-term action quality assessment: A dual-
24 modality dataset and causal cross-modal framework for trampoline gymnastics, *Sensors*
25 25 (18) (2025) 5824.
- 26 [70] P. Parmar, B. T. Morris, Action quality assessment across multiple actions, in: WACV,
27 2019, pp. 1468–1476.
- 28 [71] T. Nagai, S. Takeda, M. Matsumura et al., Action quality assessment with ignoring scene
29 context, in: ICIP, 2021, pp. 1189–1193.
- 30 [72] X. Zhang, H. Feng, M. S. Hossain et al., Scaled background swap: Video augmentation
31 for action quality assessment with background debiasing, *ACM TOMM* 21 (8) (2025).
- 32 [73] S.-J. Zhang, J.-H. Pan, J. Gao et al., Semi-supervised action quality assessment with self-
33 supervised segment feature recovery, *TCSVT* 32 (9) (2022) 6017–6028.
- 34 [74] H. Pirsivash, C. Vondrick, A. Torralba, Assessing the quality of actions, in: ECCV, 2014,
35 pp. 556–571.
- 36 [75] A. Krizhevsky, I. Sutskever, G. E. Hinton, Imagenet classification with deep convolutional
37 neural networks, *NeurIPS* 25 (2012).
- 38 [76] D. Tran, L. Bourdev, R. Fergus et al., Learning spatiotemporal features with 3d convolu-
39 tional networks, in: ICCV, 2015, pp. 4489–4497.
- 40 [77] Z. Qiu, T. Yao, T. Mei, Learning spatio-temporal representation with pseudo-3d residual
41 networks, in: ICCV, 2017, pp. 5533–5541.
- 42 [78] J. Carreira, A. Zisserman, Quo vadis, action recognition? a new model and the kinetics
43 dataset, in: CVPR, 2017, pp. 6299–6308.
- 44 [79] Z. Liu, J. Ning, Y. Cao et al., Video swin transformer, in: CVPR, 2022, pp. 3202–3211.
- 45 [80] W. Kay, J. Carreira, K. Simonyan et al., The kinetics human action video dataset, *arXiv*
46 preprint arXiv:1705.06950 (2017).
- 47 [81] J. Carreira, E. Noland, A. Banki-Horvath et al., A short note about kinetics-600, *arXiv*
48 preprint arXiv:1808.01340 (2018).
- 49 [82] X. Xiang, Y. Tian, A. Reiter et al., S3d: Stacking segmental p3d for action quality assess-
50 ment, in: ICIP, 2018, pp. 928–932.

- 1 [83] L.-J. Dong, H.-B. Zhang, Q. Shi et al., Learning and fusing multiple hidden substages for
2 action quality assessment, *Knowl.-Based Syst.* 229 (2021) 107388.
- 3 [84] H.-B. Zhang, L.-J. Dong, Q. Lei et al., Label-reconstruction-based pseudo-subscore learn-
4 ing for action quality assessment in sporting events, *Appl. Intell.* 53 (9) (2023) 10053–
5 10067.
- 6 [85] Z. Li, Y. Huang, M. Cai et al., Manipulation-skill assessment from videos with spatial
7 attention network, in: *ICCV*, 2019, pp. 4385–4395.
- 8 [86] Y. Bai, D. Zhou, S. Zhang et al., Action quality assessment with temporal parsing trans-
9 former, in: *ECCV*, Vol. 13664, 2022, pp. 422–438.
- 10 [87] J. Gao, J.-H. Pan, S.-J. Zhang et al., Automatic modelling for interactive action assess-
11 ment, *IJCV* 131 (3) (2023) 659–679.
- 12 [88] X. Ke, H. Xu, X. Lin et al., Two-path target-aware contrastive regression for action quality
13 assessment, *Inf. Sci.* 664 (2024) 120347.
- 14 [89] R. Han, K. Zhou, S. Chen et al., Calfow: Enhancing long-term action quality assessment
15 with causal counterfactual flow, *arXiv preprint arXiv:2511.21653* (2025).
- 16 [90] J. Xu, S. Yin, Y. Peng, Human-centric fine-grained action quality assessment, *TPAMI*
17 47 (8) (2025) 6242–6255.
- 18 [91] K. He, X. Zhang, S. Ren et al., Deep residual learning for image recognition, in: *CVPR*,
19 2016, pp. 770–778.
- 20 [92] X. Dong, X. Liu, W. Li et al., Interpretable long-term action quality assessment, in:
21 *BMVC*, 2024.
- 22 [93] Z. Li, L. Gu, W. Wang et al., Surgical skill assessment via video semantic aggregation, in:
23 *MICCAI*, 2022, pp. 410–420.
- 24 [94] L. Okamoto, P. Parmar, Hierarchical neurosymbolic approach for comprehensive and ex-
25 plainable action quality assessment, in: *CVPRW*, 2024, pp. 3204–3213.
- 26 [95] J. Gao, W.-S. Zheng, J.-H. Pan et al., An asymmetric modeling for action assessment, in:
27 *ECCV*, Vol. 12375, 2020, pp. 222–238.
- 28 [96] S. Zhang, W. Dai, S. Wang et al., Logo: A long-form video dataset for group action
29 quality assessment, in: *CVPR*, 2023, pp. 2405–2414.
- 30 [97] T. He, H. Liu, Y. Li et al., Collaborative weakly supervised video correlation learning for
31 procedure-aware instructional video analysis, in: *AAAI*, Vol. 38, 2024, pp. 2112–2120.
- 32 [98] T. He, H. Liu, Z. Ni et al., Achieving procedure-aware instructional video correlation
33 learning under weak supervision from a collaborative perspective, *IJCV* 133 (4) (2025)
34 2070–2095.
- 35 [99] H. Fang, W. Zhou, H. Li, End-to-end action quality assessment with action parsing trans-
36 former, in: *VCIP*, IEEE, 2023, pp. 1–5.
- 37 [100] Q. Lei, H. Zhang, J. Du, Temporal attention learning for action quality assessment in
38 sports video, *Signal, Image and Video Processing* 15 (7) (2021) 1575–1583.
- 39 [101] Y. Gao, S. S. Vedula, C. E. Reiley et al., Jhu-isi gesture and skill assessment working set
40 (jigsaws): A surgical activity dataset for human motion modeling, in: *MICCAI Workshop*,
41 Vol. 3, 2014, p. 3.
- 42 [102] J. Wang, Z. Du, A. Li et al., Assessing action quality via attentive spatio-temporal convo-
43 lutional networks, in: *PRCV*, Vol. 12306, 2020, pp. 3–16.
- 44 [103] J. Liu, H. Wang, W. Zhou et al., Adaptive spatiotemporal graph transformer network for
45 action quality assessment, *TCSVT* 35 (7) (2025) 6628–6639.
- 46 [104] Q. An, M. Qi, H. Ma, Multi-stage contrastive regression for action quality assessment, in:
47 *ICASSP*, 2024, pp. 4110–4114.
- 48 [105] M. Qi, H. Ye, J. Peng et al., Action quality assessment via hierarchical pose-guided multi-
49 stage contrastive regression, *TIP* 34 (2025) 6461–6474.
- 50 [106] D. Liu, Q. Li, T. Jiang et al., Towards unified surgical skill assessment, in: *CVPR*, 2021,

- pp. 9517–9526.
- [107] X. Dong, X. Liu, W. Li et al., Uil-aqa: Uncertainty-aware clip-level interpretable action quality assessment, *IJCV* 134 (24) (2026).
- [108] M. Fang, Y. Zhou, J. Ren et al., A teacher action quality assessment method based on label constraint strategy, in: *ICASSP*, 2025.
- [109] C. Wang, X. Ding, X. Zhao et al., Adaptive frequency-aware network for action quality assessment 31 (5) (2025).
- [110] C. Zhou, Y. Huang, H. Ling, Uncertainty-driven action quality assessment, *arXiv preprint arXiv:2207.14513* (2022).
- [111] H. Matsuyama, N. Kawaguchi, B. Y. Lim, Iris: Interpretable rubric-informed segmentation for action quality assessment, in: *ICIUI*, 2023, pp. 368–378.
- [112] H. Doughty, W. Mayol-Cuevas, D. Damen, The pros and cons: Rank-aware temporal attention for skill determination in long videos, in: *CVPR*, 2019, pp. 7854–7863.
- [113] H.-M. Qiu, H.-B. Zhang, Q. Lei et al., Learning referee evaluation and assessing action quality from coarse to fine in diving sport, *Neurocomputing* 648 (2025).
- [114] Y.-M. Li, L.-A. Zeng, J.-K. Meng et al., Continual action assessment via task-consistent score-discriminative feature distribution modeling, *TCSVT* 34 (10) (2024) 9112–9124.
- [115] P.-X. Lian, Z.-G. Shao, Improving action quality assessment with across-staged temporal reasoning on imbalanced data, *Appl. Intell.* 53 (24) (2023) 30443–30454.
- [116] H. Jain, G. Harit, A. Sharma, Action quality assessment using siamese network-based deep metric learning, *TCSVT* 31 (6) (2021) 2260–2273.
- [117] Y. Pan, C. Zhang, G. Bertasius, Basket: A large-scale video dataset for fine-grained skill estimation, in: *CVPR*, 2025, pp. 28952–28962.
- [118] X. Ding, C. Wang, X. Zhao et al., Scene-aware contrastive regression for multi-person action quality assessment, *Appl. Intell.* 55 (16) (2025).
- [119] A. Kanade, M. Sharma, M. Muniyandi, Attention-guided deep learning framework for movement quality assessment, in: *ICASSP*, 2023, pp. 1–5.
- [120] Q. Lei, H. Li, H. Zhang et al., Multi-skeleton structures graph convolutional network for action quality assessment in long videos, *Appl. Intell.* 53 (19) (2023) 21692–21705.
- [121] H. Li, Q. Lei, H. Zhang et al., Skeleton-based deep pose feature learning for action quality assessment on figure skating videos, *JVCIR* 89 (2022) 103625.
- [122] Q. Lei, H.-B. Zhang, J.-X. Du et al., Learning effective skeletal representations on rgb video for fine-grained human action quality assessment, *Electronics* 9 (4) (2020) 568.
- [123] A. Kryeem, S. Raz, D. Eluz et al., Personalized monitoring in home healthcare: An assistive system for post hip replacement rehabilitation, in: *ICCV*, 2023, pp. 1868–1877.
- [124] A. Abedi, M. Malmirian, S. S. Khan, Cross-modal video to body-joints augmentation for rehabilitation exercise quality assessment, *arXiv preprint arXiv:2306.09546* (2023).
- [125] U. Gallardo, F. Caro, E. Hernández et al., Gymetricpose: A light-weight angle-based graph adaptation for action quality assessment, in: *CBMS, IEEE*, 2024, pp. 43–50.
- [126] S.-E. Wei, V. Ramakrishna, T. Kanade et al., Convolutional pose machines, in: *CVPR*, 2016, pp. 4724–4732.
- [127] T. Wang, M. Jin, J. Wang et al., Towards a data-driven method for rgb video-based hand action quality assessment in real time, in: *Annual ACM Symposium on Applied Computing*, 2020, pp. 2117–2120.
- [128] MMPose Contributors, MMPose: Openmmlab pose estimation toolbox and benchmark, <https://github.com/open-mmlab/mmpose>, accessed: December 01, 2024 (2020).
- [129] Z. Li, H. Chen, J. Cai et al., Segmentation and quality assessment of continuous fitness movements based on vision, in: *International Conference on Intelligent Computing*, Springer, 2024, pp. 96–107.

- 1 [130] K. Iskakov, E. Burkov, V. Lempitsky et al., Learnable triangulation of human pose, in:
2 ICCV, 2019, pp. 7718–7727.
- 3 [131] K. Zheng, J. Wu, J. Zhang et al., A skeleton-based rehabilitation exercise assessment
4 system with rotation invariance, TNSRE 31 (2023) 2612–2621.
- 5 [132] S. Mehraban, V. Adeli, B. Taati, Motionagformer: Enhancing 3d human pose estimation
6 with a transformer-genformer network, in: WACV, 2024, pp. 6920–6930.
- 7 [133] C. Zimmermann, T. Brox, Learning to estimate 3d hand pose from single rgb images, in:
8 ICCV, 2017, pp. 4903–4911.
- 9 [134] D. Pavllo, C. Feichtenhofer, D. Grangier et al., 3d human pose estimation in video with
10 temporal convolutions and semi-supervised training, in: ICCV, 2019, pp. 7753–7762.
- 11 [135] C.-I. Joung, S. Byun, S. Baek, Contrastive learning for action assessment using graph
12 convolutional networks with augmented virtual joints, IEEE Access (2023).
- 13 [136] B. Garg, A. Postlmayr, P. Cosman et al., Short: Deep learning approach to skeletal per-
14 formance evaluation of physical therapy exercises, in: ACM/IEEE CHASE, 2023, pp.
15 168–172.
- 16 [137] A. Zanfir, E. Marinoiu, M. Zanfir et al., Deep network for the integrated 3d sensing of
17 multiple people in natural images, NeurIPS 31 (2018).
- 18 [138] M. Fieraru, M. Zanfir, S. C. Pirlea et al., Aifit: Automatic 3d human-interpretable feed-
19 back models for fitness training, in: CVPR, 2021, pp. 9919–9928.
- 20 [139] C. Li, X. Ling, S. Xia, A graph convolutional siamese network for the assessment and
21 recognition of physical rehabilitation exercises, in: ICANN, Springer, 2023, pp. 229–240.
- 22 [140] S. Deb, M. F. Islam, S. Rahman et al., Graph convolutional networks for assessment of
23 physical rehabilitation exercises, TNSRE 30 (2022) 410–419.
- 24 [141] L. Yao, Q. Lei, H. Zhang et al., A contrastive learning network for performance metric
25 and assessment of physical rehabilitation exercises, TNSRE 31 (2023) 3790–3802.
- 26 [142] C. Du, S. Graham, C. Depp et al., Assessing physical rehabilitation exercises using graph
27 convolutional network with self-supervised regularization, in: EMBC, 2021, pp. 281–285.
- 28 [143] J. Li, J. Xue, R. Cao et al., Finerehab: A multi-modality and multi-task dataset for reha-
29 bilitation analysis, in: CVPRW, 2024, pp. 3184–3193.
- 30 [144] M. Nekoui, L. Cheng, Enhancing human motion assessment by self-supervised represen-
31 tation learning., in: BMVC, 2021, p. 322.
- 32 [145] K. Gedamu, Y. Ji, Y. Yang et al., Visual-semantic alignment temporal parsing for action
33 quality assessment, TCSVT 35 (3) (2025) 2436–2449.
- 34 [146] T. Nagai, S. Takeda, S. Suzuki et al., Mmw-aqa: Multimodal in-the-wild dataset for action
35 quality assessment, IEEE Access (2024).
- 36 [147] S. Hirose, T. Kato, T. Yamashita et al., Action quality assessment model using special-
37 ists’ gaze location and kinematics data—focusing on evaluating figure skating jumps,
38 Sensors 23 (22) (2023) 9282.
- 39 [148] H. Yin, L. Gu, P. Parmar et al., Flex: A large-scale multi-modal multi-action dataset for
40 fitness action quality assessment, arXiv preprint arXiv:2506.03198 (2025).
- 41 [149] H. Xu, H. Wu, X. Ke et al., Mcmoe: Completing missing modalities with mix-
42 ture of experts for incomplete multimodal action quality assessment, arXiv preprint
43 arXiv:2511.17397 (2025).
- 44 [150] K. Ashutosh, T. Nagarajan, G. Pavlakos et al., Expertaf: Expert actionable feedback from
45 video, in: CVPR, 2025, pp. 13582–13594.
- 46 [151] K. Ashutosh, K. Grauman, Learning skill-attributes for transferable assessment in video,
47 in: NeurIPS, 2025.
- 48 [152] M. Qi, Y. Wu, X. Zhang et al., Explainable action form assessment by exploiting multi-
49 modal chain-of-thoughts reasoning, arXiv preprint arXiv:2512.15153 (2025).
- 50 [153] H. Xu, H. Wu, X. Ke et al., Quality-guided vision-language learning for long-term action

- quality assessment, *TMM* 27 (2025) 7326–7339.
- [154] J. Li, H. Hu, Q. Xing et al., Tai chi action quality assessment and visual analysis with a consumer rgb-d camera, in: *MMSp*, 2022.
- [155] Y. Zhang, X. Li, W. Chai et al., Pose-guided transformer for fine-grained action quality assessment, *TCSVT* 35 (8) (2025) 7940–7952.
- [156] H. Xu, X. Ke, H. Wu et al., Dancefix: An exploration in group dance neatness assessment through fixing abnormal challenges of human pose, in: *AAAI*, Vol. 39, 2025, pp. 8869–8877.
- [157] T. Moodley, D. van der Haar, I3d-ae-1stm: A 2-stream autoencoder for action quality assessment using a newly created cricket batsman video dataset, in: *WACV*, 2025, pp. 5470–5478.
- [158] W. Yun, M. Qi, F. Peng et al., Semi-supervised teacher-reference-student architecture for action quality assessment, in: *ECCV*, Vol. 15132, 2024, pp. 161–178.
- [159] L. Ye, K. Gedamu, J. Shao, Decoupling representations with quantized vectors for semi-supervised action quality assessment, in: *ICME*, 2025, pp. 1–6.
- [160] Y.-M. Li, A.-L. Wang, K.-Y. Lin et al., Techcoach: Towards technical-point-aware descriptive action coaching, *arXiv preprint arXiv:2411.17130* (2024).
- [161] Y.-M. Li, W.-J. Huang, A.-L. Wang et al., Egoexo-fitness: Towards egocentric and exocentric full-body action understanding, in: *ECCV*, 2024, pp. 363–382.
- [162] M. Fang, Y. Zhou, Q. Liu et al., A contrastive video language multimodal method for teacher action quality assessment, *ESWA* 297 (C) (2026) 129501.
- [163] A. Vakanski, H.-p. Jun, D. Paul et al., A data set of human body movements for physical rehabilitation exercises, *Data* 3 (1) (2018) 2.
- [164] Z. Chen, W. Sun, Y. Tian et al., Gaia: Rethinking action quality assessment for ai-generated videos, *NeurIPS* 37 (2024) 40111–40144.
- [165] Z. Zhang, W. Sun, X. Li et al., Human-activity agv quality assessment: A benchmark dataset and an objective evaluation metric, in: *ACM MM*, 2025, pp. 6771–6780.
- [166] L. Wang, B. Huang, Z. Zhao et al., Videomae v2: Scaling video masked autoencoders with dual masking, in: *CVPR*, 2023, pp. 14549–14560.
- [167] J. Park, H.-S. Kim, K. Ko et al., Videomamba: Spatio-temporal selective state space model, in: *ECCV*, Springer, 2024, pp. 1–18.
- [168] B. Jiang, X. Chen, W. Liu et al., Motiongpt: Human motion as a foreign language, *NeurIPS* 36 (2023) 20067–20079.
- [169] H. Zhang, X. Li, L. Bing, Video-llama: An instruction-tuned audio-visual language model for video understanding, 2023, pp. 543–553.
- [170] K. Zhou, C. Li, Q. Pan et al., Brima: Bridged modality adaptation for multi-modal continual action quality assessment, in: *CVPR*, 2026.
- [171] Y. Yuan, J. Song, U. Iqbal et al., Physdiff: Physics-guided human motion diffusion model, in: *CVPR*, 2023, pp. 16010–16021.
- [172] G. Paolo, J. Gonzalez-Billandon, B. Kégl, Position: A call for embodied ai, in: *ICML*, 2024.
- [173] Z. Feng, R. Xue, L. Yuan et al., Multi-agent embodied ai: Advances and future directions, *arXiv preprint arXiv:2505.05108* (2025).
- [174] J. L. Laughlin, L. R. Johnson, B. Ghanekar et al., Technical skills assessment in robotic surgery: A review of recent methods, *Methodist DeBakey Cardiovascular Journal* 21 (5) (2025) 49.
- [175] M. W. Boal, D. Anastasiou, F. Tesfai et al., Evaluation of objective tools and artificial intelligence in robotic surgery technical skills assessment: a systematic review, *British Journal of Surgery* 111 (1) (2024) znad331.

1 **Appendix A. Database Query Details**

2 To ensure a comprehensive and reproducible literature collection process, we con-
3 ducted a systematic search across five major academic databases, including IEEE Xplore,
4 Web of Science, Scopus, Google Scholar, and arXiv. The search focused on publica-
5 tions from 2014 to 2025, covering the rapid development period of deep learning-based
6 AQA methods.

7 The query strategy was designed to capture representative studies related to AQA,
8 skill assessment, and score prediction tasks in both sports and medical domains. The
9 core search terms included:

```
10 ("action quality assessment" OR "action scoring" OR  
11 "skill assessment" OR "performance evaluation" OR "score  
12 prediction")
```

13 Since different databases employ different indexing mechanisms, database-specific
14 query formulations were adopted as follows:

15 **IEEE Xplore**

```
16 ("All Metadata":"action quality assessment" OR "All  
17 Metadata":"action scoring" OR "All Metadata":"skill  
18 assessment")
```

19 **Web of Science**

```
20 TS=("action quality assessment" OR "action scoring"  
21 OR "skill assessment")
```

22 **Scopus**

```
23 TITLE-ABS-KEY("action quality assessment" OR "action  
24 scoring" OR "skill assessment")
```

25 **Google Scholar**

1 Keyword combinations were manually queried using phrases such as *action quality*
2 *assessment*, *skill assessment*, and *action scoring*, with relevance ranking and citation
3 impact considered during screening.

4 **arXiv**

5 `all:"action quality assessment" OR all:"skill assessment"`

6 The initial search results were further refined through title screening, abstract re-
7 view, and full-text examination according to the inclusion and exclusion criteria de-
8 scribed in Section 2. Duplicate studies, non-peer-reviewed works, papers without quan-
9 titative evaluation, and studies unrelated to score prediction were excluded to ensure
10 the final benchmark set was representative, reproducible, and technically rigorous.