

A retrieval-augmented LLM framework for PRISMA-aligned systematic reviews: a case study on action quality assessment

Authors anonymized as requested^{1*}

^{1*} Affiliations anonymized as requested.

Abstract

Systematic reviews and meta-analyses are essential for advancing scientific research, but they are time-consuming, labor-intensive, and prone to human error during literature selection and data extraction. This paper proposes a framework that integrates Large Language Models (LLMs) and retrieval-augmented generation (RAG) into the PRISMA 2020 workflow to assist and standardize key review tasks. The framework uses structured query templates to ensure reproducible bibliographic searches, an LLM-driven pipeline for title and abstract screening, and a RAG-based full-text eligibility assessment. Both automated stages rely on structured prompt templates that encode inclusion and exclusion criteria as machine-interpretable instructions. Beyond selecting relevant papers, it extracts specific information from documents, such as application details or datasets, and harmonizes it into a standardized metadata schema. To evaluate our proposal, we performed an experimental validation using an existing systematic literature review on Human Action Quality Assessment (AQA) as a reference. The framework achieved an extraction accuracy of 85.11% for domain-specific datasets and identified 25 additional data sources not reported in the original review. The article detection rate of 72.34% was primarily attributable to the dataset-oriented query design adopted in the identification phase, while the LLM-based screening and RAG-based eligibility stages retained 97.14% of the reference studies that entered the automated pipeline. These results suggest that the framework can effectively assist in identifying relevant studies and extracting information while maintaining alignment with PRISMA 2020 guidelines. The proposed approach is structured to facilitate adaptation to other research areas by declaring new query and prompt templates.

Keywords: PRISMA 2020, systematic review, large language models, retrieval-augmented generation

1 Introduction

Systematic reviews are essential to scientific research because they synthesize evidence in a structured way, support state-of-the-art assessment, and help identify current limitations and research gaps [1]. Yet, conducting them remains labor-intensive, often exceeding 1,000 person-hours due to extensive literature searches with low inclusion rates [3]. Furthermore, this process is typically sensitive to human variability, with screening errors ranging from 10% to 20% [30], and worst-case values reaching almost 40% [41]. Early pioneering work on automation in research synthesis showed that machine learning can support this process [25, 39], but many current solutions still address isolated steps (*e.g.*, screening) rather than the full workflow, or do not align with reporting standards such as PRISMA [27].

Large Language Models (LLMs), *i.e.*, transformer-based models trained on large text corpora, are increasingly used in systematic-review support tasks [43]. Early studies have shown that LLMs can attain human-like performance (*i.e.*, accuracy > 80%) only when screening full texts in imbalanced literature datasets, provided that highly reliable and carefully designed prompts are used [21].

However, many current approaches remain fragmented and do not operationalize the PRISMA workflow end-to-end. Existing studies also show limited workflow-level integration of retrieval-augmented generation (RAG) across PRISMA stages [9, 24, 37], which can reduce traceability and auditability of stage-wise decisions [33]. Although RAG improves factual grounding at the component level [14], its methodological integration in PRISMA-aligned end-to-end workflows remains limited.

In this paper, we present a unified framework that incorporates structured literature searches, LLM-driven screening, and RAG-based full-text assessment within the PRISMA 2020 workflow, an updated guideline that refines the original PRISMA [29]. We treat PRISMA as an operational framework to enforce transparent reporting, traceable decisions, and reproducible procedures across the full pipeline. In addition to identifying relevant studies, the framework extracts domain-specific information, such as datasets and application details, and organizes it into a consistent metadata schema. By integrating LLMs and RAG into a methodologically grounded PRISMA framework, our approach aims to support systematic reviews reducing manual effort and to provide a practical tool for transparent and reproducible research synthesis.

To evaluate the proposed framework, we conducted an empirical evaluation in one target domain, Human Action Quality Assessment (AQA). First, we selected a manually curated, peer-reviewed reference study [23], as the comparison baseline and replicated its PRISMA-aligned search, screening, and eligibility procedure as closely as possible. Then we applied our automated pipeline under comparable conditions. The evaluation focused on methodological reproducibility and on the framework’s ability to identify relevant studies and extract related information, not on computational performance.

We assessed how closely the resulting corpus and extracted outputs matched the reference baseline in terms of coverage and relevance. The framework achieved a reference-study detection rate of 72.34% and a reference-dataset extraction accuracy of 85.11%. Most discrepancies originated in the initial bibliographic search, while the

automated screening and eligibility stages accounted for only a marginal loss, retaining 97.14% of the reference studies available at the pipeline entry. Although validated on AQA in this study, the workflow is designed to allow adaptation across domains. By re-specifying query terms and extraction criteria in the prompt templates, the same PRISMA-aligned framework could potentially be applied to other research areas, without modifying the pipeline itself. However, this transferability remains to be empirically verified in future work.

The main contributions of this work are as follows:

- We propose an end-to-end PRISMA-aligned framework that integrates structured query design, LLM-based screening, and RAG-based full-text eligibility assessment in a single reproducible pipeline.
- We formalize operational artifacts for reproducibility, including query templates, prompt templates, harmonization rules, and machine-readable outputs for each PRISMA stage.
- We provide an empirical validation against a manually curated reference study, demonstrating substantial agreement in study retrieval and data-source extraction.
- We demonstrate that the framework can identify additional relevant data sources beyond those reported in the reference study, supporting its utility for evidence discovery in evolving research domains.

The remainder of this work is organized as follows. Section 2 discusses the relevant literature. Section 3 presents our proposed methodology, while Section 4 describes the experimental validation in a case-study setting, focusing on reproducibility and the ability to identify relevant studies and information in a manner comparable to a traditional manual review. This section also presents the reference study, its workflow and results, our partial replication of their approach, and the technical implementation of our framework. Section 5 discusses the evaluation results, and Section 6 concludes the paper with final remarks.

2 Related work

Systematizing the synthesis of scientific evidence has become a critical requirement as the volume of global literature grows exponentially [1]. The PRISMA framework was originally introduced to standardize this process [27] and subsequently updated in PRISMA 2020 to incorporate advances in existing methods and technologies [29]. Extensions such as PRISMA-S [31] have been developed to provide detailed guidance on reporting search strategies across multiple databases, supporting transparency and reproducibility. Yet, the practical application of these protocols still relies heavily on manual workflows, which are time-consuming and prone to human error, motivating research interest in automated solutions. Recent broad mappings of the field indicate that evidence-synthesis automation now spans both classical machine learning tools and LLM-based systems; however, the landscape remains heterogeneous in terms of workflow coverage, reporting practices, and validation depth [16, 20].

Early efforts in this domain focused on semi-automated title/abstract screening and structured information extraction to support evidence synthesis. Tools such as

RobotReviewer applied NLP to assist with risk-of-bias assessment in randomized controlled trials [25]. In parallel, active-learning-based screening systems (*e.g.*, ASReview) implemented a reviewer-in-the-loop workflow in which the model is retrained after each reviewer label and used to rank the remaining records for manual screening, yielding substantial reductions in screening effort at high target recall [39]. In this context, a complementary work by de la Torre-López et al. [38] describes automated pipelines for classification, clustering, and extraction based on discriminative models. Collectively, these approaches reinforced that automation in systematic reviewing must remain auditable and conservative, with clear traceability of decisions and outputs [33]. Khalil et al. [20] offered a scoping overview of available automation tools, their validation status, and limitations. More recent reviews confirm that, although the field is expanding rapidly, many solutions remain task-specific rather than fully end-to-end [8, 16].

Recent advances in LLM research extend earlier discriminative approaches; unlike prior tools focused on narrow classification, LLMs introduce generative capabilities and richer semantic understanding to the review process [43]. These technologies have paved the way for systems that integrate diverse data sources to support researchers during the selection and extraction phases [39]. At the same time, these studies highlight the variability of LLM outputs, stressing the need for validation mechanisms and traceable reporting of model decisions [28]. Empirical evaluations of LLM-assisted workflows further illustrate this trade-off. Scherbakov et al. [32] test an LLM-based systematic review workflow and report notable efficiency gains paired with the continued need for human oversight to manage errors. Broader analyses also discuss ethical and transparency concerns around assisted reviews [2, 40]. Recent studies have examined how LLMs can operationalize PRISMA principles. Thode et al. [37] and Delgado-Chaves et al. [9] studied the role of LLMs in the screening phase and reported significant increases in efficiency. Their empirical evaluations suggest that although LLMs can substantially reduce workload, human supervision remains essential to mitigate risks of hallucinations or misinterpretations of text. Complementarily, Galli et al. [13] examined the use of LLMs in abstract screening, highlighting prompt engineering, zero-shot/few-shot classification, and the scalability constraints of AI-assisted screening. To formalize this supervision, recent research has proposed hybrid methodological frameworks: Brincoveanu et al. [4] developed a collaborative human-AI environment to minimize error rates, while Malik and Terzidis [24] outlined specific operational guidelines for augmenting PRISMA workflows with AI checkpoints. Furthermore, RAG has emerged as a key technical method for grounding LLM outputs in source documents, thereby improving factuality [14, 22]. Recent surveys map the broader RAG design space, but also show that practical integration into evidence-synthesis workflows and explicit reasoning support remain limited [5, 26].

Regarding reporting standards, the emergence of Generative AI (GenAI) has highlighted the need to adapt guidelines to AI-assisted review workflows. Recent governance-oriented contributions, including the Cochrane statement and the PRISMA-trAIce checklist, emphasize explicit reporting of AI-specific design choices [11, 17]. For example, Shailendra et al. [34] present L-PRISMA, a preprint that proposes reporting extensions for AI-specific elements such as model provenance, prompt

design, and verification procedures. In parallel, Forero et al. [12] applied LLMs to assess the adherence of published reviews to PRISMA 2020 items, showing the potential of these models as auditing tools. Finally, toward broader automation, agent-based pipelines have been proposed to automate multiple stages of the review process. Wang et al. [42] proposed TrialMind, an LLM pipeline for clinical evidence synthesis integrating study search, citation screening, data extraction, and evidence synthesis, evaluated on a dedicated benchmark of 100 published systematic reviews. Cao et al. [6] introduced a multi-agent system capable of executing sequential tasks, and Susnjak [35] explored the use of fine-tuned models for specific domains (PRISMA-DFILM), with improved relevance in targeted application settings.

Taken together, these studies show substantial progress in automating individual review tasks. Table 1 summarizes representative approaches by PRISMA-stage coverage, integration level, use of RAG, and reproducibility/auditability criteria. As the comparison illustrates, the landscape remains dominated by stage-specific or partially integrated solutions, with limited simultaneous coverage of full PRISMA alignment, explicit reproducibility artifacts, and formally grounded RAG across the entire workflow. This limitation substantiates the methodological gap outlined in Section 1 and motivates the framework proposed here.

Table 1 Comparison of representative AI-assisted systematic review approaches with respect to PRISMA workflow coverage and methodological rigor.

Study	PRISMA stage(s)	Approach type	Reproducibility	Auditability	Attributes Extraction	RAG use
[25]	Inclusion	Task-specific Tool	Limited	Moderate	Moderate	No
[39]	Screening	Task-specific Tool	Moderate	Moderate	None	No
[9, 37]	Screening	Partial LLM Workflow	Limited	Moderate	None	No
[24, 34]	All stages	Theoretical Framework	None	High	None	No
[38]	Screening, Inclusion	Partial Workflow	Limited	Moderate	Moderate	No
[32]	Screening, Eligibility, Inclusion	Partial LLM Workflow	Moderate	Moderate	High	No
[13]	Screening	Partial LLM Workflow	Limited	Moderate	None	No
[4]	Screening, Eligibility	Partial LLM Workflow	Moderate	High	Moderate	No
[12]	N/A	Audit Tool	Moderate	High	None	No
[42]	All stages	End-to-end LLM Workflow	Moderate	High	High	No
[6]	Screening, Eligibility, Inclusion	Multi-agent LLM Framework	Limited	Moderate	High	No
[35]	Screening, Eligibility	Task-specific Tool	Moderate	Moderate	Moderate	No
This work	All stages	End-to-end LLM Framework	High	High	High	Yes

Note: Reproducibility, Auditability, and Attributes Extraction are rated on a four-level scale: *None* (no support), *Limited* (partial or ad-hoc), *Moderate* (systematic but incomplete), *High* (fully supported and documented). RAG use indicates whether Retrieval-Augmented Generation is employed. N/A denotes tools not directly mapped to a PRISMA stage.

3 Proposed methodology

In this section, we describe our proposal for a methodology that extends the PRISMA 2020 workflow for systematic reviews [29] to integrate LLM and RAG components. Our workflow, shown in Figure 1, follows the original four phases of the PRISMA 2020 methodology (*i.e.*, identification, screening, eligibility, and inclusion) while integrating automated and semi-automated components. Each automated component is associated with a set of parameters and outputs that must be reported to ensure reproducibility. In conventional PRISMA-aligned reviews, data extraction is formally

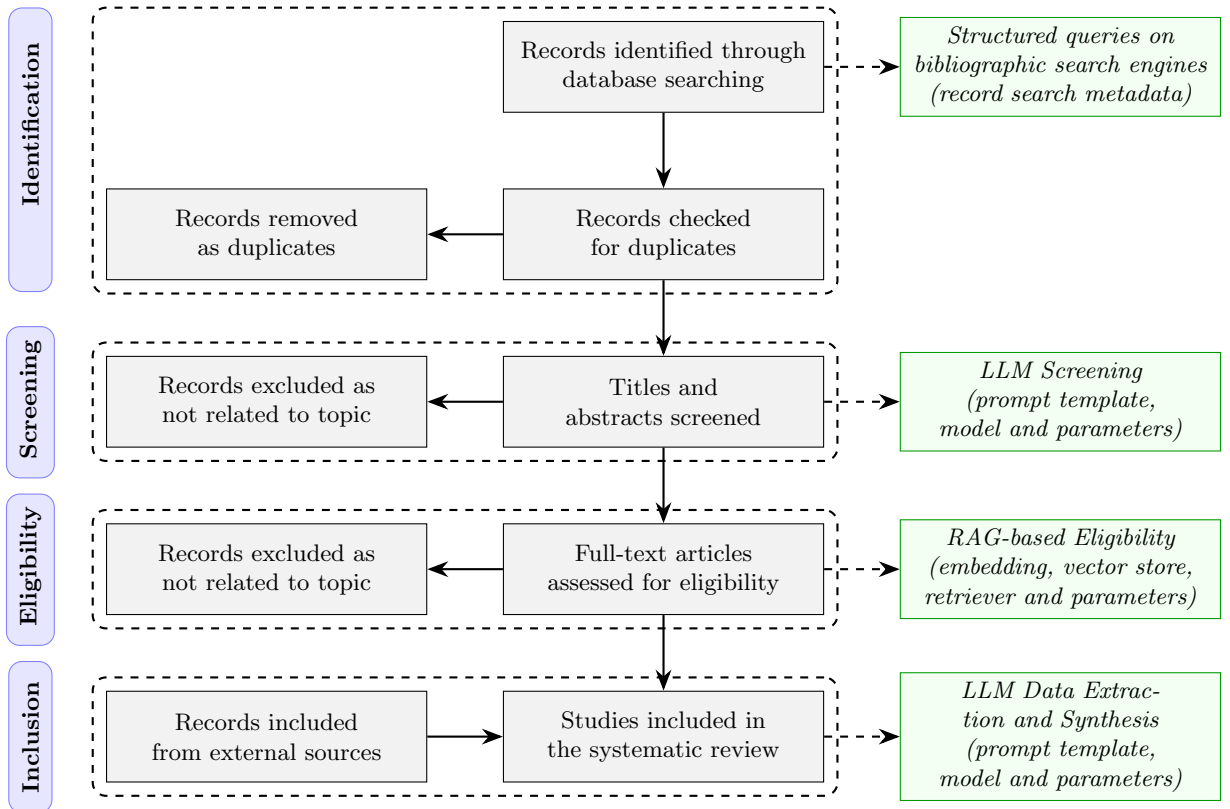


Fig. 1 Proposed workflow, with integrated LLM and RAG components, in reference to the traditional PRISMA 2020 diagram [29].

prescribed as a post-selection step, conducted once the final corpus has been established. In the proposed framework, attribute extraction is instead embedded directly within the screening and eligibility phases, allowing extracted information to propagate across pipeline stages and reducing redundant effort on records that will ultimately be excluded. In the following, we provide a detailed description of each phase, focusing on the theoretical and procedural alignment with PRISMA 2020 requirements and highlighting the novel procedural elements introduced by this framework.

3.1 Identification: information sources and search strategy

The identification phase aims to design and execute search queries across multiple bibliographic search engines. Queries are composed of a fixed set of declarative components that together delimit the scope of the search and standardize the query formulation process. For each search, the exact query string should be declared along with the date it was executed and the bibliographic search engine used. Declaring these details follows the PRISMA-S recommendations for reporting complete search

```

TITLE-ABS-KEY ( <DOMAIN_TERMS> )
AND TITLE-ABS-KEY ( <OBJECT_TERMS> )
AND PUBYEAR > <START_YEAR> AND PUBYEAR < <END_YEAR>
AND ( LIMIT-TO ( DOCTYPE , <DOCTYPE_FILTERS> ) )
AND ( LIMIT-TO ( SUBJAREA , <SUBJAREA_FILTERS> ) )
AND ( LIMIT-TO ( LANGUAGE , <LANGUAGE_FILTERS> ) )

```

Fig. 2 Topic-agnostic Scopus search template for systematic review searches.

strategies [31], ensuring transparency and supporting reproducibility. Formally, we can define a Query Q as:

$$Q = \langle D, O, P, QF \rangle \quad (1)$$

where:

- D (Domain / Subdomain): the principal research area and any narrower subfields, guiding the domain-specific terminology used in the query;
- O (Object of study): the phenomenon, method, or element of interest that drives core topical terms of the search;
- P (Period of reference): explicit temporal bounds (publication year range) used to control temporal bias and support reproducible results;
- QF (Qualifiers): optional filters such as publication type, geographic focus, language, or other attributes that restrict the search to the review objective.

To operationalize Equation (1), we propose the use of query templates. A query template is a reusable structure composed of named fields corresponding to the components D , O , P , and QF , together with angle-bracketed placeholders (*i.e.*, strings intended to be replaced by specific terms, *e.g.*, <DOMAIN_TERMS>) that are instantiated with keywords and Boolean operators. Templates also encode the syntax rules of the target bibliographic search engine and preserve the logical structure of the query. Each field in a template is populated with all relevant terms corresponding to its component, including synonyms and acronyms, to represent the domain or object of study. We motivate the use of query templates through four objectives. First, they provide a structured framework reusable and adaptable in future systematic reviews on related topics. Second, they organize queries into distinct components, thereby clearly delimiting the research objective at this stage of the workflow. Third, they increase robustness to lexical variation by systematically covering both domain-specific and object-specific terminology. Fourth, they ensure transparency and reproducibility by making the exact executed queries explicit and auditable.

As an example, a topic-agnostic Scopus query template is presented in Figure 2, illustrating the query structure without specifying a particular research topic. The angle-bracketed placeholders are replaced with the terms specific to the review under consideration, as defined in Equation (1).

Executing queries on search engines returns a list of documents (*i.e.*, papers) stored in a centralized repository. Each document record is normalized to a structured metadata schema to ensure consistency across sources and support data processing. The schema presented in Table 2 is a practical metadata model proposed by the

Table 2 Metadata schema required for each record in the identification phase.

Metadata field	Description	Purpose	Required
Identifier	Persistent and unique identifier (<i>e.g.</i> , DOI, PMID)	Ensure unambiguous identification and traceability across all systematic review stages	Yes
Title	Full record title	Support initial eligibility assessment in the PRISMA screening phase	Yes
Publication year	Year of publication	Enable verification of reference period and reproducibility of the review	Yes
Abstract	Study summary	Support PRISMA screening and subsequent data extraction	Yes
Source	Bibliographic source or search engine	Verify provenance and ensure traceability of the resource	Yes
Authors	List of authors as reported by the source	Support provenance tracking and optional verification during deduplication	Optional
Keywords	Author or indexer assigned keywords	Assist eligibility filtering and verification in later stages	Optional
Document type	Type of document (<i>e.g.</i> , article, conference paper, review)	Assist eligibility filtering and verification in later stages	Optional
Language	Language of publication	Assist eligibility filtering and verification in later stages	Optional

authors to capture key bibliographic details of each record, including the identifier, title, publication year, abstract, source, and other optional attributes. While PRISMA 2020 recommends reporting the number of retrieved records and the search strategies adopted, it does not prescribe a specific metadata schema. The schema proposed here aligns with standard bibliographic metadata conventions and provides a reproducible basis for the subsequent phases.

As a final step in the identification phase, deduplication is required by PRISMA 2020. The goal is to remove all redundant copies of the same record. When persistent identifiers (*e.g.*, DOI or PMID) are available, simple deterministic rules are applied to remove duplicates. For records lacking reliable identifiers, which may occur when integrating results from grey literature or non-standard sources, reproducible fuzzy-matching heuristics may be applied, or records may be resolved through manual assessment.

3.2 Screening: title and abstract evaluation

The screening phase performs an initial eligibility assessment at the title and abstract level using an LLM-driven classification pipeline (hereafter referred to as the "screening pipeline"), as illustrated in Figure 3. This stage corresponds to the PRISMA screening phase and is designed to reduce the candidate set while preserving alignment between the article content and the predefined inclusion and exclusion criteria for the systematic review. The screening stage operates exclusively on bibliographic metadata generated in the identification stage and does not require access to full-text content.

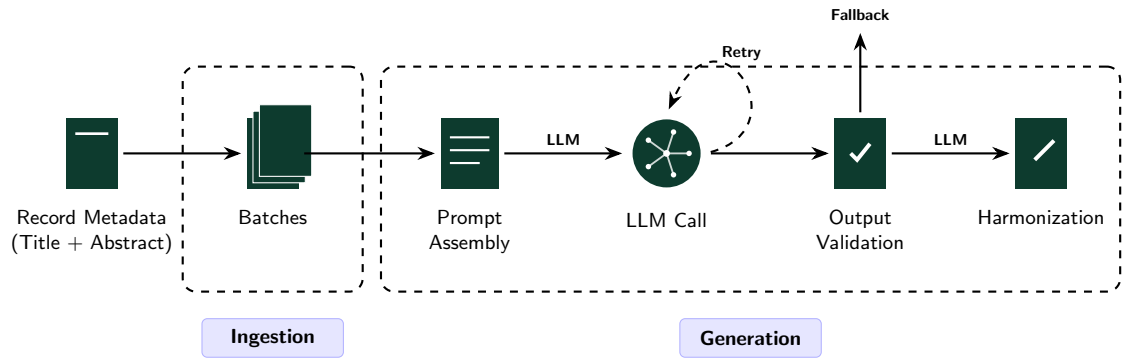


Fig. 3 Screening pipeline supporting title and abstract evaluation under PRISMA workflow.

Methodologically, the screening pipeline is organized into two main conceptual steps:

1. **Ingestion:** Bibliographic records are prepared for the LLM evaluation by extracting the title and abstract of each study and grouping them into batches. A batch is a fixed-size subset of records that are jointly submitted to the LLM in a single prompt. Batch size is a parameter that must be explicitly declared: smaller batches improve output consistency but increase processing time and API calls, while larger batches improve efficiency but may degrade screening quality if the cumulative input approaches the LLM’s context window limit. The batch size must therefore be chosen in relation to the context window of the selected model and reported as part of the screening metadata.
2. **Generation:** Using the batches prepared in the ingestion step, each group of records is evaluated by providing the LLM with the title and abstract, along with a prompt template that encodes the inclusion and exclusion criteria as machine-interpretable instructions. Each batch is dynamically concatenated into the prompt at runtime, ensuring that the template itself remains fixed. The prompt is applied uniformly across all records to ensure reproducibility and methodological consistency. The LLM is required to return a structured response that conforms to the *Output schema* field specified in Table 3. Following LLM evaluation, responses are automatically validated against the schema (*i.e.*, output validation). Finally, selected attributes are standardized in a subsequent harmonization step to ensure consistency across records.

A key distinction of this workflow compared to manual screening is the shift from subjective human assessment to an automated process. Upon executing these steps, the pipeline produces a structured object for each screened record. This output includes:

- a categorical screening decision (*include*, *exclude*, or *uncertain*);
- a short justification that links the decision to the stated eligibility criteria;
- any extracted study-level fields required for subsequent analyses (*e.g.*, named data sources or other identifying attributes).

To ensure the reliability of this automated process, all LLM calls use model parameters configured to minimize output variability, and responses are validated against the expected *Output schema* prior to acceptance. Outputs that fail schema validation, are incomplete, or cannot be reliably parsed are not forcefully classified. Instead, such cases are conservatively labeled as *uncertain* and either passed to subsequent workflow steps or flagged for human review. A conservative fallback policy is therefore intrinsic to the screening design: when automated classification is not completed with sufficient confidence, records are kept rather than discarded. This approach minimizes the risk of false negatives during the screening stage.

In addition to the screening decision, selected attributes can be extracted from the LLM output for use in later phases of the workflow. These fields undergo a dedicated harmonization step to reduce variability in attribute names, such as differences in capitalization, punctuation, abbreviations, or descriptive modifiers. Typical normalization rules include lowercasing, trimming whitespace, removing diacritics such as accent marks and similar signs (*e.g.*, "Intérnational" to "International"), standardizing punctuation, and expanding common abbreviations (*e.g.*, "ML" to "Machine Learning").

Robustness measures are integrated to handle model errors and execution failures. The screening pipeline includes retry and error-management policies to prevent data loss during temporary service interruptions. Human supervision is explicitly integrated as a complementary component of the screening phase. First, manual review is reserved for all records labeled *uncertain*. Second, a configurable sample of automatically accepted and discarded records, typically in the range from 5% to 10% of each group, is examined for quality control purposes. Finally, manual review is conducted for cases where the automated harmonization results are ambiguous. The frequency and scope of manual checks are recorded and may be used to refine screening criteria or calibration strategies.

Table 3 summarizes the key items reported from the use of an LLM-based screening procedure at the title and abstract level. Most items capture metadata about the phase as a whole, while provenance logging is recorded for each individual record.

3.3 Eligibility: full-text assessment via RAG

The automated eligibility phase implements the full-text assessment under PRISMA 2020 requirements using a RAG-based pipeline. Unlike the standard PRISMA methodology, which assumes manual reading of full-text documents by the researchers involved, this framework introduces automated retrieval as the primary mechanism for evidence selection. Full-text articles are first acquired from their original sources and stored in a repository. Each document is then transformed into a numerical vector representation of text (*i.e.*, embeddings) using an embedding model, and these embeddings are indexed in a persistent vector store (*i.e.*, a database optimized for searching information) to enable efficient semantic similarity search. During assessment, relevant context windows are retrieved from the vector store and provided as input to the LLM. The choice of embedding model is detailed in Section 4.3, while the retrieval configuration and similarity search specifications are described in Section 4.4.3.

Table 3 Reporting items required for the LLM-assisted title and abstract screening phase.

Item	Description	Reporting Purpose	Required
LLM Specifics	Identifier, version and parameters of the language model	Ensure LLM provenance and enable comparison of automated decisions across runs	Yes
Prompt template	Instruction template encoding the inclusion and exclusion criteria for screening	Support reproducibility of eligibility rules and related instructions	Yes
Output schema	Machine-readable response format (<i>e.g.</i> , JSON Schema or Pydantic)	Ensure consistent parsing and alignment of expected outputs	Yes
Fallback policy	Rules for fallback labels, recovery, and manual review routing	Specify handling of non-conforming outputs and possible escalation to human review	Yes
Batching strategy	Batch size and grouping logic for processing	Support efficient processing and consistency of results	Optional
Harmonization rules	Normalization rules for extracted fields	Document standardization and mapping of extracted fields	Optional
Provenance logging	Metadata recorded for screening decisions at record level	Enable traceability of decisions in the screening phase	Optional

Methodologically, the eligibility pipeline is organized into three conceptual steps, as illustrated in Figure 4:

1. Representation: Full-text articles are segmented into retrievable units, called chunks (*e.g.*, logical sections, paragraphs, or fixed-length segments), and encoded into dense vector embeddings. Embeddings are stored in a vector store alongside basic provenance metadata, such as the document identifier and chunk origin, enabling robust semantic retrieval that is insensitive to lexical variation.
2. Retrieval: Using the vector store populated in the previous step, a similarity-based retriever identifies the most relevant information using a retrieval query (*i.e.*, a natural language string or set of keywords derived from the inclusion and exclusion criteria). For this query, the retriever selects a bounded set of candidate chunks per document that are most likely to contain information relevant to the prespecified inclusion and exclusion criteria. Retrieval outputs include ranked relevance scores and explicit provenance pointers.
3. Generation: Based on the candidate chunks provided by the retriever, the LLM analyzes a prompt template that integrates the chunks with the inclusion and exclusion criteria. The model then produces structured responses that can be used to assign a categorical eligibility decision at the full-text level and to extract standardized study attributes that are useful for the review.

During execution of the eligibility pipeline, several constraints are enforced to maintain rigor and mitigate LLM limitations. LLM outputs are required to conform to a predefined schema (*e.g.*, a JSON schema specifying fields for decision and extracted attributes). Consistent with the screening phase, outputs are automatically validated against this schema, and any response that fails validation is flagged and routed to

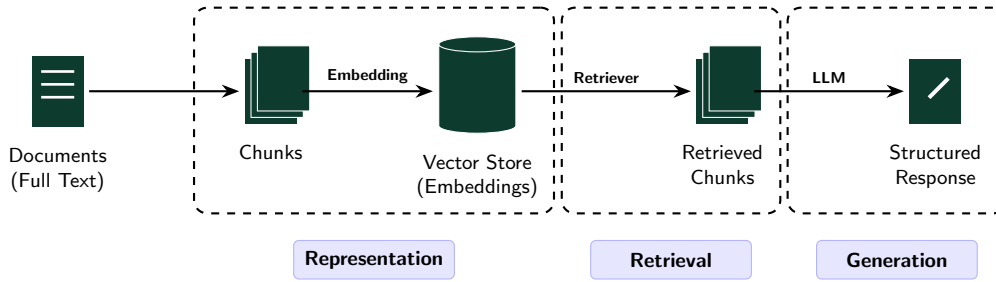


Fig. 4 Eligibility pipeline supporting full-text eligibility assessment under PRISMA workflow. For the sake of space, the Generation phase is shown here in a simplified form; its detailed architecture is provided in Figure 3.

Table 4 Additional reporting items for RAG-based full-text eligibility.

Item	Description	Reporting Purpose	Required
Embedding Specifics	Identifier, version, and parameters of the embedding model	Enable reproducibility of semantic representations and consistency of retrieval behavior	Yes
Document acquisition	Procedure for obtaining the full texts	Ensure provenance and traceability of source documents	Yes
Segmentation strategy	Method for dividing documents into chunks	Support appropriate retrieval granularity and contextual coverage	Yes
Retriever configuration	Set of parameters such as similarity metric and other filtering rules	Specify how candidate evidence chunks are selected for generation	Yes

human review in the following inclusion phase rather than automatically accepted. All operational parameters that could influence LLM behavior must be declared to ensure reproducibility. In the proposed framework, these parameters are reported as part of the workflow metadata; their specific values are detailed in Section 4.3.

The same robustness measures as in Section 3.2 are integrated at key steps of the pipeline — including retry policies, rate limiting, output validation, and conservative fallback routing — with the addition of document-level provenance tracking, which records the origin of each retrieved chunk at query time.

All reporting items listed in Table 3 also apply to the full-text eligibility assessment. Table 4 further specifies the additional items that govern the steps of the eligibility pipeline, as they define how full-text documents are acquired, segmented, embedded, and retrieved to provide the evidence base for automated eligibility decisions.

3.4 Inclusion

The inclusion phase is the final step of the proposed workflow and corresponds directly to the inclusion phase as defined in PRISMA 2020 [29]. It consolidates the set of

studies that have passed all screening and eligibility assessments and prepares them for data synthesis.

In the proposed framework, this phase takes as input the records that were assessed as eligible in the previous phase (*i.e.*, those for which the RAG-based pipeline produced a complete, validated, and harmonized structured response) and marks them as included in the systematic review. Records for which extraction was incomplete or schema validation failed are not automatically excluded; instead, they are flagged for manual verification and curation, in line with PRISMA 2020’s requirement for transparent reporting of the selection process and human oversight of uncertain cases.

In addition to confirming eligibility, the inclusion phase performs supplementary functions. First, it supports the integration of records from external sources, such as domain knowledge or complementary searches, provided that these records undergo the same screening and eligibility pipeline as the primary corpus. This is consistent with the PRISMA 2020 guidelines for reporting additional sources. Second, it performs a final manual validation of the attributes extracted incrementally during the screening and eligibility phases (*e.g.*, study identifiers, methods, results, or domain-specific fields) to ensure their correctness before synthesis. Third, it maintains full traceability for each included study, linking the original bibliographic metadata, full-text sources, retrieved chunks, and generative outputs into a complete provenance trail.

Therefore, the proposed framework does not alter the scope or criteria of the PRISMA 2020 inclusion phase, but operationalizes them through structured information extraction. This approach is a procedural refinement that supports automated processing while preserving full compliance with reporting requirements.

4 Experimental validation

To evaluate the quality of our proposal, we conduct an experimental validation within a specific research domain. This validation emphasizes the reproducibility of the proposed approach rather than the performance of its software implementation. The methodology aims to identify relevant papers and related datasets (as an example of specific information) in alignment with a traditional, manually curated process. The expected outcome need not exactly match manual review decisions; instead, it should achieve a comparable level of coverage and relevance in the resulting collection of studies and associated data.

In the remainder of this section, we present the reference study described in [23], including its workflow and results. Additionally, we replicate the methodology described by the authors in the original work to the best of our ability, using all available information from the study, and report those results as well. We then describe the technical implementation of our methodology and, finally, how we applied it to the same domain as the reference study. The results of this final step, along with their comparison to the reference study, are presented in Section 5.

```

TITLE-ABS-KEY ( "action quality assessment" OR "action
    assessment" OR "action quality evaluation" OR "human
    action evaluation" OR "movement quality assessment" )
AND ( LIMIT-TO ( DOCTYPE , "ar" ) OR LIMIT-TO ( DOCTYPE , "cp"
    ) )
AND ( LIMIT-TO ( SUBJAREA , "COMP" ) )
AND ( LIMIT-TO ( LANGUAGE , "English" ) )

```

Fig. 5 Search query of the reference study.

4.1 Reference case study

As a reference, we adopted the systematic review by Liu et al.[23]. This review focuses on AQA and examines aspects such as existing datasets, applications, equipment, models, and evaluation metrics. We selected this study because it provides a clearly defined object of interest, explicitly adheres to the PRISMA methodology, employs multiple established bibliographic databases, and reports detailed metadata that enable a meaningful comparison. Additionally, it is a recent peer-reviewed study, ensuring alignment with current literature and terminology. From a difficulty standpoint, AQA represents a moderately specialized domain with partially standardized terminology, making exhaustive retrieval non-trivial and therefore a meaningful test case for the proposed framework. For our evaluation, we focused on extracting and analyzing existing datasets as additional domain data. In the following paragraphs, we provide the details relevant to comparing our methodology. For further information, the reader is referred to the original paper.

The reference review relied on three major bibliographic databases: IEEE Xplore, Scopus, and Web of Science. Figure 5 reports the search query template adopted in the study, using Scopus as a representative example. Equivalent queries, adapted to the syntax of the respective platforms, were employed for the other databases. The search was executed on 16 July 2024, with no additional time constraints specified in the query formulation.

The selection process in the reference study was guided by the following inclusion criteria: (a) publication in a journal or conference proceedings, (b) related to computer science, and (c) published in English. The exclusion criteria comprised: (a) duplicated articles, (b) articles that presented results of surveys or reviews, (c) articles not related to using vision-based methods, and (d) articles not related to human AQA.

The results of the inclusion and exclusion steps of the PRISMA methodology are summarized in the first column of Table 5, directly extracted from [23]. At the end of the workflow, 96 studies were included in the systematic review. Among these, two were identified by the authors through a backward reference search of the selected papers. An internal verification subsequently revealed that these two papers corresponded to updated versions of studies already present in the set obtained from the three bibliographic databases. To ensure consistency and restrict comparisons to results derived from the three engines, these two additional articles were excluded from the reference set used as the benchmark for the evaluation. Consequently, the remaining 94 articles were considered as the reference corpus for evaluation.

From these studies, 47 AQA datasets were identified and extracted. Datasets explicitly reported as *self-created* in the reference study were excluded from the comparison, as they are not publicly available and cannot be independently reproduced. Both the selected papers and the corresponding datasets constitute the basis for comparison with the results produced by the proposed framework.

4.2 Replication of the reference case study

To establish a baseline, we replicated the identification phase of the reference study, adding temporal boundaries to account for the period elapsed since the original publication. The remaining phases were not replicated, as they rely on subjective human judgment that cannot be systematically reproduced from the information reported in the original study, a constraint inherent to all manual systematic reviews. The replicated search returned 244 documents, reduced to 151 after duplicate removal. These results are detailed in the second column of Table 5. As discussed in Section 5.1, minor differences with respect to the original figures are consistent with the known temporal variability of bibliographic database indices and do not compromise the validity of the comparison.

4.3 Framework implementation and toolchain

The experimental evaluation employed LLMs from the Gemini 2.5 family [36] for both screening and full-text eligibility phases. Gemini 2.5 Flash was selected as the primary inference engine for large-scale tasks due to its favorable trade-off between speed, cost, and output accuracy. Its use enables efficient processing of large collections of papers while maintaining adequate accuracy. Gemini 2.5 Pro was reserved for harmonization steps that require stronger reasoning capabilities and higher accuracy. The combination of these two LLMs supports a scalable architecture in which higher-capacity inference is applied only where it is most beneficial.

To ensure fine-grained control over response determinism and minimize variability across repeated executions, all LLM invocations were performed under decoding configurations designed to approximate deterministic behavior. For the selected LLMs, this control is primarily achieved through the temperature parameter, which governs randomness in token selection by modulating the probability distribution over candidate tokens; as temperature decreases, probability mass becomes more concentrated on high-likelihood tokens, reducing output variability. In the limiting case, setting $T = 0$ yields greedy decoding, in which the highest-probability token is selected at each generation step. Determinism was further reinforced through Top- k and Top- p constraints: Top- k sampling restricts token selection to the k most probable candidates, while Top- p (nucleus) sampling restricts selection to the smallest candidate set whose cumulative probability mass exceeds a threshold p . In our experiments, decoding parameters were set to $T = 0$, Top- $k = 1$, and Top- $p = 1$, substantially reducing variability in the generated outputs.

It is important to note, however, that strict determinism cannot be fully guaranteed in LLM deployments, even under greedy decoding. Factors such as distributed execution environments and limited control over internal random seeds may introduce

```

TITLE-ABS-KEY ("action quality" OR "action assessment" OR "
  action evaluation" OR AQA)
AND TITLE-ABS-KEY (dataset OR "data set" OR database OR "data
  catalogue" OR "data repository" OR "data sharing" OR "open
  data")
AND PUBYEAR > 2012 AND PUBYEAR < 2025
AND ( LIMIT-TO ( DOCTYPE , "ar" ) OR LIMIT-TO ( DOCTYPE , "cp"
  ) )
AND ( LIMIT-TO ( SUBJAREA , "COMP" ) )
AND ( LIMIT-TO ( LANGUAGE , "English" ) )

```

Fig. 6 Search query of the proposed framework.

minor run-to-run differences. Achieving full determinism would generally require local deployment with complete control over random seeds and execution environments, typically at the cost of access to larger-capacity models.

For the full-text eligibility assessment, documents were encoded in dense vector embeddings using the text-embedding-004 model [15]. This embedding model was selected for its ability to produce compact and semantically informative representations suitable for similarity-based retrieval. All embeddings were indexed in a Facebook AI Similarity Search (FAISS) vector store [10, 18], chosen due to its maturity, robustness, and ability to handle high-dimensional vector data with low latency, making it ideal for the workflow.

All experimental workflows were implemented on a personal computer equipped with 32 GB of RAM, a 1 TB SSD, and a 14-core Intel Core i7 CPU running at 3.1 GHz, with Windows 11 as the operating system. The codebase was developed in Python 3.11, chosen for its widespread adoption in the scientific community, extensive library ecosystem, and ease of customization. The entire workflow was orchestrated using LangChain as a unified framework. pandas was used for structured data handling, while pydantic was used to define validated input and output schemas. PDF documents were loaded using PyPDFLoader and segmented into chunks with RecursiveCharacterTextSplitter. Interactions with Gemini models were handled through the google.genai client library. Finally, text embeddings were generated using GoogleGenerativeAIEmbeddings and indexed with the faiss library.

4.4 Workflow execution in the case-study setting

This subsection describes the concrete execution of the proposed methodology within the experimental setting, explicitly mapping each step to the corresponding PRISMA phase in the reference study.

4.4.1 Identification

The proposed workflow starts from a query executed on the same database engines considered in the reference study. The query adopted in the proposed workflow is reported in Figure 6. Compared with the original query shown in Figure 5, our query is more permissive in referring to the research domain and is therefore expected to

retrieve a broader set of results. In particular, only minimal and widely used terms were included to cover the possible variations used to refer to AQA. This choice accounts for the fact that the terminology associated with the domain has progressively evolved and become standardized over time, whereas earlier works often adopted heterogeneous formulations.

In addition, the query explicitly requires the presence of terms related to data sources. This constraint reflects the objective of the case study, which targets the identification and analysis of AQA datasets: by anchoring the search to articles that mention the use, release, or discussion of data, the query filters out purely methodological contributions that would not yield relevant dataset information.

Particular attention was paid to the definition of the reference period. The temporal window was restricted to publications between 2013 and 2024 (inclusive). The lower bound was set for demonstration purposes, as no relevant articles were identified prior to 2013, and the first study included in the reference review also dates back to that year. The upper bound was introduced to limit the query to studies published up to the final year covered by the reference study. However, since the reference study executed its queries on 16 July 2024, all publications released after that date were excluded from the proposed workflow to ensure a fair and consistent comparison. For IEEE Xplore and Web of Science, this filtering was performed by exploiting the publication date metadata returned by the respective database engines. In contrast, Scopus does not consistently provide a reliable publication date field in the query metadata. Consequently, for Scopus results, a manual verification step was performed solely for comparison purposes, inspecting the publication dates reported in conference proceedings or journals. Given potential indexing delays affecting Scopus, a conservative strategy was adopted, excluding all papers published in July 2024 to avoid introducing a temporal advantage over the reference study. All other optional query constraints, including document type, subject area, and language, were kept consistent with those adopted in the reference study. No additional filtering criteria were introduced.

4.4.2 Screening

The automated screening and subsequent harmonization of the extracted data sources were executed sequentially. The screening stage was driven by the prompt reported in Figure 7. In the prompt, the screening decisions are encoded as *yes*, *no*, and *maybe* to simplify LLM output parsing, corresponding respectively to the *include*, *exclude*, and *uncertain* labels defined in Section 3.2. It is important to note that the inclusion and exclusion criteria are equivalent to those of the reference study. However, criteria referring to the identification phase, such as duplication verification, have been removed from the prompt to avoid redundancy with steps already performed upstream. On the other hand, an additional inclusion criterion has been introduced to retain only articles that explicitly mention the use or production of datasets, in line with the dataset-focused objective of this case study. The corpus of articles was processed in fixed-size batches of 15 records. Batching was chosen as a practical trade-off between accuracy and efficiency. For each batch, the pipeline instantiates the screening prompt by inserting the <Title> and <Abstract> fields of the corresponding articles, enforces

Task. Evaluate each record based on the following Inclusion and Exclusion Criteria and return ONLY JSON that matches the provided schema.

Inclusion Criteria:

- The article mentions the use or production of datasets, databases, or other data sources.
- The article is related to computer science.

Exclusion Criteria:

- The article states that the data are not available, must be requested from the authors or are not free (cost is required).
- The article presents results of surveys or reviews.
- The article is not related to using vision-based methods.
- The article is not related to Human Action Quality Assessment (AQA).

Constraint. ALL the requirements should be met.

Output Instructions:

- For each provided record, return an object with: id, decision ('yes'—'no'—'maybe'), reason, data_source.
- If title or abstract are missing or empty use decision='no' and reason='Title or abstract is missing', data_source=''.

Records to evaluate: <Title> and <Abstract> of each article.

Fig. 7 Prompt template used to instruct the LLM during the automated screening phase.

a rate-limit check before issuing the request, invokes the LLM, and then maps the parsed response back to the original records through the provided identifiers.

To comply with the rate limits imposed by the Google Gemini API under the free usage tier adopted for this experimental validation, the implementation enforces a simple rate limiter that allows at most 2 requests within a sliding window of 60 seconds. This constraint increases total processing time but does not affect the methodological validity of the results, and would be substantially reduced in production deployments with higher API quotas. LLM calls are only executed when the current call count falls below this threshold. LLM outputs must conform to a structured JSON format. The returned responses are then validated against the corresponding Pydantic schema. To improve robustness, an exponential backoff retry policy is applied on failures, with an initial delay of 10 seconds, a doubling strategy across retries, and a maximum number of attempts. If a batch cannot be successfully processed after all retries, the pipeline adopts a conservative fallback strategy: when responses cannot be parsed or are missing after all retry attempts, the affected records are assigned an indeterminate outcome (`decision = uncertain`).

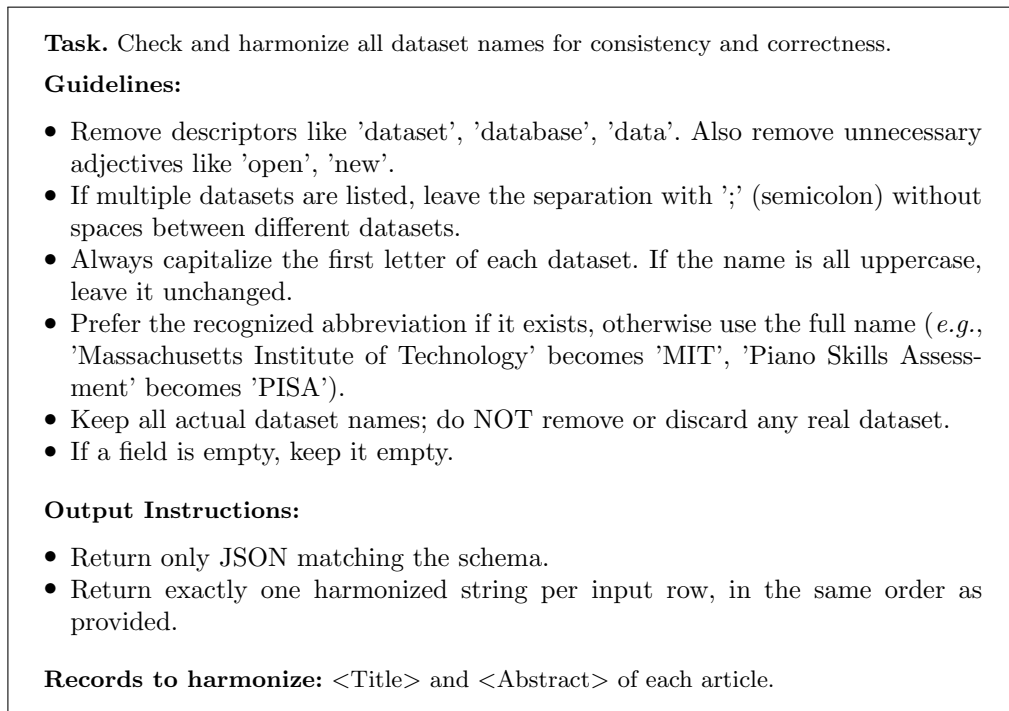


Fig. 8 Prompt template used to instruct the LLM to harmonize dataset names.

Following screening, dataset name harmonization is executed as a separate phase using the prompt shown in Figure 8. Unlike screening, this step does not employ batching; however, the same rate-limiting, retry, and output validation policies described above are applied.

4.4.3 Eligibility

The eligibility phase uses the same operational policies as screening (rate limiting, LLM model and decoding parameters, retry/backoff, and Pydantic validation). However, unlike screening, articles are processed individually at the document level rather than in fixed-size batches, and a RAG system is introduced. The use of RAG at this stage is motivated by the fact that scientific articles often exceed the effective context window of the LLM, and relevant information such as dataset names may be scattered across distant sections. Retrieving the most relevant chunks before prompting allows the generative step to focus on the portions most likely to contain the information needed for eligibility decisions. This process is applied only to articles for which no dataset was extracted during the screening phase, to improve computational efficiency and reduce processing costs. As a more comprehensive analysis, we suggest extending this procedure to all articles that successfully passed the screening phase.

Full texts are downloaded directly from the respective sources and loaded as PDFs into a local repository. PDFs are then cleaned and split using a character-based splitter

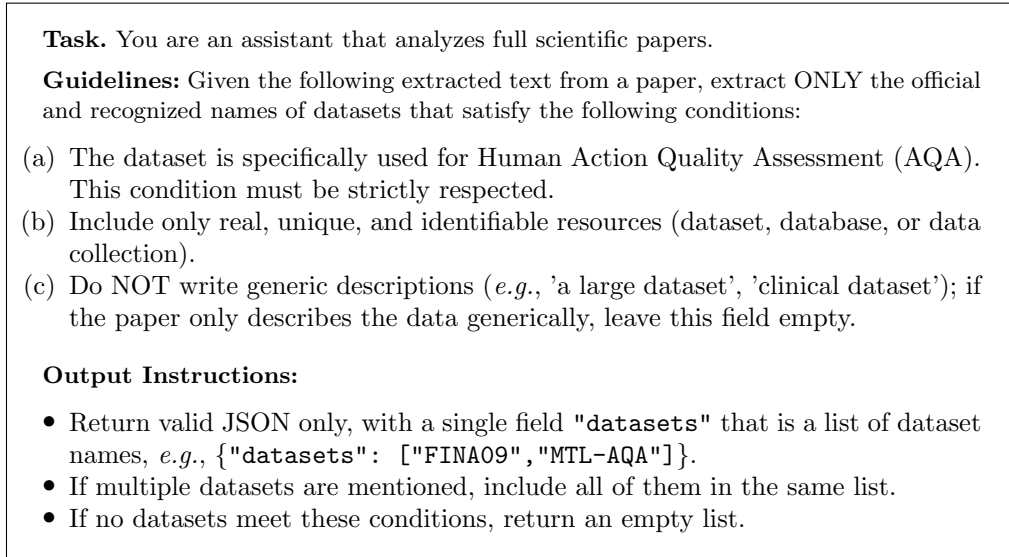


Fig. 9 Prompt template used to instruct the LLM to extract dataset names from the full text of scientific papers.

with a maximum chunk size of 2000 characters and an overlap of 400 characters, which helps preserve contextual continuity across adjacent chunks. Embeddings are computed with the configured embedding model and stored in a FAISS vector store, as explained previously.

For each document, the retrieval is driven by the fixed query "Human Action Quality Assessment (AQA) datasets", which is designed to focus on data sources relevant to AQA. A max-marginal-relevance (MMR) strategy is adopted rather than a standard similarity search [7, 19]. This approach is motivated by the need to balance relevance and diversity, reducing redundancy among chunks: MMR explicitly penalizes similar content and favors the selection of complementary information. The parameters that control the returned chunks are set proportionally to the total number of chunks, with $k = \lfloor N/3 \rfloor$ and $fetch_k = \lfloor N/2 \rfloor$, where N denotes the total number of chunks in the index. Here, $fetch_k$ defines the size of the initial candidate pool retrieved by semantic similarity search, while k specifies the number of chunks finally selected after MMR re-ranking. The trade-off between relevance and diversity is controlled by the MMR trade-off parameter λ . A high value ($\lambda = 0.9$) was adopted to strongly prioritize semantic relevance to the query while still introducing a small diversification effect.

For each paper, up to five candidate chunks are concatenated and passed to the LLM using the extraction prompt shown in Figure 9, with the records to harmonize adjusted for the full text. Extracted dataset names are then normalized with the same harmonization prompt shown in Figure 8.

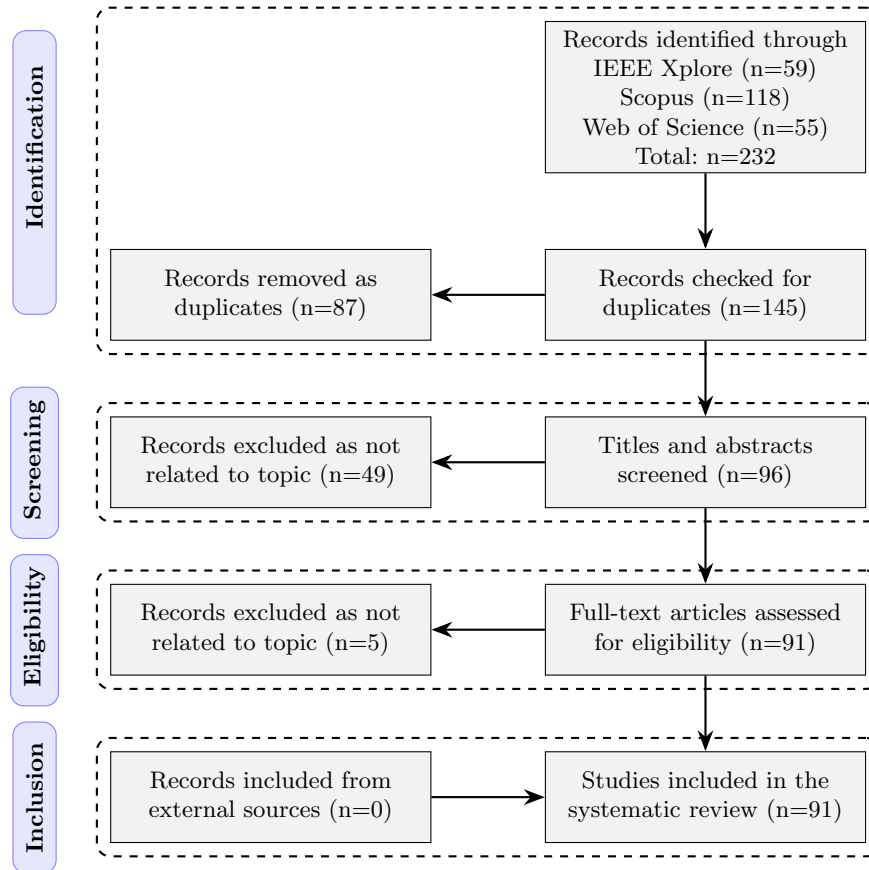


Fig. 10 Diagram of the PRISMA workflow for the proposed framework, where n denotes the number of papers.

4.4.4 Inclusion

Only articles for which at least one dataset was identified were considered eligible and included in the systematic review, forming the final corpus for subsequent analysis. Finally, Figure 10 shows the number of articles identified for each stage of the workflow. Specifically, a total of 91 articles were included in the systematic review. From these studies, 65 unique AQA datasets were extracted and are compared in the following section.

5 Results and discussion

The evaluation examines the framework’s ability to reproduce the reference study by retrieving relevant papers and associated datasets with coverage and relevance comparable to a manual review. These results are presented and discussed in two parts. First, we examine the selection of papers by comparing the set retrieved and processed by the framework with those included in the reference study. Second, we assess the

dataset selection to determine whether the automated pipeline captures information comparable to that of the manual process used in the original study. In all cases, exact matches are not required; rather, the goal is to achieve similar coverage and relevance to that of the manual review.

5.1 Results on study identification

The proposed framework ultimately identifies 91 articles for inclusion in the systematic review, 68 of which are also present among the 94 articles in the reference study. This corresponds to a detection rate of 72.34%, computed after excluding the two studies added through backward reference search in the reference set, as previously discussed. To account for the remaining gap, it is necessary to examine where the divergence occurs along the pipeline. As shown in Table 5, the screening and eligibility stages contribute only marginally to the loss: the screening stage correctly excludes 49 irrelevant articles while retaining all reference-study records that reach it, whereas the eligibility stage removes only two reference articles, neither of which contributes unique datasets. Overall, the automated stages retained 68 of the 70 reference studies that entered the pipeline, corresponding to a combined recall of 97.14%. Therefore, most of the discrepancy originates in the identification phase, *i.e.*, before any automated decision-making by the LLM and RAG components is applied.

This finding suggests that the LLM-based components of the framework operate with high fidelity, while the main limitation lies upstream in the bibliographic search stage. The missing records can be explained by factors intrinsic to the identification phase. First, and most importantly, the query was intentionally designed to retrieve only articles that explicitly refer to datasets, since the primary objective of the review is dataset identification rather than broad methodological coverage. As a consequence, articles discussing AQA methods that do not mention data sources in the title, abstract, or keywords are excluded by design. Second, terminological heterogeneity in the AQA literature, especially in earlier works, where non-standardized formulations were more common, may have led some records to fall outside the query’s lexical coverage. Third, the dynamic nature of bibliographic databases, including delayed indexing, metadata corrections, and changes in record availability over time, introduces a residual source of variability that cannot be fully controlled.

Despite the limitations observed in the identification phase, the overall detection rate remains broadly consistent with error rates reported for manual systematic reviews, as discussed in Section 1, where even experienced reviewers are subject to omissions. This suggests that the proposed framework can achieve comparable coverage while providing a more consistent and repeatable procedure, with lower operational effort than a fully manual workflow.

A comparison with our replication of the reference study further highlights the instability of bibliographic search results. Despite adopting the same declared queries, the replication produced different record counts at the identification stage. In particular, after duplicate removal, it yielded 27 more articles than the reference study. At the same time, neither corpus is a strict superset of the other: 92 of the 94 articles included in the final corpus of the reference study were also retrieved in the replicated

Table 5 Comparison of the PRISMA workflow between the reference study, the replicated study, and the proposed framework.

PRISMA Step	Reference Study	Replicated Study	Proposed Framework
Identification			
Records initially identified	230	244	232
IEEE Xplore	50	60	59
Scopus	119	131	118
Web of Science	61	54	55
Records removed as duplicates	105	93	87
Records after identification	125	151 (92)	145 (70)
Screening			
Records excluded during screening	24		49
Records after screening	101		96 (70)
Eligibility			
Records excluded during eligibility	7		5
Records after eligibility	94		91 (68)
Inclusion			
Records integrated from reference search	2		0
Final records included	96		91 (68)

Note: The replicated study reports only the Identification phase; empty cells denote steps not performed. Numbers in parentheses indicate the records in common with the final set of articles included in the reference study, excluding the two studies identified by the authors through backward reference search.

study, leaving a discrepancy of two articles. These differences confirm that the identification phase is the least stable and least reproducible step in a PRISMA-aligned workflow, whether manual or automated.

For completeness, the full list of the 91 articles identified by the proposed framework for inclusion in the systematic review is provided in [A.1](#). This table details the bibliographic references (DOI, year, title, and authors) along with the corresponding datasets extracted by the automated pipeline for each article.

5.2 Results on dataset identification

The proposed framework identified 40 of the 47 AQA datasets reported in the reference study, corresponding to a recall of 85.11%. A quantitative and qualitative synthesis of the discrepancies between the two studies is provided in [Table 6](#), which compares the 7 datasets missed by the framework with the 25 newly identified ones. To understand the nature of these differences, a manual expert validation was conducted by the authors to characterize both groups of datasets. Specifically, each resource was analyzed based on its relationship to the AQA domain, its availability status (public vs. private), and the inclusion status of related articles. Furthermore, the scientific impact was assessed by collecting citation counts from Google Scholar (as of December 20, 2025). This metric was chosen to evaluate whether the automated framework can retrieve influential resources that might have been overlooked in the original review, or vice versa. Crucially, this citation analysis was performed on two levels: we examined both the article identified during the workflow and the original publication where the dataset

Table 6 Summary comparison between datasets missed by the proposed framework and newly identified datasets.

Category	Feature	Missed by Framework (A.2)	Newly Identified (A.3)
Quantity	Total Datasets	7	25
Relation	AQA-ready	6	11
	AQA-related	1	14
Availability	Public	3	14
	Private	4	11
Inclusion	Found paper, missed dataset ^a	1 (14.3%)	8 (32.0%)
	Paper not found ^b	6 (85.7%)	17 (68.0%)
Impact	Total Citations (Articles)	197	967
	Avg. Citations (Articles)	28.14	38.68
	Total Citations (Original)	64	27,413
	Avg. Citations (Original)	21.33	1,713.31

Note: **a** The article is present in the study’s final collection, but the specific dataset was not extracted. For newly identified, it means the article was in the reference study, but the dataset name was omitted. **b** The article was not retrieved during the identification phase of the respective study/framework. The replicated version is used as the reference study.

was first introduced. This dual perspective enables a more robust evaluation of the framework’s ability to identify relevant sources.

A primary focus of this validation was the relation of each dataset to the domain. The resources have been classified into three levels of suitability:

- AQA-ready: Datasets fully prepared for AQA, containing ground truth quality scores, error labels, or clinical metrics required for action assessment;
- AQA-related: Datasets not natively AQA-ready but adaptable (*e.g.*, action recognition or motion analysis) with additional annotations;
- Not Applicable to AQA: Resources lacking the necessary quality variance or temporal granularity for being adaptable to AQA, belonging to an unaligned domain, or those not belonging to the dataset category (*e.g.*, methods or tools).

The manual inspection confirmed the high precision of both workflows, as no *“Not Applicable to AQA”* datasets were included in either the reference study or in the proposed framework. Notably, while almost all datasets from the reference study, not captured by our framework, were classified as *AQA-ready* (6 out of 7), the proposed framework introduced a more balanced distribution, identifying both *AQA-ready* (11) and *AQA-related* (14) resources.

Regarding dataset availability, the two systems performed similarly, extracting a comparable mix of public and private datasets. As shown in the summary table, most

undetected datasets (85.7%) are associated with articles that the pipeline failed to retrieve during the initial identification phase. Conversely, the automated pipeline identified a substantial number of additional resources not included in the original review. Notably, 32% of these newly identified datasets were extracted from articles already included in the reference study, but whose dataset names had been omitted from the mapping. Moreover, these newly identified datasets exhibit higher citation impact, particularly compared with the original papers. For a detailed breakdown of these findings, including individual metadata and citation counts, the reader is referred to [A: A.1](#) provides the full list of papers included in the systematic review, [A.2](#) lists the datasets missed by the framework, while [A.3](#) summarizes those newly identified.

6 Conclusions

This work demonstrates that LLMs and RAG can be integrated into a PRISMA-aligned workflow to automate systematic review tasks without sacrificing transparency or reproducibility. The experimental validation showed that the primary bottleneck lies not in the automated components, but in the bibliographic identification phase, a finding that has direct implications for the design of future automated review pipelines.

In the AQA case study used for validation, the framework retrieved 91 articles, including 68 of the 94 studies reported in the reference review (72.34% detection rate). The screening and eligibility components introduced only a marginal loss (97.14% combined recall) on the reference studies that reached the automated pipeline. At the dataset level, it extracted 65 distinct datasets, including 40 of the 47 reference datasets (85.11% extraction accuracy) and 25 additional datasets not reported in the original review. The newly discovered datasets further highlight the framework’s potential to surface relevant information not captured in the reference review.

This study has a number of limitations that should be noted. The framework was validated on a single domain (AQA), and it remains unclear how well it would transfer to fields with substantially different terminological conventions, publication structures, or data annotation practices. Moreover, the observed detection rate means that over a quarter of the reference articles were missed, largely as a consequence of the dataset-oriented query design adopted in the case study. Broader query formulations could improve recall, though likely at the expense of precision. On the modeling side, all experiments relied on a single LLM family (Gemini 2.5) under one configuration, and the sensitivity of the results to alternative models or pipeline components has not been assessed. It is also worth noting that the 25 newly identified datasets were validated by the authors themselves, without involving independent annotators. Lastly, no quantitative evidence on time savings or workload reduction was collected, as the evaluation was deliberately based on coverage and reproducibility rather than computational performance.

These results show that the framework can support systematic reviews without compromising methodological rigor and provides a complementary tool for evidence synthesis and targeted information retrieval. Beyond this case-study validation, the proposed pipeline is intended to be adaptable across domains by re-specifying the query template, eligibility criteria, and extraction schema while preserving the same

PRISMA-aligned process. Future work will target three main directions. First, the bibliographic identification phase will be refined through query expansion and the integration of additional bibliographic engines to improve recall. Second, the RAG-based retrieval and metadata harmonization components will be optimized and additional LLMs will be benchmarked across detection performance, extraction quality, and computational efficiency. Third, the framework will be validated across diverse research domains, which is essential to substantiate the transferability of the proposed approach. By making systematic review procedures more accessible, reproducible, and scalable, the proposed framework aims to lower the barrier for evidence synthesis in research domains where manual reviews remain prohibitively resource-intensive.

Declarations

- **Funding:** The authors declare that no funds, grants, or other support were received during the preparation of this manuscript.
- **Competing interests:** The authors have no competing interests to declare that are relevant to the content of this article.
- **Ethics approval and consent to participate:** Not applicable.
- **Consent for publication:** Not applicable.
- **Data availability:** The datasets supporting the results of this study will be made publicly available upon acceptance of the manuscript.
- **Materials availability:** Not applicable.
- **Code availability:** The code supporting the results of this study will be made publicly available upon acceptance of the manuscript.
- **Author contributions:** Omitted for double-blind review. Full author contributions are provided in the separate Title Page.

References

- [1] Bastian H, Glasziou P, Chalmers I (2010) Seventy-five trials and eleven systematic reviews a day: how will we ever keep up? *PLoS medicine* 7(9):e1000326
- [2] Bolanos F, Salatino A, Osborne F, et al (2024) Artificial intelligence for literature reviews: opportunities and challenges: F. bolaños et al. *Artificial Intelligence Review* 57(10):259
- [3] Borah R, Brown AW, Capers PL, et al (2017) Analysis of the time and workers needed to conduct systematic reviews of medical interventions using data from the prospero registry. *BMJ open* 7(2):e012545
- [4] Brîncoveanu C, Carl KV, Witzki A, et al (2026) Augmenting systematic literature reviews: A human-ai collaborative framework. In: *German Conference on Artificial Intelligence (Künstliche Intelligenz)*, Springer, pp 3–17
- [5] Brown A, Roman M, Devereux B (2025) A systematic literature review of retrieval-augmented generation: Techniques, metrics, and challenges. *arXiv*

- [6] Cao C, Arora R, et al (2025) Automation of systematic reviews with large language models. medRxiv pp 2025–06
- [7] Carbonell J, Goldstein J (1998) The use of mmr, diversity-based reranking for reordering documents and producing summaries. In: Proceedings of the 21st annual international ACM SIGIR conference on Research and development in information retrieval, pp 335–336
- [8] Clark J, Barton B, Albarqouni L, et al (2025) Generative artificial intelligence use in evidence synthesis: A systematic review. Research Synthesis Methods pp 1–19
- [9] Delgado-Chaves FM, Jennings MJ, Atalaia A, et al (2025) Transforming literature screening: The emerging role of large language models in systematic reviews. Proceedings of the National Academy of Sciences 122(2):e2411962122
- [10] Douze M, Guzhva A, Deng C, et al (2025) The Faiss library. IEEE Transactions on Big Data
- [11] Flemyng E, Noel-Storr A, Macura B, et al (2025) Position statement on artificial intelligence (ai) use in evidence synthesis across cochrane, the campbell collaboration, jbi and the collaboration for environmental evidence 2025. Cochrane Database of Systematic Reviews (10)
- [12] Forero DA, Abreu SE, Tovar BE, et al (2025) Large language models and the analyses of adherence to reporting guidelines in systematic reviews and overviews of reviews (PRISMA 2020 and PRIOR). Journal of Medical Systems 49(1):80
- [13] Galli C, Gavrilova AV, Calciolari E (2025) Large language models in systematic review screening: opportunities, challenges, and methodological considerations. Information 16(5):378
- [14] Gao Y, Xiong Y, Gao X, et al (2023) Retrieval-augmented generation for large language models: A survey. arXiv preprint arXiv:231210997 2(1)
- [15] Google (2024) Text-Embedding-004: latest english text embedding model. URL <https://cloud.google.com/vertex-ai/generative-ai/docs/model-reference/text-embeddings>
- [16] Harasgama S, Pearce H, Appel C, et al (2026) Artificial intelligence tools for automating evidence synthesis: Scoping review. Journal of Medical Internet Research 28:e81597
- [17] Holst D, Moenck K, Koch J, et al (2025) Transparent reporting of ai in systematic literature reviews: Development of the prisma-traice checklist. JMIR AI 4:e80247

- [18] Johnson J, Douze M, Jégou H (2019) Billion-scale similarity search with GPUs. *IEEE Transactions on Big Data* 7(3):535–547
- [19] Karpukhin V, Oguz B, Min S, et al (2020) Dense passage retrieval for open-domain question answering. In: *Proceedings of the 2020 conference on empirical methods in natural language processing (EMNLP)*, pp 6769–6781
- [20] Khalil H, Ameen D, Zarnegar A (2022) Tools to support the automation of systematic reviews: a scoping review. *Journal of Clinical Epidemiology* 144:22–42
- [21] Khraisha Q, Put S, Kappenberg J, et al (2024) Can large language models replace humans in systematic reviews? evaluating gpt-4’s efficacy in screening and extracting data from peer-reviewed and grey literature in multiple languages. *Research Synthesis Methods* 15(4):616–626
- [22] Lewis P, Perez E, Piktus A, et al (2020) Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in neural information processing systems* 33:9459–9474
- [23] Liu J, Wang H, Stawarz K, et al (2025) Vision-based human action quality assessment: A systematic review. *Expert Systems with Applications* 263:125642
- [24] Malik FS, Terzidis O (2025) A hybrid framework for creating artificial intelligence-augmented systematic literature reviews. *Management Review Quarterly* pp 1–27
- [25] Marshall IJ, Kuiper J, Wallace BC (2016) RobotReviewer: evaluation of a system for automatically assessing bias in clinical trials. *Journal of the American Medical Informatics Association* 23(1):193–201
- [26] Miao Y, Zhao Y, Luo Y, et al (2025) Improving large language model applications in the medical and nursing domains with retrieval-augmented generation: Scoping review. *Journal of Medical Internet Research* 27:e80557
- [27] Moher D, Liberati A, Tetzlaff J, et al (2009) Preferred reporting items for systematic reviews and meta-analyses: the PRISMA statement. *Bmj* 339
- [28] Oami T, Okada Y, Nakada Ta (2024) Performance of a large language model in screening citations. *JAMA network open* 7(7):e2420496–e2420496
- [29] Page MJ, McKenzie, et al (2021) The PRISMA 2020 statement: an updated guideline for reporting systematic reviews. *BMJ* 372
- [30] Pérez J, Díaz J, Garcia-Martin J, et al (2020) Systematic literature reviews in software engineering—enhancement of the study selection process using cohen’s kappa statistic. *Journal of Systems and Software* 168:110657
- [31] Rethlefsen ML, Kirtley S, Waffenschmidt S, et al (2021) PRISMA-S: an extension to the PRISMA statement for reporting literature searches in systematic reviews.

- [32] Scherbakov D, Hubig N, Jansari V, et al (2025) The emergence of large language models as tools in literature reviews: a large language model-assisted systematic review. *Journal of the American Medical Informatics Association* 32(6):1071–1086
- [33] Scotti KL, Young S, Gainey MA, et al (2025) Artificial intelligence and automation in evidence synthesis: An investigation of methods employed in cochrane, campbell collaboration, and environmental evidence reviews. *Cochrane Evidence Synthesis and Methods* 3(5):e70046
- [34] Shailendra S, Kadel R, Sharma A, et al (2026) L-prisma: An extension of prisma in the era of generative artificial intelligence (genai). *Authorea Preprints*
- [35] Susnjak T (2023) Prisma-dflm: An extension of prisma for systematic literature reviews using domain-specific finetuned large language models. *arXiv preprint arXiv:230614905*
- [36] Team G, Anil R, Borgeaud S, et al (2023) Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:231211805*
- [37] Thode L, Iftikhar U, Mendez D (2025) Exploring the use of llms for the selection phase in systematic literature studies. *Information and Software Technology* p 107757
- [38] de la Torre-López J, Ramírez A, Romero JR (2023) Artificial intelligence to automate the systematic review of scientific literature. *Computing* 105(10):2171–2194
- [39] Van De Schoot R, Bruin D, et al (2021) An open source machine learning framework for efficient and transparent systematic reviews. *Nature machine intelligence* 3(2):125–133
- [40] Wagner G, Lukyanenko R, Paré G (2022) Artificial intelligence and the conduct of literature reviews. *Journal of Information Technology* 37(2):209–226
- [41] Wang Z, Nayfeh T, Tetzlaff J, et al (2020) Error rates of human reviewers during abstract screening in systematic reviews. *PloS one* 15(1):e0227742
- [42] Wang Z, Cao L, Danek B, et al (2025) Accelerating clinical evidence synthesis with large language models. *npj Digital Medicine* 8(1):509
- [43] Zhao WX, Zhou, et al (2023) A survey of large language models. *arXiv preprint arXiv:230318223* 1(2)

Appendix A

Table A.1: Overview of all identified papers with all datasets extracted.

DOI	Year	Title	Authors	Data Source
10.1007/ 978-3-319-10599-4_36	2014	Assessing the quality of actions	Pirsiavash, H.; Vondrick, C.; Torralba, A.	MIT-Skate; MIT-Dive
10.1109/CVPRW. 2017.16	2017	Learning to Score Olympic Events	P. Parmar; B. T. Morris	FINA09
10.1109/SSCI.2017. 8285270	2017	The open online repository of karate motion capture data: A tool for scientists and sport educators	Hachaj, T.; Ogiela, M.R.; Piekarczyk, M.	self-created
10.1007/ 978-3-030-00767-6_12	2018	End-to-end learning for action quality assessment	Li, Y.; Chai, X.; Chen, X.	MIT-Skate; UNLV-Dive; UNLV-Vault
10.1016/j.patcog.2017. 12.007	2018	Motion analysis: Action detection, recognition and evaluation based on motion capture data	Patrona, F.; Chatzitofis, A.; Zarpalas, D.; Daras, P.	MSR-Action3D; MSRC-12
10.1109/ICIP.2018. 8451364	2018	S3D: Stacking Segmental P3D for Action Quality Assessment	X. Xiang; Y. Tian; A. Reiter; G. D. Hager; T. D. Tran	UNLV-Dive
10.1109/WACV.2019. 00161	2019	Action Quality Assessment Across Multiple Actions	P. Parmar; B. Morris	AQA-7
10.1109/CVPR.2019. 00039	2019	What and How Well You Performed? A Multitask Learning Approach to Action Quality Assessment	P. Parmar; B. T. Morris	MTL-AQA
10.1007/ 978-3-030-58577-8_14	2020	An Asymmetric Modeling for Action Assessment	Gao, J.; Zheng, W.-S.; Pan, J.-H.; Gao, C.; Wang, Y.-W.; Zeng, W.; Lai, J.-H.	JIGSAWS; TASD-2; AQA-7
10.1109/ ICMEW46912.2020. 9106049	2020	Efficient Fitness Action Analysis Based on Spatio-Temporal Feature Encoding	J. Li; H. Cui; T. Guo; Q. Hu; Y. Shen	Fitness-AQA
10.1145/3394171. 3413560	2020	Hybrid Dynamic-static Context-aware Attention Network for Action Assessment in Long Videos	Zeng, L.-A.; Hong, F.-T.; Zheng, W.-S.; Yu, Q.-Z.; Zeng, W.; Wang, Y.-W.; Lai, J.-H.	Rhythmic Gymnastics
10.3390/ electronics9040568	2020	Learning effective skeletal representations on RGB video for fine-grained human action quality assessment	Lei, Q.; Zhang, H.; Du, J.; Hsiao, T.-C.; Chen, C.-C.	MIT-Skate; MIT-Dive; UNLV-Skate; UNLV-Dive; UNLV-Vault
10.1109/ TransAI49837.2020. 00030	2020	Skeleton-Based Detection of Abnormalities in Human Actions Using Graph Convolutional Networks	B. X. B. Yu; Y. Liu; K. C. C. Chan	UI-PRMD
10.1145/3341105. 3374092	2020	Towards a data-driven method for RGB video-based hand action quality assessment in real time	Wang, T.; Jin, M.; Wang, J.; Wang, Y.; Li, M.	Origami Video
10.1109/CVPR42600. 2020.00986	2020	Uncertainty-Aware Score Distribution Learning for Action Quality Assessment	Y. Tang; Z. Ni; J. Zhou; D. Zhang; J. Lu; Y. Wu; J. Zhou	AQA-7; MTL-AQA

Continued on next page

Table A.1 – Continued from previous page

DOI	Year	Title	Authors	Data Source
10.1109/TCSVT.2020.3017727	2021	Action Quality Assessment Using Siamese Network-Based Deep Metric Learning	Jain, Hiteshi; Harit, Gaurav; Sharma, Avinash	MIT-Dive; MTL-AQA; UNLV-Dive
10.1109/WACV48630.2021.00044	2021	EAGLE-Eye: Extreme-pose Action Grader using detail bird's-Eye view	M. Nekoui; F. O. Tito Cruz; L. Cheng	AQA-7; ExPose; MIT-Skate
10.1109/ICCV48922.2021.00782	2021	Group-aware Contrastive Regression for Action Quality Assessment	X. Yu; Y. Rao; W. Zhao; J. Lu; J. Zhou	AQA-7; MTL-AQA; JIGSAWS
10.1145/3462244.3479891	2021	Improving the Movement Synchrony Estimation with Action Quality Assessment in Children Play Therapy	Li, J.; Bhat, A.; Barmaki, R.L.	PT13
10.1016/j.knosys.2021.107388	2021	Learning and fusing multiple hidden substages for action quality assessment	Dong, L.-J.; Zhang, H.; Shi, Q.; Lei, Q.; Du, J.; Gao, S.	UNLV-Dive
10.1109/MMSP53017.2021.9733638	2021	Piano Skills Assessment	P. Parmar; J. Reddy; B. Morris	PISA
10.1109/ITME53901.2021.00048	2021	Skeleton Based Action Quality Assessment of Figure Skating Videos	H. -Y. Li; Q. Lei; H. -B. Zhang; J. -X. Du	MIT-Skate
10.1016/j.patcog.2021.108095	2021	Skeleton-based human action evaluation using graph convolutional network for monitoring Alzheimer's progression	Yu, B.X.B.; Liu, Y.; Chan, K.C.C.; Yang, Q.; Wang, X.	UI-PRMD
10.1007/s11263-021-01486-4	2021	SportsCap: Monocular 3D Human Motion Capture and Fine-Grained Understanding in Challenging Sports Videos	Chen, X.; Pang, A.; Yang, W.; Ma, Y.; Xu, L.; Yu, J.	AQA-7; Diving48; FineGym; MTL-AQA; SMART
10.1007/s11760-021-01890-w	2021	Temporal attention learning for action quality assessment in sports video	Lei, Q.; Zhang, H.; Du, J.	AQA-7
10.1145/3453892.3461624	2021	Towards Improved and Interpretable Action Quality Assessment with Self-Supervised Alignment	Roditakis, K.; Makris, A.; Argyros, A.	MTL-AQA
10.1145/3474085.3475438	2021	TSA-Net: Tube Self-Attention Network for Action Quality Assessment	Wang, S.; Yang, D.; Zhai, P.; Chen, C.; Zhang, L.	AQA-7; MTL-AQA
10.1016/j.jvcir.2021.103304	2021	What and how well you exercised? An efficient analysis framework for fitness actions	Li, J.; Hu, Q.; Guo, T.; Wang, S.; Shen, Y.	Fitness-28
10.1007/978-3-031-19772-7_25	2022	Action Quality Assessment with Temporal Parsing Transformer	Bai, Y.; Zhou, D.; Zhang, S.; Wang, J.; Ding, E.; Guan, Y.; Long, Y.; Wang, J.	AQA-7; JIGSAWS; MTL-AQA
10.3390/electronics11193051	2022	An Efficient Motion Registration Method Based on Self-Coordination and Self-Referential Normalization	Ren, Y.; Zhang, B.; Chen, J.; Guo, L.; Wang, J.	KTH
10.1155/2022/9402195	2022	Design of a Professional Sports Competition Adjudication System Based on Data Analysis and Action Recognition Algorithm	Lin, Q.; Zou, J.	Northwestern-UCLA; MSRC-12; CAD-60

Continued on next page

Table A.1 – Continued from previous page

DOI	Year	Title	Authors	Data Source
10.1007/978-3-031-19839-7.7	2022	Domain Knowledge-Informed Self-supervised Representations for Workout Form Assessment	Parmar, P.; Gharat, A.; Rhodin, H.	Fitness-AQA
10.1109/CVPR52688.2022.00296	2022	FineDiving: A Fine-grained Dataset for Procedure-aware Action Quality Assessment	J. Xu; Y. Rao; X. Yu; G. Chen; J. Zhou; J. Lu	FineDiving
10.1038/s41597-022-01188-7	2022	Functional movement screen dataset collected with two Azure Kinect depth sensors	Xing, Qing-Jun; Shen, Yuan-Yuan; Cao, Run; Zong, Shou-Xin; Zhao, Shu-Xiang; Shen, Yan-Fei	self-created
10.1007/978-3-031-04881-4.46	2022	Improving Action Quality Assessment Using Weighted Aggregation	Farabi, S.; Himel, H.; Gazzali, F.; Hasan, M.B.; Kabir, M.H.; Farazi, M.	MTL-AQA
10.1016/j.patrec.2022.04.015	2022	Learning time-aware features for action quality assessment	Zhang, Y.; Xiong, W.; Mi, S.	MTL-AQA
10.1109/CVPR52688.2022.00323	2022	Likert Scoring with Grade Decoupling for Long-term Action Assessment	A. Xu; L. -A. Zeng; W. -S. Zheng	FIS-V; JIGSAWS; Rhythmic Gymnastics
10.1007/978-3-031-19772-7.27	2022	Pairwise Contrastive Learning Network for Action Quality Assessment	Li, M.-Z.; Zhang, H.; Lei, Q.; Fan, Z.; Liu, J.; Du, J.	AQA-7; MTL-AQA
10.1109/ISCAS48785.2022.9937262	2022	RARN: A Real-Time Skeleton-based Action Recognition Network for Auxiliary Rehabilitation Therapy	M. Shen; H. Lu	RDSD
10.1109/ICIP46576.2022.9897932	2022	Representation Learning Using Rank Loss for Robust Neurosurgical Skills Evaluation	B. Baby; M. Chasmai; T. Banerjee; A. Suri; S. Banerjee; C. Arora	JIGSAWS; NETS
10.1155/2022/5430463	2022	Research on Tennis Motion Evaluation Method and System Based on Deep Learning	Chang, Huan	MSR-Action3D
10.1109/TCSVT.2022.3143549	2022	Semi-Supervised Action Quality Assessment With Self-Supervised Segment Feature Recovery	S. -J. Zhang; J. -H. Pan; J. Gao; W. -S. Zheng	MTL-AQA; Rhythmic Gymnastics
10.1016/j.jvcir.2022.103625	2022	Skeleton-based deep pose feature learning for action quality assessment on figure skating videos	Li, H.; Lei, Q.; Zhang, H.; Du, J.; Gao, S.	MIT-Skate; FIS-V
10.1007/978-3-031-18913-5.17	2022	Skeleton-Based Action Quality Assessment via Partially Connected LSTM with Triplet Losses	Wang, X.; Li, J.; Hu, H.	UMONS-TAICHI; Walking Gait
10.1109/MMSP55362.2022.9949464	2022	Tai Chi Action Quality Assessment and Visual Analysis with a Consumer RGB-D Camera	J. Li; H. Hu; Q. Xing; X. Wang; J. Li; Y. Shen	TaiChi-24
10.1145/3581783.3613774	2023	A Figure Skating Jumping Dataset for Replay-Guided Action Quality Assessment	Liu, Y.; Cheng, X.; Ikenaga, T.	RFSJ
10.1007/978-3-031-44216-2.19	2023	A Graph Convolutional Siamese Network for the Assessment and Recognition of Physical Rehabilitation Exercises	Li, C.; Ling, X.; Xia, S.	UI-PRMD; IRDS

Continued on next page

Table A.1 – Continued from previous page

DOI	Year	Title	Authors	Data Source
10.1117/12.2685368	2023	A novel blind action quality assessment based on Multi-headed GRU network and attention mechanism	Sun, W.; Hu, Y.; Zhang, B.; Chen, X.; Hao, C.; Gao, Y.	AQA-7; JIGSAWS
10.1155/2023/3649217	2023	A Novel Model for Intelligent Pull-Ups Test Based on Key Point Estimation of Human Body and Equipment	Liu, G.; Wang, J.; Zhang, Z.; Liu, Q.; Ren, Y.; Zhang, M.; Chen, S.; Bai, P.	self-created
10.1109/EIECC60864.2023.10456639	2023	A Temporal Action Evaluation Algorithm for Medical Clinical Skill Operations	Z. Wang; J. Ruan; J. Zhang; J. Yue	Thumos14; ActivityNet1.3
10.1109/TVCG.2023.3247092	2023	A Video-Based Augmented Reality System for Human-in-the-Loop Muscle Strength Assessment of Juvenile Dermatomyositis	K. Zhou; R. Cai; Y. Ma; Q. Tan; X. Wang; J. Li; H. P. H. Shum; F. W. B. Li; S. Jin; X. Liang	JDM; 3D Animation
10.1109/ICMLC58545.2023.10327994	2023	Action Quality Assessment for ASD Behaviour Evaluation	D. Zhang; D. Zhou; H. Liu	AQA-7; BEST; DREAM; EPIC-Skills; FIS-V; Infinite Grasp; JIGSAWS; MIT-Dive
10.1109/ACCESS.2023.3316009	2023	AI Trainer: Autoencoder Based Approach for Squat Analysis and Correction	M. Chariar; S. Rao; A. Irani; S. Suresh; C. S. Asha	Countix; Fitness-AQA; Kinetics 700; MM-Fit; UCF-101
10.1109/TNSRE.2023.3317411	2023	A Contrastive Learning Network for Performance Metric and Assessment of Physical Rehabilitation Exercises	L. Yao; Q. Lei; H. Zhang; J. Du; S. Gao	UI-PRMD; KIMORE
10.1007/s00521-023-09068-w	2023	Auto-encoding score distribution regression for action quality assessment	Zhang, Boyu; Chen, Jiayuan; Xu, Yinfei; Zhang, Hui; Yang, Xu; Geng, Xin	AQA-7; MTL-AQA; JIGSAWS
10.1007/s11263-022-01695-5	2023	Automatic Modelling for Interactive Action Assessment	Gao, J.; Pan, J.-H.; Zhang, S.-J.; Zheng, W.-S.	JIGSAWS; TASD-2; PaSk; AQA-7
10.1109/ACCESS.2023.3305372	2023	Contrastive Learning for Action Assessment Using Graph Convolutional Networks With Augmented Virtual Joints	C. -I. Jeong; S. Byun; S. Baek	AI-Hub Fitness; Squat
10.1109/TIP.2023.3331212	2023	Fine-Grained Spatio-Temporal Parsing Network for Action Quality Assessment	K. Gedamu; Y. Ji; Y. Yang; J. Shao; H. T. Shen	FineDiving; AQA-7; MTL-AQA
10.1007/s40747-022-00892-6	2023	Gaussian guided frame sequence encoder network for action quality assessment	M.-Z. Li; Zhang, H.; Dong, L.-J.; Lei, Q.; Du, J.	AQA-7; MTL-AQA
10.1109/TCSVT.2023.3281413	2023	Hierarchical Graph Convolutional Networks for Action Quality Assessment	K. Zhou; Y. Ma; H. P. H. Shum; X. Liang	AQA-7; MTL-AQA; JIGSAWS
10.1007/s10489-023-05166-3	2023	Improving action quality assessment with across-staged temporal reasoning on imbalanced data	Lian, P.-X.; Shao, Z.-G.	FineDiving
10.1007/s10489-022-03984-5	2023	Label-reconstruction-based pseudo-subscore learning for action quality assessment in sporting events	Zhang, H.; Dong, L.-J.; Lei, Q.; Yang, L.-J.; Du, J.	AQA-7; MTL-AQA

Continued on next page

Table A.1 – Continued from previous page

DOI	Year	Title	Authors	Data Source
10.1145/3581783.3613795	2023	Localization-assisted Uncertainty Score Disentanglement Network for Action Quality Assessment	Ji, Y.; Ye, L.; Huang, H.; Mao, L.; Zhou, Y.; Gao, L.	FineFS
10.1109/CVPR52729.2023.00238	2023	LOGO: A Long-Form Video Dataset for Group Action Quality Assessment	S. Zhang; W. Dai; S. Wang; X. Shen; J. Lu; J. Zhou; Y. Tang	LOGO
10.1145/3650400.3650642	2023	Martial arts Scoring System based on U-shaped networkWushu intelligent scoring systemLearning to score Chinese Wushu	Li, M.; Tian, F.; Li, Y.	AGF-Olympics; FIS-V; Fri2023; LOGO; MTL-AQA
10.3390/systems11010021	2023	MLA-LSTM: A Local and Global Location Attention LSTM Learning Model for Scoring Figure Skating	Han, C.; Shen, F.; Chen, L.; Lian, X.; Gou, H.; Gao, H.	FIS-V; MIT-Skate
10.1145/3577190.3614117	2023	MMASD: A Multimodal Dataset for Autism Intervention Analysis	Li, J.; Chheang, V.; Kullu, P.; Brignac, E.; Guo, Z.; Bhat, A.; Barner, K.E.; Barmaki, R.L.	MMASD
10.1109/MMSP59012.2023.10337711	2023	MR-STGN: Multi-Residual Spatio Temporal Graph Network Using Attention Fusion for Patient Action Assessment	Y. Mourchid; R. Slama	UI-PRMD
10.1145/3622896.3622916	2023	Multi-Stage Action Quality Assessment Method	Liu, L.; Zhai, P.; Zheng, D.; Fang, Y.	FineDiving
10.1007/s10489-023-04613-5	2023	Multi-skeleton structures graph convolutional network for action quality assessment in long videos	Lei, Q.; Li, H.; Zhang, H.; Du, J.; Gao, S.	MIT-Skate; Rhythmic Gymnastics
10.1145/3606038.3616150	2023	Video-based Skill Assessment for Golf: Estimating Golf Handicap	Ingwersen, C.K.; Xarles, A.; Clapés, A.; Madadi, M.; Jensen, J.N.; Hannemose, M.R.; Dahl, A.B.; Escalera, S.	Golf Swing Video
10.1109/TMM.2022.3222681	2023	View-Invariant Center-of-Pressure Metrics Estimation With Monocular RGB Camera	C. Du; S. Graham; C. Depp; T. Nguyen	self-created
10.26555/ijain.v9i1.919	2023	Who danced better? Ranked tiktok dance video dataset and pairwise action quality assessment method	Hipiny, I.; Ujir, H.; Alias, A.A.; Shanat, M.; Ishak, M.K.	Ranked TikTok
10.1109/ICAACE61206.2024.10548995	2024	Action Quality Assessment with Multi-scale Temporal Attention Mechanism	W. Wang; H. Wang; Y. Hao; Q. Wang	AQA-7; MTL-AQA
10.1109/TMM.2023.3294800	2024	Adaptive Stage-Aware Assessment Skill Transfer for Skill Determination	S. -J. Zhang; J. -H. Pan; J. Gao; W. -S. Zheng	AQA-7; BEST; EPIC-Skills; JIGSAWS
10.1007/s10489-024-05349-6	2024	Assessing action quality with semantic-sequence performance regression and densely distributed sample weighting	Huang, F.; Li, J.	UNLV-Dive; AQA-7
10.1109/TPAMI.2024.3378753	2024	EGCN++: A New Fusion Strategy for Ensemble Learning in Skeleton-Based Rehabilitation Exercise Assessment	Yu, B.X.B.; Liu, Y.; Chan, K.C.C.; Chen, C.W.	UI-PRMD; KIMORE; EHE

Continued on next page

Table A.1 – Continued from previous page

DOI	Year	Title	Authors	Data Source
10.1109/CVPR52733.2024.01386	2024	FineParser: A Fine-grained Spatio-temporal Action Parser for Human-centric Action Quality Assessment	Xu, J.; Yin, S.; Zhao, G.; Wang, Z.; Peng, Y.	FineDiving
10.1109/CVPRW63382.2024.00324	2024	FineRehab: A Multi-modality and Multi-task Dataset for Rehabilitation Analysis	Li, J.; Xue, J.; Cao, R.; Du, X.; Mo, S.; Ran, K.; Zhang, Z.	FineRehab
10.1109/CBMS61543.2024.00016	2024	GYMetricPose: A light-weight angle-based graph adaptation for action quality assessment	Gallardo, U.; Caro, F.; Hernandez, E.; Espinosa, R.; Ochoa-Ruiz, G.	Fitness-AQA
10.1109/TMM.2023.3328180	2024	Learning Semantics-Guided Representations for Scoring Figure Skating	Z. Du; D. He; X. Wang; Q. Wang	OlympicFS
10.1109/TIM.2024.3398072	2024	Learning Sparse Temporal Video Mapping for Action Quality Assessment in Floor Gymnastics	S. Zahan; G. Mubashar Hassan; A. Mian	AGF-Olympics
10.1016/j.combiomed.2024.108382	2024	LightPRA: A Lightweight Temporal Convolutional Network for Automatic Physical Rehabilitation Exercise Assessment	Sardari, S.; Sharifzadeh, S.; Daneshkhah, A.; Loke, S.W.; Palade, V.; Duncan, M.J.; Nakisa, B.	UI-PRMD; KIMORE
10.1109/ACCESS.2024.3423462	2024	MMW-AQA: Multimodal In-the-Wild Dataset for Action Quality Assessment	T. Nagai; S. Takeda; S. Suzuki; H. Seshimo	MMW-AQA
10.1109/ICASSP48485.2024.10447069	2024	Multi-Stage Contrastive Regression for Action Quality Assessment	Q. An; M. Qi; H. Ma	FineDiving
10.1109/TIP.2024.3362135	2024	Multimodal Action Quality Assessment	L. -A. Zeng; W. -S. Zheng	FIS-V; FineDiving; Rhythmic Gymnastics
10.1109/CVPR52733.2024.01744	2024	Narrative Action Evaluation with Prompt-Guided Multimodal Interaction	Zhang, S.; Bai, S.; Chen, G.; Chen, L.; Lu, J.; Wang, J.; Tang, Y.	MTL-AQA; FineGym
10.1109/WACV57701.2024.00012	2024	PECoP: Parameter Efficient Continual Pretraining for Action Quality Assessment	A. Dadashzadeh; S. Duan; A. Whone; M. Mirmehdi	JIGSAWS; MTL-AQA; FineDiving; PD4T
10.1007/978-3-031-61678-5_6	2024	Sport Action Evaluation Based on Human Pose Estimation: A Case of Evaluating Golf Swing Action	Li, X.; Tao, R.; Wang, Y.; Ding, Y.; Luo, X.	Golf Swing Action
10.1016/j.ins.2024.120347	2024	Two-path target-aware contrastive regression for action quality assessment	Ke, X.; Xu, H.; Lin, X.; Guo, W.	MTL-AQA; FineDiving; AQA-7; JIGSAWS
10.1109/ICASSP48485.2024.10448158	2024	Which is the Better Teacher Action? A New Ranking Model and Dataset	M. Fang; X. Du; Q. Liu; Y. Zhou; Q. Liang; S. Liu	TAQR

Table A.2: Summary of datasets extrapolated from the reference study that were not included in the framework.

Dataset	Relation	Avail.	Article DOI	Original DOI	Incl.	Cit.	Orig.
CPR-Coach	AQA-ready	Public	10.1109/LSP.2023.3346207	10.1109/CVPR52733.2024.01777	No	8	7

Continued on next page

Table A.2 – Continued from previous page

Dataset	Relation	Avail.	Article DOI	Original DOI	Incl.	Cit.	Orig.
FS1000	AQA-ready	Public	10.1109/TMM. 2023.3328180	10.1609/aaai.v37i3. 25392	Yes	25	30
Perturbed Workout	AQA-ready	Private	10.1145/2505323. 2505330	–	No	55	–
SkatingVerse	AQA-ready	Public	10.1049/cvi2.12287	–	No	4	–
Stroke Rehabilitation	AQA-ready	Private	10.1109/CVPRW. 2013.82	–	No	42	–
TGC20ReId	AQA-ready	Private	10.1007/ 978-3-031-06433-3- 21	10.1016/j.patrec. 2020.08.003	No	8	27
SU-41 Gesture	AQA-related	Private	10.1145/2505323. 2505330	–	No	55	–
Average Total						28.14 197	21.33 64

Note: *Relation*: AQA-ready, AQA-related (adaptable), Not applicable. *Incl.*: Article included in the framework. *Cit./Orig.*: Citation counts from Google Scholar (Dec 20, 2025).

Table A.3: Summary of datasets extrapolated with the proposed framework that were not included in the reference study.

Dataset	Relation	Avail.	Article DOI	Original DOI	Incl.	Cit.	Orig.
3D Animation	AQA-ready	Private	10.1109/TVCG. 2023.3247092	–	Yes	25	–
ActivityNet1.3	AQA- related	Public	10.1109/ EIECC60864. 2023.10456639	10.1109/CVPR. 2015.7298698	No	0	3,536
AGF-Olympics	AQA-ready	Public	10.1109/TIM. 2024.3398072	–	No	15	–
CAD-60	AQA- related	Public	10.1155/2022/ 9402195	10.1109/ICRA. 2012.6224591	No	3	754
Countix	AQA- related	Public	10.1109/ACCESS. 2023.3316009	10.48550/arXiv. 2006.15418	Yes	31	183
Diving48	AQA- related	Public	10.1007/ s11263-021-01486-4	10.1007/ 978-3-030-01231-1- 32	No	73	224
ExPose	AQA-ready	Public	10.1109/ WACV48630.2021. 00044	10.1109/ CVPRW50498. 2020.00458	Yes	39	41
FINA09	AQA-ready	Private	10.1109/CVPRW. 2017.16	10.1007/ 978-3-642-22822-3- 27	No	283	17
FineRehab	AQA-ready	Public	10.1109/ CVPRW63382. 2024.00324	–	No	18	–
Fitness-AQA	AQA-ready	Private	10.1109/ACCESS. 2023.3316009	10.48550/arXiv. 2202.14019	Yes	31	42
Golf Swing Action	AQA- related	Private	10.1007/ 978-3-031-61678-5- 6	–	No	0	–
Golf Swing Video	AQA- related	Private	10.1145/3606038. 3616150	–	Yes	7	–

Continued on next page

Table A.3 – Continued from previous page

Dataset	Relation	Avail.	Article DOI	Original DOI	Incl.	Cit.	Orig.
Infinite Grasp	AQA-ready	Private	10.1109/ ICMLC58545. 2023.10327994	10.48550/arXiv. 1901.02579	Yes	0	76
Kinetics 700	AQA- related	Public	10.1109/ACCESS. 2023.3316009	10.48550/arXiv. 1907.06987	Yes	31	681
KTH	AQA- related	Public	10.3390/ electronics11193051	10.48550/arXiv. 1907.06987	No	0	5,397
MM-Fit	AQA- related	Public	10.1109/ACCESS. 2023.3316009	10.1145/3432701	Yes	31	91
MSR-Action3D	AQA-ready	Public	10.1016/j.patcog. 2017.12.007	10.1109/CVPRW. 2010.5543273	No	128	1,936
MSRC-12	AQA- related	Public	10.1016/j.patcog. 2017.12.007	10.1109/CVPR. 2011.5995316	No	128	4,935
Northwestern- UCLA	AQA- related	Public	10.1155/2022/ 9402195	10.1109/CVPR. 2011.5995729	No	3	360
Ranked TikTok (CDRG- UNIMAS)	AQA-ready	Private	10.26555/ijain. v9i1.919	–	No	7	–
RDSB	AQA- related	Private	10.1109/ ISCAS48785.2022. 9937262	–	No	5	–
SMART	AQA-ready	Public	10.1007/ s11263-021-01486-4	–	No	73	–
TAQR	AQA-ready	Public	10.1109/ ICASSP48485. 2024.10448158	–	No	5	–
Thumos14	AQA- related	Public	10.1109/ EIECC60864. 2023.10456639	10.1016/j.cviu. 2016.10.018	No	0	747
UCF-101	AQA- related	Public	10.1109/ACCESS. 2023.3316009	10.48550/arXiv. 1212.0402	Yes	31	8,393
Average						38.68	1,713.31
Total						967	27,413

Note: Relation: AQA-ready, AQA-related (adaptable), Not applicable. Incl.: Article included in the reference study. Cit./Orig.: Citation counts from Google Scholar (Dec 20, 2025).