# A Literature Review of Literature Reviews in Pattern Analysis and Machine Intelligence

Penghai Zhao, Xin Zhang, Ming-Ming Cheng, Jian Yang, Xiang Li*

**Abstract**—By consolidating scattered knowledge, the literature review provides a comprehensive understanding of the investigated topic. However, excessive reviews, especially in the booming field of pattern analysis and machine intelligence (PAMI), raise concerns for both researchers and reviewers. In response to these concerns, this Analysis aims to provide a thorough review of reviews in the PAMI field from diverse perspectives. First, large language model-empowered bibliometric indicators are proposed to evaluate literature reviews automatically. To facilitate this, a meta-data database dubbed RiPAMI, and a topic dataset are constructed, which are utilized to obtain statistical characteristics of PAMI reviews. Unlike traditional bibliometric measurements, the proposed article-level indicators provide real-time and field-normalized quantified assessments of reviews without relying on user-defined keywords. Second, based on these indicators, the study presents comparative analyses of different reviews, unveiling the characteristics of publications across various fields, periods, and journals. The newly emerging AI-generated literature reviews are also appraised, and the observed differences suggest that most AI-generated reviews still lag behind human-authored reviews in several aspects. Third, we briefly provide a subjective evaluation of representative PAMI reviews and introduce a paper structure-based typology of literature reviews. This typology may improve the clarity and effectiveness for scholars in reading and writing reviews, while also serving as a guide for AI systems in generating well-organized reviews. Finally, this Analysis offers insights into the current challenges of literature reviews and envisions future directions for their development.

**Index Terms**—literature review, pattern analysis, machine intelligence, bibliometric, AI-generated scholar content.

◆

## 1 INTRODUCTION

*"The entropy of the universe tends to a maximum."*

— Rudolf Clausius

THE entropy of almost all natural and artificial systems in the universe exhibits a continuous increase. This principle also applies to the knowledge system of humanity, which similarly undergoes a perpetual escalation in entropy. Such escalating disorder within the knowledge system can lead to redundant endeavors, adversely affecting the generation of new knowledge. Analogous to the role of gravity in shaping the early universe filled with diverse particles, literature review plays a vital role in the knowledge system by consolidating scattered knowledge. Essentially, a literature review is a scholarly composition that not only demonstrates an understanding of the academic literature on a given topic but also situates this knowledge within a broader context. It compiles the most relevant and significant publications related to a specific research area, thereby offering a comprehensive overview of that field.

Almost all the field boast their own literature reviews, especially for the rapidly developing field of pattern analysis and machine intelligence, including image classification [20], [61], [62], [64], image segmentation [35], [67], [70], [91], object detection [21], [47], [81], [116], natural language processing [36], [63], [69], etc. As reported by the AI Index Report [65], there has been a striking surge in artificial intelligence publications, soaring

from 200,000 in 2010 to nearly 500,000 by 2021. This exponential increase has subsequently led to a proliferation of related literature reviews. This trend is clearly illustrated in Fig. 1, which shows a marked increase in the annual publication of reviews, underscoring the growing prevalence of literature reviews in the field.

While comprehensive literature reviews are valuable, excessive reviews may lead to information overload and redundant effort. These issues emerge when multiple reviews address the same topic, creating significant redundancy and presenting readers with the same studies and ideas repeatedly. Additionally, writing and peer-reviewing literature reviews imposes a considerable burden on both authors and reviewers, as it involves a broader and more thorough examination of the literature than standard research papers typically demand. Additionally, with the advent of AI-generated literature review systems, the landscape is further complicated. The differences between AI-generated and human-authored reviews remain an area for exploration.

Up to now, far too little attention has been paid to the above-mentioned issues. This paper concentrates on reviews within the field of PAMI, with three distinct objectives: (1) to endeavor an automated evaluation of these literature reviews in a relatively scientific manner; (2) to offer a thorough review of existing human-authored and AI-generated literature reviews from diverse perspectives; (3) to pinpoint limitations while providing future insights for literature review.

### 1.1 Scope

In this study, we exclusively examine review articles. Grant's comprehensive review [31] concluded that there are at least fourteen review types characterized by methods used, which highlights the fact that different reviews serve different purposes. For instance, a critical review strives to exhibit that the author has carried

* *Corresponding author. https://github.com/IMPlus-PCALab*

- *P. Zhao, X. Zhang, M. Cheng, J. Yang and X. Li are with VCIP, the College of Computer Science, Nankai University. Email: xiang.li.implus@nankai. edu.cn, zhaopenghai@mail.nankai.edu.cn*
- *This work was supported by the Young Scientists Fund of the National Natural Science Foundation of China (Grant No.62206134) and the Tianjin Key Laboratory of Visual Computing and Intelligent Perception (VCIP). Computation is Supported by the Supercomputing Center of Nankai University (NKSC).*
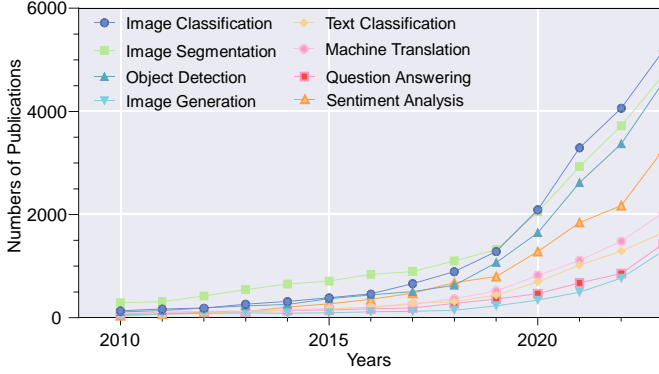
Fig. 1. Annual publication trends of literature review in the field of PAMI. A notable rising trend can be observed since 2015, reflecting an increasing scholarly focus and growing recognition of the importance of review articles in synthesizing the state of research in PAMI. The data is collected from the Google Scholar search engine.

out extensive research on literature and evaluated it critically. A systematic review aims to systematically explore, assess and integrate research evidence, usually in adherence with review guidelines. However, these types of reviews are relatively rare in the field of PAMI. Given the predominance of literature reviews in this area, we entitle this paper "A Literature Review of Literature Reviews."

Although the quantity of reviews is considerably smaller compared to that of normal papers, it remains impractical to analyze reviews within every field. Therefore, this paper will focus only on reviews in the PAMI field. The reviews being analyzed in this paper are sourced from various publication venues, including conferences, journals, and pre-prints. Reviews from various sources exhibit a diverse range of writing styles, article lengths, and citation counts, providing a representative cross-section of the PAMI field. Considering the recent emergence of deep learning in PAMI, most reviews being investigated were published in the last decade. However, for a comprehensive understanding, we also analyze a few reviews published over ten years ago.

## 1.2 Contribution

In summary, the main contributions of this survey are as follows:

- We conduct a comprehensive and detailed review of literature reviews in the field of PAMI. To the best of our knowledge, there have been no scholarly attempts to systematically review the literature reviews.
- A biography database dubbed RiPAMI is built and released. This database contains meta-data for 2904 literature reviews, including titles, authors, citations and references detail, etc. In addition, we construct a topic key phrase dataset comprising 201 different domain reviews through manual annotation. This dataset can be utilized to evaluate the prompt effectiveness in identifying the paper topic or to conduct further analyses for the community.
- Subjective evaluations of seminal literature reviews in PAMI are presented. Alongside this, we introduce an organizational structure-based review typology designed to guide both human authors and AI systems in the methodical crafting of literature reviews.
- We propose the impact and quality indicators for quantitatively evaluating literature reviews. Diverging from tradi-

tional, our approaches offer real-time, article-level, and field-normalized quantitative assessment of literature reviews. These automated indicators may help address concerns arising from excessive reviews and provide support during the appraisal stage of AI-generated review systems.
- Based on the proposed evaluation indicators, we thoroughly assess the quality of the selected human-authored literature reviews in the PAMI field. Furthermore, a comparison between human-authored and AI-generated reviews is carried out, which highlights the enduring strengths of human-authored literature reviews and uncovers the limitations of existing AI-generated review systems.
- We briefly discuss the challenges and provide a preliminary analysis of future directions for literature reviews in the PAMI field, highlighting potential areas for improvement.

All the data and code framework used in this paper are publicly available at https://sway.cloud.microsoft/2TXEuPuNlDKEmC9p.

## 1.3 Organization of the Paper

The remainder of the paper is organized as follows. Section 2 briefs the evolution of the literature review and details the database constructed for this paper. Section 3 introduces the typology of reviews and provides subjective comments on selected reviews. Section 4 presents the quantitative measures used to gauge the impact and quality of the publications respectively, and compares literature reviews from various journals and fields based on these evaluation indicators. The characteristics of human-authored and AI-generated literature reviews are discussed in Section 5 and a case study is carried out to elucidate what are the advantages and shortcomings of these two types of literature reviews. In Section 6, We further discuss the challenges and future of the literature review. Finally, this paper is concluded in section 7.

## 2 BACKGROUND AND SETTING

### 2.1 Background of Literature Review

Given the historical evolution and disciplinary variations of literature reviews, it is challenging to discern the very first peer-reviewed literature review in history. However, there is limited research to suggest that contemporary literature reviews can be traced back to the 17th and 18th centuries [37] when the scientific method and the concept of peer review began to take shape. Initially, literature reviews were focused more on assimilating existing knowledge as opposed to following a formalized structure. Scholars used prior works to direct their research and grounded their contributions in an established understanding of the subject. Literature reviews during this period served to prevent duplicated efforts and build upon the work of predecessors.

As scholarly communication advanced and the scientific method became more rigorous, the literature review experienced a notable transformation in both methodology and purpose. Nowadays, a well-executed literature review is more than a basic summary of references; it goes beyond that by organizing and combining the information through analysis and synthesis. Reading a good literature review may benefit both novice and seasoned researchers. For novice researchers, literature review is one of the most important ways to gain knowledge in a specific field. They can quickly learn the basics, grasp fundamental concepts, and gain a comprehensive understanding of major theories. For seasoned researchers, literature reviews may help them keep up

with the latest research, identify gaps in the current research field, and avoid duplicating work that has already been done.

## 2.2 Database

Literature reviews usually cite regular papers in a particular field as references; however, our literature review investigates literature reviews in various fields as the material to be analyzed and synthesized. This section describes how the reviews are selected and provides further details on the construction of the database.

### 2.2.1 Data Source

Reliable data sources for analyzing extensive reviews are fundamentally important. Based on the means of data acquisition and storage, existing scientific scholar data sources may be classified into two main categories: web-based and snapshot-based sources. Web-based source data refers to the meta-data that can be retrieved from the data provider in real-time with the use of the web crawler or the API (e.g. Semantic Scholar, arXiv, CrossRef). Such an approach would only consume a small amount of storage on the local machine but needs to query meta-data every single time. On the contrary, the offline snapshot consumes a larger amount of storage space but eliminates the need for frequent API queries. Since the snapshot is a mirror image of relevant papers before a specific time, its data remains unchanged over time compared to the API-based sources. As a result, snapshot ensures a consistent dataset in all of the experiments, avoiding issues of irreproducibility caused by changes in provided web retrieving service.

Tab. 1 compares various most commonly used data sources, where "Counts only" in the citations column means that the source only records the citation counts, while "Complete" signifies that the data source provides a complete list of citations. Unfortunately, none of these sources is perfect. For example, arXiv is a free distribution service and an open-access archive for millions of scholarly articles in various fields. Users can utilize arXiv's APIs to access all of the paper meta-data stored in arXiv databases. Meta-data retrieved from arXiv contains valuable information but fails to query the publication venue, citations, and references. Semantic Scholar seems promising, but it suffers a lower update rate than arXiv and a narrower search scope than Google Scholar. Google Scholar is an online search engine that indexes scholarly literature from a wide range of disciplines and publication formats. It employs automated programs to retrieve files for inclusion in the search results (which are not limited to academic papers, but also include patents, books, etc.). Despite its widespread use, Google Scholar still encounters challenges. Beel [11], [12], [13] argues that Google Scholar places a high weight on citation counts in its ranking algorithm and has therefore been criticized for exacerbating the Matthew effect. Moreover, the citation counts displayed on Google Scholar are subject to manipulation by complete nonsense articles indexed on Google Scholar (e.g. citations from AI-generated pre-print papers published on arXiv should have been ignored). Therefore, a promising engineering solution is to leverage the strengths of the different approaches to overcome the weaknesses of each approach, as demonstrated in this paper.

### 2.2.2 Database construction

We first detail the database construction approach used for this paper (as illustrated in Fig. 2). To avoid potential copyright and licensing issues, papers are retrieved and downloaded using the arXiv's API. Calls to the API are made by means of HTTP requests to a certain URL. The responses will be cached and stored in an SQL-based database.

We employed a simple yet effective technique to ensure that the type of paper retrieved is review paper and highly relevant to the field of pattern recognition. First, we identify 106 keywords based on the scopes of related journals and conferences. As can be seen in Fig. 2, these selected keywords include, but are not limited to, speech recognition, optical character recognition, and self-supervised learning. Next, we conduct a primary retrieving and filter the results following a clear rule, i.e. the title of the paper must contain "survey" or "review", and the keyword should be included in the abstract. Additionally, we filter out noisy data through both ChatGPT-based and manual double-checking processes.

As mentioned in Sec. 2.2.1, we suggest enriching the meta-data of papers by leveraging a combination of disparate data sources. Considering the potential legal risks of crawling to obtain academic data from Google Scholar, Semantic Scholar API was employed to obtain additional meta-data, such as citation and reference details which are not provided by the arXiv API.

Finally, a total of 2904 eligible papers with meta-data are retrieved and collected. To ensure reproducible experiments and prevent overburdening the server, we construct an SQL-based database dubbed RiPAMI (Reviews in Pattern Analysis and Machine Intelligence). This database stores information related to the paper such as title, abstract, date of publication, venue, citaion and reference details, etc.

It should be noted that this paper employs the arXiv API for retrieving literature reviews. As such, there might be potential biases and issues related to incomplete retrieval. That is, papers that are not published on arXiv but meet the criteria will not be included in the database. However, we believe that such a problem is unavoidable. On the one hand, it is difficult for any researcher to guarantee that he or she is able to retrieve the entire relevant literature. On the other hand, most researchers tend to select high-level articles as references to ensure the quality of the literature review, which also leads to the bias issue to some extent. Considering that most articles published on arXiv will also be published subsequently in different conferences and journals, the arXiv API-based retrieval can be regarded as a sample of the full set of literature reviews. We expect that the dataset constructed in this manner shares similar statistical characteristics with the full set of literature reviews.

Given that citation counts vary over time, the date for retrieving citation details is January 1, 2024.

### 2.2.3 Database Statistics

As aforementioned, The database consists of more than 2900 literature reviews from a variety of sources, publication years, and fields. To elucidate the characteristics of the RiPAMI database, we conduct a statistical analysis and plot the result in Fig. 3.

**Years of Publication** Figure 3 (a) illustrates the distribution of publication years of literature reviews contained in the RiPAMI database. What can be clearly seen in this figure is a growing trend in the number of reviews. This trend has similar characteristics to the one obtained in Fig. 1, e.g. both are steadily increasing and showing a fairly significant increase between 2019 and 2020. This to some extent reflects the consistency between our sample data and the original data in terms of statistics.

| Database | Title & Authors | Venue | Abstract | Citations | References | Source Types | Charge |
|---|---|---|---|---|---|---|---|
| arXiv | ✓ | ✗ | ✓ | ✗ | ✗ | API-based | ✗ |
| CrossRef | ✓ | ✓ | ✓ | Counts Only | ✓ | API-based | ✗ |
| Google Scholar | ✓ | ✓ | ✓ | Complete | ✗ | Crawler-based | ✗ |
| IEEE Xplore | ✓ | ✓ | ✓ | Counts Only | ✗ | API-based | ✗ |
| Semantic Scholar | ✓ | ✓ | ✓ | Complete | ✓ | API-based | ✗ |
| Web of Science | ✓ | ✓ | ✓ | Complete | ✓ | API-based | ✓ |
| Scopus | ✓ | ✓ | ✓ | Complete | ✓ | API-based | ✓ |
| arXiv Data File | ✓ | ✗ | ✓ | ✗ | ✗ | Snapshots | ✗ |
| CrossRef Data File | ✓ | ✓ | ✓ | Counts Only | ✓ | Snapshots | ✗ |
| **RiPAMI(Ours)** | ✓ | ✓ | ✓ | Complete | ✓ | Snapshots | ✗ |

TABLE 1
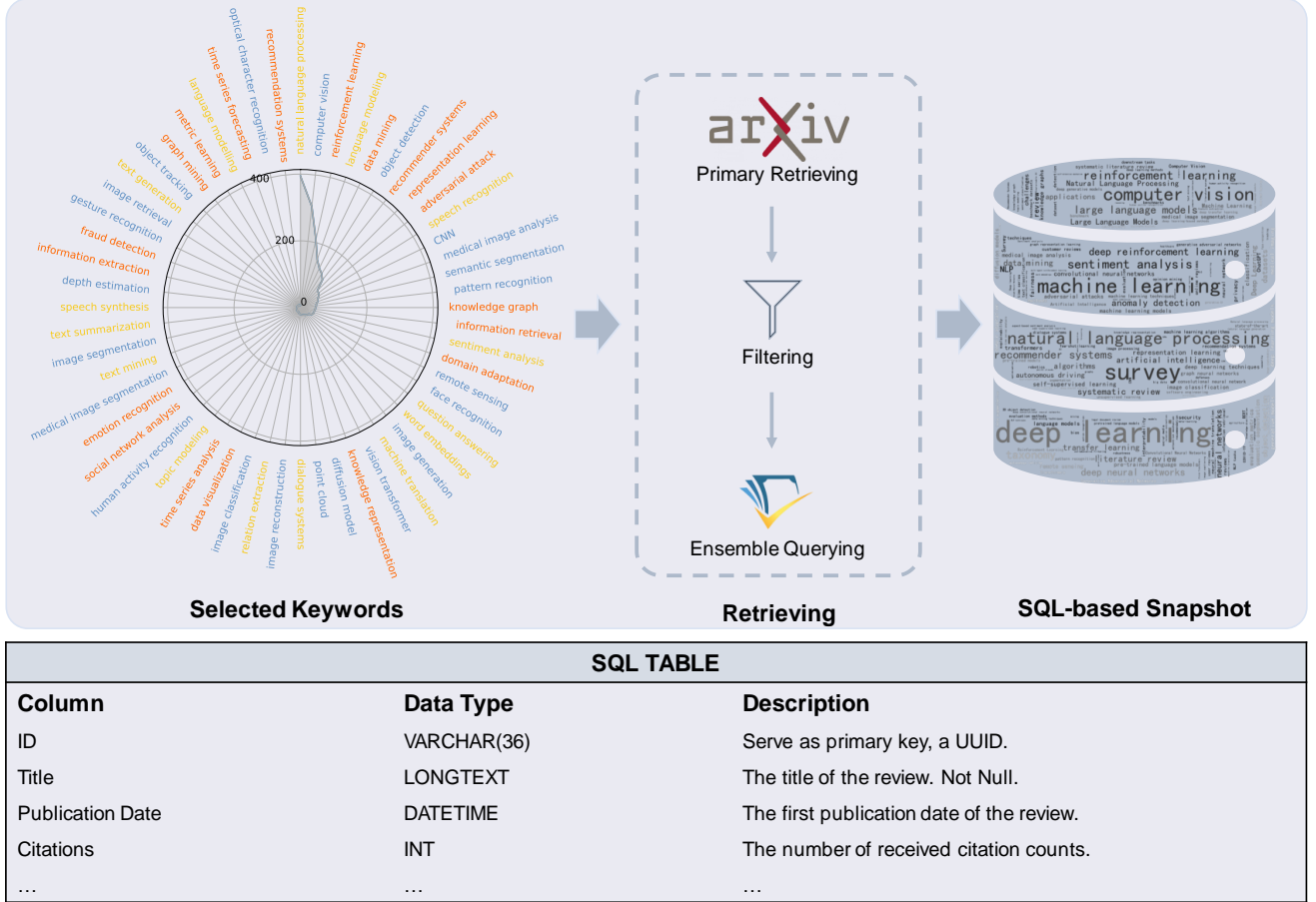Comparisons between various data source.



Fig. 2. Schematic representation of the database construction process. From keywords to a SQL-based snapshot RiPAMI, there are three major steps to ensure that the data in RiPAMI is clean, accurate, and reliable. For visual clarity, the polar figure displays only a part of retrieving keywords.

**Number of the References** The number of references in a survey paper may influence its credibility and reliability. As shown in 3 (b), the distribution of literature review references in the RiPAMI database follows a log normal pattern. The average number of cited references is approximately 140, while the median number is 123.

**Number of the Authors** A review paper typically aims to provide a comprehensive overview of a specific topic. Involving multiple authors with diverse expertise could enhance the depth and breadth of the review. Fig. 3 (c) indicates that the majority of reviews are written by fewer than 10 authors.

**Number of the Citations** We count the citations of papers in

RiPAMI and plot them in Fig. 3 (d). A power law distribution of received citations could be found. This phenomenon where a small subset of papers receives the majority of citations is sometimes referred to as the "Matthew Effect" or the "Pareto principle", as reported in [17], [68].

## 3 A SUBJECTIVE EVALUATION OF LITERATURE REVIEWS

### 3.1 Popular Paper Structure

The review's structure is alternatively referred to as the framework of the review. It is the outline that authors consider when starting
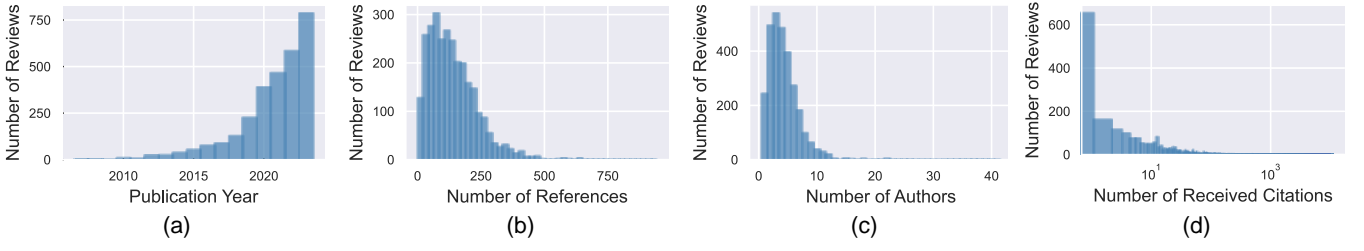
Fig. 3. Statistics of the RiPAMI database. The samples in the RiPAMI database are characterized by a diversity of publication dates, scholar impact, reference numbers, etc., offering a comprehensive reflection of the state of reviews in the PAMI field.

to conduct the survey. A well-designed structure is believed to enhance the paper's readability and facilitate the reader's comprehension of the paper's concepts and knowledge. Typically, this framework encompasses foundational sections such as the introduction, methodology, discussion, and conclusion, each fulfilling a specific function.

Similar to the research article, the introduction section of the literature review usually includes contextualizing the research topic, identifying the knowledge gap, defining the scope and objectives, and laying the foundation knowledge for wider audiences. "Introduction" usually lies at the very beginning of the paper. Most reviews first introduce the definition of the topic to acquaint readers with a basic understanding of the field. For instance, the term "Named Entity Recognition" may be unfamiliar to many scholars outside the field of natural language processing. A clever introduction might provide the origins of the terminology and inform the readers what is named entity recognition, as has been done in Ref [73]. A brief definition of the field in the introduction is also welcomed for some relatively popular fields, e.g., the literature review [84] examines the concept of image categorization in the first sentence.

In addition to contextualizing the research topic, many introductions also highlight the existing research and identify gaps or challenges in the current understanding of the topic. It sets the stage for the literature review by explaining why the research being reviewed is necessary and what gaps are expected to be filled. Wang et al. [103] emphasize that while there are comprehensive surveys of self-supervised learning in computer vision, there is a lack of a similar overview specifically tailored to the remote sensing community.

Readers need to evaluate if the article is worth reading before investing more time in it. A statement of the scope or contributions of the article is a way for readers to quickly assess its relevance. In the first section, Wangkhade et al. [104] provide us with several important contributions including analyzing well-known technologies, proposing taxonomies of approaches, and summarizing benefits and challenges of sentiment analysis.

Preliminaries and problem formulations are also popular with readers, as they both provide profound background knowledge. Given that the Gumbel-max involves considerable mathematical concepts and calculations, a "Preliminaries" section in the review of the Gumbel-max trick [39] serves the purpose of providing background information and basic understanding related to the Gumbel-max trick.

The middle part of the survey paper, also known as the review part, presents a detailed examination of the relevant research studies, methodologies, findings, and theories related to the chosen theme of the review. Beyond that, The middle part synthesizes the information from related studies to provide a cohesive and integrated understanding of the topic. This synthesis not only aligns the disparate studies but also evaluates their contributions and the interrelations among them. It often includes comparative analyses, highlighting similarities and differences in approaches, and may identify limitations in the existing research that warrant further investigation. This section is instrumental in demonstrating how individual studies collectively advance understanding of the subject matter.

The manner of organization and synthesis in the review part is subject to the choice of topic. Based on the organization of the review part, we propose a typology for literature reviews which includes three major types: method clustering-based, challenge-oriented, and hybrid literature review. More details about the typology are presented in Section 3.2.

The subsequent ending part serves as a succinct conclusion to the reviewed studies. Within this section, the authors prefer to sum up the key findings to answer the question in the beginning part. It usually highlights both the advantages and disadvantages of reviewed literature to conclude research gaps and suggest potential avenues for future research.

One of the main objectives of the concluding part is to provide a relatively concise summary of the main insights derived from the reviewed studies and highlight their significance and relevance to the research topic. Some literature reviews conduct comprehensive analysis and synthesis, which in turn run up to a considerable number of pages. The paper [51] consists of 51 double-column pages in the main body, which poses a challenge for readers who may not have the time to peruse it thoroughly. Fortunately, this paper also includes a concise summary which provides readers with an efficient understanding of the content and allows them to navigate to relevant sections of interest.

When discussing further in the ending part, the authors may attempt to identify gaps and point out future directions through the analysis of existing literature. Some of the papers present the gaps and future directions separately. Hassain et al. [40] first discussed major challenges of multi-view video summarization such as lack of synchronization, instability of camera, and crowded scenes individually. Where there are challenges, there are future research directions. In addition to challenges, the authors also provide recommendations and future directions from various perspectives including models, benchmark datasets, and agents-based MVS, etc. Conversely, some papers choose to combine discussions of current issues with emerging research trends, as seen in the 'Future Directions' section of Ref [110].

For literature reviews delving into pragmatic methodologies,

the inclusion of an applications section is quite fitting. Ref [48] offers a comprehensive review of various surveys on convolutional neural networks, dedicating a distinct section to the typical applications of 1-D, 2-D, and multidimensional CNNs. Similarly, Ref [59] illustrates the deployment of few-shot learning techniques across disciplines like computer vision, natural language processing, and reinforcement learning.

## 3.2 Typology of the Literature Review

To investigate more comprehensively how papers are organized and synthesized, we propose a typology for the literature review based on the organization of the middle part. Note that the terms typology and taxonomy are often used interchangeably, but there is a subtle difference. According to Ref. [93], typology creates categories based on conceptual dimensions and idealized types, providing a systematic basis for comparison. Conversely, taxonomy categorizes items based on observable and measurable characteristics. Hence, we opt for "typology" to categorize literature reviews in this context, considering their content and structure.

The typologies for review have been well developed [23], [24], [31], [78]. Most of these articles mainly focus on categorizing literature reviews based on their purpose or the analytical framework used, such as the SALSA framework (Search, Appraisal, Synthesis, and Analysis). For example, literature reviews can be categorized as narrative reviews, critical reviews, etc., depending on their purpose. However, few studies have investigated how the review part of a literature review is organized in a systematic way. Investigating the organization of the review part may enhance the understanding of the literature review, thereby offering valuable guidance for researchers and AI-generated review systems. Based on the writing style and paper structure, we categorized existing literature reviews into three main types: method clustering-based, challenge-oriented, and hybrid literature review.

### 3.2.1 Method Clustering-based Literature Review

Method Clustering-based literature review refers to grouping methods according to their technical characteristics and presenting them in separate sections or subsections. Authors need to identify different sections depending on the topic of the article and arrange the same type of methods as closely as possible.

Featuring a clear and well-organized article structure, method clustering-based literature reviews are widely favored by researchers. The reader can gain a comprehensive understanding of the typology, details, and advancements of technology within a specific field which facilitates a foundational comprehension of the research field. However, such a type of literature review is less tailored for readers who are completely unfamiliar with the subject. This is because readers who lack basic knowledge of the field may not understand the connection between the methods mentioned in different sections, or are even confused to the select a proper algorithm for a certain task.

By way of illustration, Minaee *et al.* conducted a comprehensive survey on deep learning-based image segmentation models [70] in a method clustering-based manner. They grouped the models into 10 categories based on the adopted model architectures, such as fully convolutional models, encoder-decoder-based models, and attention-based models. In this way, the reader can quickly gain an understanding of the typologies of deep learning-based image segmentation models and how they differ from each other.

### 3.2.2 Challenge-oriented Literature Review

Known for its practicality, challenge-oriented literature reviews focus on presenting research methods and findings that address a specific challenge or task. In this type of review, each section or subsection focuses on a specific challenge or task and centers the discussion on how to address that challenge. The references cited within each section are primarily directly related to addressing that challenge to provide support and evidence.

Challenge-oriented literature reviews are well-suited to application-oriented fields, where readers at different levels of expertise can easily and quickly seek out appropriate solutions within the literature to address the problem they are facing. Despite the advantages, challenge-oriented literature reviews may fail to provide a comprehensive elaboration of fundamentals and methods details given the limited paper length.

Bandini *et al.* divided their survey about egocentric vision hand analysis into several sections including hand segmentation, hand detection, and hand identification [10]. Each section contains multiple subsections, and each of the subsections investigates a certain challenge. (e.g. subsections entitled "Robustness to Illumination Changes" and "Lack of Pixel-Level Annotations" were arranged in the "Hand Segmentation" section ).

### 3.2.3 Hybrid Literature Review

By integrating characteristics of method clustering-based and challenge-oriented review, a hybrid literature review strikes a fair balance between comprehensiveness and practicality.

Typically, the hybrid literature review refers to a challenge-oriented framework incorporating a method clustering-based sub-framework. One advantage of a hybrid review is that it simultaneously offers multiple solutions to a certain problem along with the classification of the respective methodologies. Such a hybrid appears to be a perfect framework for review parts, but unfortunately, not all topics are suitable for the hybrids, particularly regarding those more theoretical fields. Therefore, it is advisable to adopt different frameworks for different topics.

An example of this is the study carried out by Guo *et al.* [33] in which authors divide sections by tasks and introduce similar approaches group by group.

## 3.3 Review of Selected Reviews

In this section, we present a subjective evaluation of selected reviews in different areas. In contrast to the previous surveys, the literature reviews investigated in this section tend to be representative. This means that these surveys are of fairly high quality and popular with a wide range of researchers, evident in the number of citations they received.

### 3.3.1 Computer Vision

Computer vision is one of the most popular sub-fields of pattern recognition and machine intelligence. This section will focus on the literature reviews within the realm of computer vision.

**Image Classification** refers to the task of assigning a label or a category to an input image, which is one of the most renowned tasks in the field of computer vision. A survey by Rawat [84] explores the development and advancements of deep convolutional neural networks (CNNs) in the field of image classification. The paper covers the historical context, their role in the deep learning renaissance, and the notable contributions and challenges faced in recent years. It highlights the remarkable progress of CNNs

in image classification, while also acknowledging the ongoing research efforts to address challenges and provide recommendations for future exploration. Schmarje *et al.* [85] provides a comprehensive survey on semi-, self-, and unsupervised learning methods for image classification. The survey compares and analyzes 34 different methods based on their performance and commonly used ideas, highlighting the trends and research opportunities in the field. Through comprehensive analysis, the authors reveal the potential of semi-supervised methods for real-world applications and identify challenges such as class imbalance and noisy labels. Furthermore, the paper emphasizes the importance of combining different techniques from various training strategies to improve overall performance. In addition to CNNs, there exist alternative techniques for image classification. The paper [19] presents a comprehensive analysis of Support Vector Machines (SVM) in image classification. It discusses various techniques that can enhance classification accuracy and highlights its advancements. Liu *et al.* [57] investigate more than 100 different visual Transformers comprehensively in three fundamental CV tasks including classification, detection, and segmentation. They also propose a taxonomy to categorize various transformers into six groups.

**Object Detection** entails identifying and localizing objects of interest within an image or video. Liu *et al.* [51] offers a comprehensive survey on the advancements in deep learning-based generic object detection. This paper discusses an extensive range of issues, including detection frameworks, taxonomies, feature depiction, training strategies, and evaluation metrics. Though there have been significant advancements in generic object detection, the detection of small objects, which focuses on identifying objects with a small size, still presents challenges. The review conducted by Cheng *et al.* [21] investigates 181 literature, constructs two large-scale datasets (SODA-D and SODA-A), and evaluates the performance of mainstream small object detection methods. Object detection demonstrates the utility and effectiveness across multiple domains. Li *et al.* [47] and Litjens *et al.* [49] investigate numerous methods and applications of object detection in remote sensing and medical image analysis respectively, showing that these methods have the flexibility to be applied in various scenarios and meet different needs.

**Image Segmentation** is the process of dividing an image into meaningful and distinct regions to facilitate analysis and understanding. This work referenced 196 papers and received 1900 citations. As described earlier, Minaee *et al.* [70] proposed a taxonomy for image segmentation methods which divides models into 11 categories. In addition to the taxonomy, the authors evaluate the quantitative performances of various methods on popular benchmarks. The paper also identifies open challenges and proposes promising research directions for future advancements in deep-learning-based image segmentation. Given that most image segmentation algorithms heavily rely on expensive pixel-level annotations, interest in weakly supervised image segmentation methods has increased. Ref [89] surveys label-efficient deep image segmentation methods. According to the paper, weakly supervised segmentation approaches can be categorized into four hierarchical types, ranging from no supervision to inaccurate supervision. The authors investigated each of these four methods in separate sections, highlighting the strategies used to bridge the gap between weak supervision and dense prediction. Image segmentation techniques have a wide range of applications in the field of medical image processing, as introduced in [60], [83], [91], [108].

### 3.3.2 Natural Language Processing

Acclaimed as the jewel of the artificial intelligence crown, natural language processing (NLP) stands as a pivotal domain within the field of PAMI. Here, we provide a further discussion on several popular NLP research directions.

**Named Entity Recognition** (NER) involves identifying and classifying named entities in text, such as person names, organizations, locations, and dates. The survey by Li *et al.* [46] begins by introducing NER resources, including tagged NER corpora and off-the-shelf NER tools. Then, authors categorize existing works based on a taxonomy that considers distributed representations for input, context encoder, and tag decoder. The paper surveys representative methods for applying deep learning in various NER tasks, and provides a valuable reference for designing deep learning-based NER models. NER serves as the foundation technique for various natural language applications, such as relation extraction [74], knowledge graph [2], etc. Due to the linguistic variance of different languages, NER methods may also vary from language to language. Surveys about various language-specific NER could be found in [52], [82], [105].

**Sentiment Analysis** focuses on determining the sentiment or emotion expressed in text, such as positive, negative, or neutral. Yadav *et al.* introduce the process of gathering and analyzing people's opinions and sentiments from various sources such as social media platforms and blogs in their paper [109]. The paper evaluates and compares different approaches used in sentiment analysis, with a focus on supervised machine learning methods like Naive Bayes and SVM algorithms. The common application areas of sentiment analysis and the challenges involved in accurately interpreting sentiments are also reported. The paper by Yue [114] categorizes and compares a large number of techniques and methods from three different perspectives: task-oriented, granularity-oriented, and methodology-oriented. It also explores different types of data and advanced tools for research, highlighting their strengths and limitations.

**Language Modeling** involves training models to understand and generate human language. As early attempts, recurrent neural networks achieved desirable performance and wide application at that time, despite some shortcomings. The paper [113] specifically focuses on RNNs and long short-term memory (LSTM) cells. The authors highlight the limitations of traditional RNNs and emphasize the significance of LSTM in handling long-term dependencies. They discuss various LSTM cell variants and their performance on different characteristics and tasks. Furthermore, the paper also categorizes LSTM networks into two major types: LSTM-dominated networks which optimize connections between inner LSTM cells, and integrated LSTM networks which incorporate advantageous features from various components. Recently, Large Language Models (LLMs) have drawn widespread attention. LLMs demonstrate significant performance improvements and unique abilities such as in-context learning, setting them apart from smaller-scale models. By investigating more than 600 works, Zhao *et al.* [121] conduct a comprehensive review of the recent advancements in LLMs. The authors discuss the evolution of language modeling techniques, from statistical models to neural models, and highlight the emergence of pre-trained language models as a powerful approach in NLP tasks. The survey focuses on LLMs with a parameter scale exceeding 10 billion and explores four key aspects: pre-training, adaptation tuning, utilization, and capacity evaluation. The paper also presents available resources

for developing LLMs and discusses important implementation guidelines. Overall, this survey serves as an up-to-date and valuable reference for researchers and engineers interested in the field of LLMs.

### 3.3.3 Others

Reviews in other popular sub-fields will be investigated in this section.

The paper by Zhou *et al.* [123] covers the evolution of pre-trained foundation models from BERT to ChatGPT and highlights their significance as parameter initializations for downstream tasks. The survey explores popular pre-trained foundation models in text, image, and graph modalities, discussing their components, pre-training methods, and advancements thoroughly. The paper also addresses topics including model efficiency, compression, security, and privacy, while offering valuable insights into scalability, logical reasoning ability, and cross-domain learning. Another survey paper [112] presents a comprehensive review of the state-of-the-art in self-supervised recommendation (SSR). The paper proposes an exclusive definition of SSR and develops a taxonomy that categorizes existing SSR methods into four categories: contrastive, generative, predictive, and hybrid. It further introduces an open-source library called SELFRec, which incorporates a wide range of SSR models and benchmark datasets. Through rigorous experiments and empirical comparison, the paper derives significant findings related to the selection of self-supervised signals for enhancing recommendation. The conclusion highlights the limitations and outlines future research directions in the field of self-supervised recommendation.

## 4 AN OBJECTIVE EVALUATION OF LITERATURE REVIEWS

Objective evaluation of literature reviews, or any other academic work, is a crucial aspect of academic research, as it enables the quantitative assessment of a certain academic publication in a relatively fair, transparent, and reproducible manner. It provides a basis for analysis and comparison, which may help readers, authors, reviewers, and editors to ensure the reliability and credibility of the paper.

In this section, we introduce both "external" impact and "internal" quality indicators designed to assess the academic impact and the content reliability of literature reviews. Our methodologies are cost-effective, transparent, and fully automated, facilitating real-time, reproducible evaluations across various dimensions of review papers. We have appraised a wide range of reviews using the proposed measurements, and further provide comparative and visual analysis for reviews in diverse fields, venues, and periods.

### 4.1 Impact Indicators

Perhaps assessing the impact of a paper sounds very simple, but simple things are always the most difficult. Bibliometrics is a research field that uses quantitative analysis and statistics to appraise the impact of scholarly publications. It is believed to play an important role in an individual's academic career, such as grant proposals or candidates for academic positions [45]. However, most existing metric methods suffer from several limitations including unfair comparison, misuse, and manipulation as introduced in the "San Francisco Declaration on Research Assessment" (DORA) [29] and the Leiden Manifesto [38].

| Metric | Assessing Level | Normalized | Pre-defined Keywords Free |
|---|---|---|---|
| Citation Counts | A/J | × | ✓ |
| Impact Factor | J | × | × |
| FNCSI [96] | J | Field | × |
| CiteScore [95] | J | × | × |
| SNIP [72] | J | Field | × |
| FWCI [30] | A | Filed | × |
| RCR [41] | A | Filed | × |
| TNCSI(*Ours*) | A/J | Filed and Value | × |
| aTNCSI(*Ours*) | A/J | Filed and Value | ✓ |

TABLE 2
**Metrics for Evaluating Scholar Impact of Papers:** "A" and "J" stand for article-level and journal-level. Filed normalized signifies that the metric can be utilized across fields.

As presented in Tab. 2, many works have tried to address the above-mentioned limitations. For example, Altmetric [98] collects evidence from the social web (X or Twitter, Facebook, etc.) to track and analyze the online attention and engagement that research outputs receive. While Altmetric provides valuable insights into the wider impact and engagement of research, bias toward the English language and none peer-reviewed further limit its broader application. Paper [80] compares two well-known article-level field-independent citation metrics, Field-Weighted Citation Impact (FWCI) and Relative Citation Ratio (RCR), and suggests they perform equally well in normalizing citations across research fields. However, these metrics require a pre-defined field (e.g. Scopus All Science Journal Classification category) and thus fail to accurately evaluate papers in newly emerging subfields. The Field Normalized Citation Success Index (FNCSI), as proposed in Ref [90], [96], is defined as the probability that a paper published in Journal A is cited more than a randomly selected paper published in Journal B. While FNCSI is robust, its reliance on pre-defined topic keywords, as well as its exclusive suitability for journal assessment, should be emphasized.

As the Leiden Manifesto [38] states: metrics should "*Account for variation by field in publication and citation practices*". We develop the concept of impact indicators, which is a measure used to gauge the impact of a certain paper in its field. The reason for using the term "indicator" rather than "metric" is we believe that the impact of a certain paper cannot be fully and accurately measured. While metrics like citation counts, h-index, or journal impact factors could indicate a paper's influence within the academic community, they fail to capture all aspects of its impact. For instance, a research paper might lead to significant advancements in theory, methods, or understanding in its field, none of which would necessarily be reflected in academic metrics. Similarly, a paper might contribute novel concepts or techniques that become influential over time but are not initially evident in citation counts. Furthermore, these metrics don't capture the quality of the research itself, such as the soundness of its methodology or the validity of its conclusions. Hence, the term "indicator" is favored, as it suggests a signal without asserting to encompass the entirety of the academic paper's contribution.

Based on the discussion, we propose the $TNCSI$ and $IEI$.

### 4.1.1 Topic Normalized Citation Success Index (TNCSI)

The Topic Normalized Citation Success Index (TNCSI) is a field-normalized article-level index that aims to assess the impact

Fig. 4. Distribution of scholarly citations: a comparative histogram of paper counts versus received citations across popular topics in PAMI. In the field of PAMI, citation distribution typically follows an exponential decay pattern, with a small subset of publications receiving the majority of citations. For a better presentation, the data obtained with a citation count greater than 5000 is ignored.

of research publications on a specific topic by normalizing the citation count to a scale ranging from 0 to 1.

Based on the selected $k$ papers and the corresponding citation counts $p_c$ for each paper $p$, we can calculate the discrete citation frequency distribution of the $k$ papers in a certain topic. We may further consider the distribution as a probability mass function:

$$P(X = x) = \frac{\text{Citation}_x}{k} \tag{1}$$

where $Citation_x$ represents the number of papers with $x$ citations. The $k$ papers related to the topic and their meta-data could be retrieved using a pre-defined topic keyword. through online scholar search engines or API, such as Semantic Scholar, CrossRef, or Google Scholar, as mentioned in Sec. 2.2. In addition, we can restrict the selection of $k$ papers by filtering out those not published within the specified timeframe to obtain a collection of $k'$ papers. This would allow the paper to be evaluated only with papers published in a specific time period (e.g., the same year), thus indicating the relative impact of the paper over a certain period.

We count the citations of papers across various topics with the help of Semantic Scholar API and plot them in Fig. 4. For each keyword, we retrieved up to 1000 (limited by Semantic Scholar) relevant literature using the API. Considering that the number of papers with $x$ citations generally follows an exponential decay, we utilize the maximum likelihood estimation method to fit $P(X = x)$ and obtain the probability density function (PDF) of a continuous exponential decay distribution:

$$f(x) = \lambda e^{-\lambda x}, x \geq 0 \tag{2}$$

where $f(x)$ represents the probability density at the value $x$, and $\lambda$ is the results obtained from the maximum likelihood estimation, representing the scale parameter controlling the scaling. Finally, the definite integral of $f(x)$ over the interval $[0, citeNum]$ gives us the desired $TNCSI$:

$$TNCSI = \int_0^{\text{citeNum}} \lambda e^{-\lambda x} \, dx \tag{3}$$

Specifically, $TNCSI_s$ indicates the relative impact of a paper compared to others published in the same year.

The $TNCSI$ demonstrates favorable mathematical properties and interpretability. The $TNCSI$ algorithm employs maximum likelihood estimation to convert the probability mass function into a probability density function. This process ensures that, in theory, the $TNCSI$ differentiates between papers with distinct citation counts, avoiding the assignment of identical values to them. $TNCSI$ has a physical meaning; it is the probability that the citation of the specific paper is greater than the citation of any other paper on the same topic. For example, a paper with a $TNCSI$ of 0.5 means it has more citations than half of the papers within the same topic.

Despite the advantages of $TNCSI$, it still faces a similar challenge to FWCI and RCR: the need to manually pre-define the topic. To address such a limitation, we propose adopting the ChatGPT [75] (gpt-3.5-turbo by default) to generate the topic keyword and retrieve related papers according to the keyword. Then, these retrieved $k$ papers will be used to calculate the discrete citation frequency distribution as afore-mentioned. We denote the resulted $TNCSI$ as $aTNCSI$, where "a" stands for adaptive.

ChatGPT is one of the most advanced and influential large language models in the field of natural language processing [121]. Equipped with state-of-the-art language understanding capabilities, ChatGPT has revolutionized the way we interact with AI-powered conversational systems by simply setting "system", "user", and "assistant" roles. The "system" role sets the conversation's behavior and initial context. It provides instructions to guide the assistant's responses. The "user" role represents the individual interacting with ChatGPT, who inputs messages to the assistant. The "assistant" role is the ChatGPT model itself which would respond based on the provided instructions and user input.

As illustrated in Fig. 5, we ask ChatGPT to identify the most representative topic keyword with the paper's title and abstract. In most scenarios, the generated topic word is sufficient to meet expectations, which can be further used as the keyword to retrieve papers from online scholar search engines. Optionally, one can set "System", "User", and "Assistant" roles before the final query to improve response quality and create more tailored interactions with the ChatGPT. In other word, a few-shot user-assistant pair prompts the ChatGPT with context on topic granularity. For example, a paper about the classification of irises by an improved CNN may have different perspectives. Some researchers focus more on the algorithm of the improved CNN, while others may be interested in classifying irises. Such ambiguity would likewise make it difficult for ChatGPT to identify the most representative topic keyword as expected. However, this could be addressed by providing the context which consists of (1) the identical prompt template with the replaced title and abstract as user input, and (2) the expected topic keyword as assistant output. By default, we adopt a well-known paper [26] as an example to guide models to generate the topic keyword for all papers.

Similar to other LLMs, ChatGPT performs various NLP tasks with user-provided natural language prompts. However, natural language prompts could be ambiguous, and even minor modifications can result in significantly different outputs. Thus, we follow the practices of LLM prompt engineering and carefully design the prompt to optimize the desired output. To determine the optimal prompt, we construct a dataset by manually annotating the topic keywords of 201 papers from various domains published later than the ChatGPT being trained and then compare the
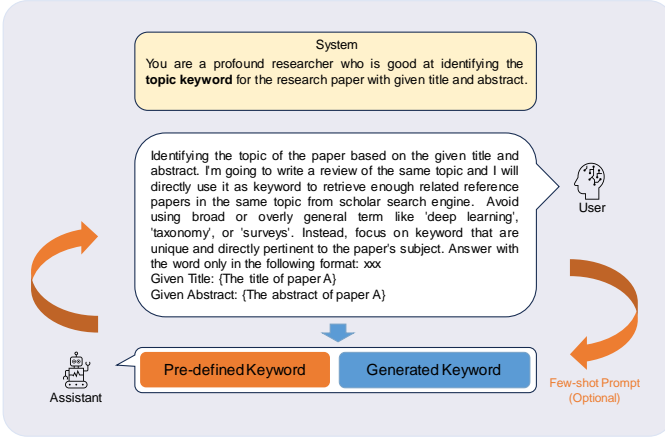
Fig. 5. Conceptual illustration of topic keyword generation process. Few-shot prompting may enhance the response quality of the large language model.

performance of multiple prompts on this dataset. The normalized edit distance [115] is adopted to measure the similarity between the GPT-generated keyword and our annotated keyword, where a lower value indicates a higher quality of the prompt. As can be seen from Tab. 3, some of the designed prompts achieve decent NED scores for papers in various domains.

### 4.1.2 Impact Evolution Index (IEI)

Imagine a scenario where two papers, A and B, receive the same number of citations. The number of new citations per month for A remains steady, whereas the number of new citations for B grows exponentially. In this context, while acknowledging the importance of Paper A, it is generally assumed that Paper B holds a greater reference value. Analyzing the popularity or citation trends of the literature may help researchers stay informed about the latest developments and identify potential areas for future research.

Most existing approaches treat the estimation of future citations as a sequence modeling task. Abrishami and Aliakbary propose employing the artificial neural network to predict long-term citations of a paper based on the number of its citations in the first few years after publication [1]. Zhao and Feng utilize graph structure representation and recurrent neural network modules to predict paper citation counts [120]. A more direct approach is to adopt polynomial fitting for scatter data, and calculate the sum of derivatives at each point. However, In practice, it has been observed that when dealing with data that lacks discernible distribution patterns, the polynomial fitting method tends to be sensitive to outliers, which can significantly compromise its robustness. Consequently, the numerical values obtained may not accurately reflect the underlying citation trend. In contrast to these series analysis-based methods, we propose a morphological theory-grounded Impact Evolution Index ($IEI$), which converts the citation trend into a clear and interpretable numerical value.

To calculate the $IEI$, it is necessary to obtain the citations of the paper first(See more in Sec. 2.2). Once the citation data is retrieved, we may create a sequence $Seq_{\text{citation}}$ about the number of citations. The $i \in \{0, 1, 2..., l\}$ item in the sequence $Seq_{\text{citation}}[i]$ represents the number of received citations in the $i_{\text{th}}$ month after the publication, where $l$ represents the number of months allocated for trend observation. Then, a sequence $Seq_{\text{time}}$ of the same

length as $Seq_{\text{citation}}$ is generated by enumerating from 0 to $l-1$, Typically, the minimum recommended value for $l$ is 6 or higher. This ensures that the data used for the analysis is adequately representative and the results are reliable. We match the items at the same positions in $Seq_{\text{time}}$ and $Seq_{\text{citation}}$ to determine a set of discrete coordinates $\{(Seq_{\text{time}}[i], Seq_{\text{citation}}[i])\}$, which serve as the control points for shaping the Bézier curve. A Bézier curve is a mathematical representation of smooth curves commonly used in computer graphics, image editing, and design software. The curve starts at the first control point and ends at the last control point, while the intermediate control points influence the curvature and direction of the curve. The number of control points determines the degree $n = l - 1$ of the curve.

$$C(t) = \sum_{i=0}^{n} B_{i,n}(t)P_i \tag{4}$$

$$B_{i,n}(t) = \binom{n}{i}(1-t)^{n-i}t^i, t \in [0,1] \tag{5}$$

where $B_{i,n}(t)$ represents the coefficient of the Bézier curve at a given parameter value $t$, which determines the position along the curve ($t = 0$ means the start and $t = 1$ means the end). $\binom{n}{i}$ is the binomial coefficient, also known as "n choose i". It represents the number of ways to choose $i$ elements from a set of $n$ elements. $P_i$ stands for the $i_{\text{th}}$ control point of the curve.

Given the continuity of the Bézier curve, we can compute its derivative as follows:

$$C'_{i,n}(t) = n \sum_{i=0}^{n-1} \left( \binom{n-1}{i}(1-t)^{n-1-i}t^i \right)(P_{i+1} - P_i) \tag{6}$$

Finally, the $IEI_{L_l}$ can be obtained by averaging the derivatives of $l = n + 1$ points, as shown in Eq. (7). Note that the value of $l$ can be configured flexibly to meet actual demands. In general, the longer the period analyzed, the more stable the citation trend becomes. We usually prefer to analyze the most recent 6 months of citations when constructing a Bézier curve of degree 5. Furthermore, different months can contribute differently to the $IEI$. For instance, if we desire closer months to have a greater impact, we can achieve this by adjusting the weighting coefficients, $w_i$, of the derivatives at different points to calculate their weighted averages (See in Eq. (8)). In addition, the instantaneous trend could be regarded as the derivative of the last month in the sequence. It can be obtained by setting $w_l = 1$ and the other weighting coefficients to 0. We denote the $IEI$ focused on the last month among the latest $l$ months (excluding the current month) as $IEI_{I_l}$ (See in Eq. (9)).

$$IEI_{L_l} = (C'_{0,l-1}(0) + \sum_{i=0}^{l-1-2} C'_{i,l-1}(1))/(l-1) \tag{7}$$

$$\begin{aligned} IEI_{W_l} = w_0 C'_{0,l-1}(0) + w_1 C'_{0,l-1}(1) \\ + \ldots + w_l C'_{l-1-2,l-1}(1) \end{aligned} \tag{8}$$

$$IEI_{I_l} = C'_{l-1-2,l-1}(1) \tag{9}$$

In general, the longer the period being analyzed, the more stable the citation trend becomes. We usually prefer to analyze the most recent 6 months of citations when constructing a Bézier curve of degree 5 and utilize the curve to calculate $IEI_{L_6}$, $IEI_{W_6}$, and $IEI_{I_6}$.

| NO. | User Prompt Content | Few-shot | NED↓ |
|---|---|---|---|
| 1 | Please analyze the title and abstract provided below and identify the main topic or central theme of the review paper. Focus on key term and the overall subject matter to determine the primary area of research or discussion.The output should be formatted as following: xxx | × | 0.75 |
| 2 | Given title and abstract, please provide the searching key phrase for me so that I can use it as keyword to search highly related papers from Google Scholar or Semantic Scholar. Please avoid responding with overly general keyword such as deep learning, taxonomy, or surveys, etc. Answer with the words only in the following format: xxx | × | 0.40 |
| 3 | Identifying the topic of the paper based on the given title and abstract. Avoid using broad or overly general term like 'deep learning', 'taxonomy', or 'surveys'. Instead, focus on keyword that is unique and directly pertinent to the paper's subject. Answer with the word only in the following format: xxx | × | 0.36 |
| 4 | Identifying the topic of the paper based on the given title and abstract. So that I can use it as keyword to search highly related papers from Semantic Scholar. Avoid using broad or overly general term like 'deep learning', 'taxonomy', or 'surveys'. Instead, focus on keyword that is unique and directly pertinent to the paper's subject. Answer with the word only in the following format: xxx | × | 0.32 |
| 5 | Identifying the topic of the paper based on the given title and abstract. I'm going to write a review of the same topic and I will directly use it as keyword to retrieve enough related reference papers in the same topic from scholar search engine. Avoid using broad or overly general term like 'deep learning', 'taxonomy', or 'surveys'. Instead, focus on keyword that are unique and directly pertinent to the paper's subject. Answer with the word only in the following format: xxx | × | 0.29 |
| 6 | Identifying the topic of the paper based on the given title and abstract. I'm going to write a review of the same topic and I will directly use it as keyword to retrieve enough related reference papers in the same topic from scholar search engine. Avoid using broad or overly general term like 'deep learning', 'taxonomy', or 'surveys'. Instead, focus on keyword that are unique and directly pertinent to the paper's subject. Answer with the word only in the following format: xxx | ✓ | **0.28** |

TABLE 3
Effectiveness of prompt engineering: comparison of various user prompts.

## 4.2 Quality Indicators

Assessing the quality and credibility of a literature review plays a pivotal role in the objective evaluation. Despite its significance, there have been minimal scholarly attempts to quantify the quality of literature reviews or any other types of publications. The paucity of such attempts may be attributed to the complexity and subjectivity involved in the process. However, the emergence of LLMs makes the quantitative quality evaluation for literature reviews no longer out of reach. By adopting an objective evaluation approach, we may move beyond the limitations of traditional metrics (such as simply regarding the citation numbers or the journal impact factor as the paper quality), thus providing a more objective and reproducible evaluation method.

We propose a set of indicators to evaluate literature reviews from the perspectives of reference quality and update urgency.

### 4.2.1 Reference Quality Measurement (RQM)

A literature review, in its essence, can not fabricate insights from a void. It fundamentally relies on the substance of existing references. Without a solid foundation of credible and high-quality sources, a literature review may lack the necessary building blocks to construct a meaningful analysis or argument. These sources provide the empirical evidence and theoretical context that ground the review, making the role of references indispensable in the creation of a substantial literature review.

The reference quality of a literature review is a multifaceted concept. It usually involves several direct factors such as credibility, relevance, breadth, and depth, etc. Quantitatively evaluating these direct elements poses significant challenges. On the one hand, the relevance of scientific literature is difficult to be precisely defined. Relying on co-citation analysis [92] or the similarity of paper embeddings both suffer from conceptual limitations. On the other hand, practical limitations are evident. For instance, the concept of breadth, also known as coverage, within a literature review, can theoretically be quantified as the ratio of the number of references to the total number of relevant references. However, accurately determining this ratio is challenging. It is difficult to

ascertain the complete number of relevant references either by keyword search or by citation network.

Due to the challenges of quantifying the reference quality using the above-mentioned direct factors, we consider an indirect quality indication of each reference with the assistance of the $TNCSI$ presented in Sec 4.1.1. Here, the $TNCSI$ can be regarded as a reference quality indicator based on user voting, which is a statistical result of numerous researchers' comprehensive analyses of those direct factors of a certain paper.

Timeliness also matters. A current and up-to-date literature review ensures that the most recent advancements, developments, and perspectives in a particular field are taken into consideration. This temporal relevance enhances the accuracy and effectiveness of research outcomes, as it reflects the current state of knowledge and understanding. By focusing on the currency of the references, we can gauge the extent to which the literature review incorporates the latest research and developments in the field.

To simultaneously consider the quality and timeliness of cited references, we propose modifying the Gompertz function to model the reference quality (see in Eq. (10)). The Gompertz function is characterized by its sigmoidal, which indicates a slow growth rate at the start and end of a time period, with a more rapid growth in the middle phase. This pattern is often observed in natural phenomena, such as species dynamics [16], tumor growth [99], etc.

$$RQM = 1 - e^{-\beta \cdot e^{-(1-ARQ) \cdot S_{mp}}} \tag{10}$$

where $\beta$ is the shift parameter, $ARQ$ stands for average reference quality, and $S_{mp}$ represents the median semester count of the reference age [94], defined as the period spanning from the publication dates of the cited references to the issuance of the review.

The calculation procedures of $ARQ$ are as follows: The first step is the extraction of the cited reference list. For most publications, their reference lists could be provided by Semantic Scholar API. For a small number of reviews, the reference list provided by Semantic Scholar may contain errors. In this case, al-

though there are powerful computer vision-based algorithms [14], [22] available for extracting the reference list within PDFs, our requirements are relatively simple and can be effectively met by relying on the heuristic algorithms or ChatGPT. More specifically, for literature review with a relatively fixed citation format, we can use the PDFMiner [79] to read text from PDF files and use heuristic rules to match citations. Alternatively, the text can be analyzed using ChatGPT to extract in-text citations. The second step is similar to the approach presented in Sec. 4.1.1, where the ChatGPT and a well-designed prompt (as presented in Fig. 5) are utilized to obtain the topic keyword of the review. Next, we calculate the $TNCSI$ for each reference in the list. To conserve computational resources, we avoid using the ChatGPT to generate keywords for each reference. Instead, the $TNCSI$ of all cited literature is calculated using a sharing topic keyword. Finally, the coverage can be further calculated in Eq. (11):

$$ ARQ = \frac{\sum_{i=1}^{N_R} TNCSI(Ref_i)}{N_R} \quad (11) $$

where $TNCSI(\cdot)$ refers to the $TNCSI$ value of the $i_{th}$ cited reference, and $N_R$ stands for the number of the reference. In certain instances, it has been noted that calculating $TNCSI_s$ for each cited literature is also reasonable. However, this paper primarily emphasizes the current impact of the cited references, hence the utilization of $TNCSI$ in this context.

The shift parameter $\beta$ can be set empirically or obtained statistically. For statistical calculation, we first examine the distribution of $S_{mp}$ and $ARQ$ across all papers within the RiPAMI database, and obtain their respective mean values $\overline{S_{mp}}$ and $\overline{ARQ}$. Then, the problem of asserting for $\beta$ is reconceptualized as an optimization problem. As shown in Eq. (12), the objective here is to identify the value of $\beta$ that maximizes the derivative of $RQM(\overline{S_{mp}}; \beta, \overline{ARQ})$, subject to the constraints of $\overline{S_{mp}} = 8$ and $\overline{ARQ} = 0.6$. Such an approach endows $RQM$ with a more discriminative nature. It should be noted that different fields may result in distinct values of $\beta$. For this study, the $\beta$ has been established as 20.

$$ \beta_{\text{opt}} = \arg \max_{\beta} \left( RQM'(\overline{S_{mp}}; \beta, \overline{ARQ}) \right) \quad (12) $$

The range of $RQM$ extends from 0 to 1, where values closer to 1 signify a higher quality of the referenced literature. As illustrated in Fig. 6 (a), when the $ARQ$ of a paper remains constant, an increase in the variable $S_mp$ will lead to a decrease in the $RQM$ value. Conversely, when $S_mp$ remains constant, a higher $ARQ$ will elevate the $RQM$ value.

### 4.2.2 Review Update Index (RUI)

The $RUI$ (Review Update Index) refers to the measure of the extent to which a literature review is required to be updated due to the iteration of technology, theory, etc. The index is related to both the literature itself and the research interests of the topic. Generally, a high update index suggests that a literature review is in need of an immediate update. Conversely, a lower update index implies that few advances have been made to the investigated field and the review is still up-to-date.

To evaluate the $RUI$, we may start with the coverage of references before and after publication. This coverage ratio can, to some extent, indicate the extent to which a review requires updating within its field. However, as mentioned earlier, accessing the coverage of a review is difficult. Fortunately, this problem is subtly avoided in calculating the ratio of relevant papers before and after publication. Assuming that the ratio of references containing the topic keyword in the title to all references is $R_k$, the total number of relevant articles can be estimated by dividing the number of articles containing those keywords retrieved from a search engine by $R_k$. Note that $R_k$ generally remains consistent before and after the publication, the CDR (Coverage Difference Ratio) can then be calculated in Eq. (13). The theoretical value range of CDR is greater than 0 to positive infinity. When the CDR of a review equals 1, it indicates that the current field has yielded new publications sufficient to constitute half of the literature referenced in the review.

$$ CDR = \frac{N_{pc} \cdot R_k}{R_k \cdot N_{mp}} = \frac{N_{pc}}{N_{mp}} \quad (13) $$

where $N_{mp}$ and $N_{pc}$ denote the number of relevant literature from the median publication date of the cited references to the publication date of the review, and from the publication date of the review to the current time, respectively.

In addition, similar to the inevitable process of biological aging, literature reviews also undergo a gradual aging process throughout time. Such passage of time bestows upon literature reviews increasing aging progress, where the degree of aging can be conceptualized as a normalized value of the academic impact already achieved. To further explore the aging of reviews in the field of PAMI, we conducted a statistical analysis of the yearly number of newly received citations for reviews published between 2015-2017 in RiPAMI. In contrast to earlier findings, however, the distribution of received citations of reviews over time follows a t-distribution rather than a log-normal distribution of the regular paper, as previously reported in Ref. [27], [66], [71]. Due to the insufficient duration of published sample data, the observation of citation-time trend curves is incomplete. Therefore, we conducted a three-degree polynomial fitting on the limited 6-year citation trend data and transformed the positive segment of the fitted curve into a PDF. To obtain the corresponding cumulative distribution function (CDF), we employed the cumulative trapezoidal numerical integration method for an approximate estimation. Thus, the review aging degree is given by:

$$ RAD(M_{pc}) = \int_0^{M_{pc}/12} (px^3 + qx^2 + rx + s)\,dx \quad (14) $$

where $M_{pc}$ denotes the duration in months from the publication of the review to the present, $p = -0.003$, $q = 0.001$, $r = 0.1267$, $s = 0.0129$ are the coefficients obtained by polynomial fitting. Please note that the integral symbol used here is for illustrative purposes only. The strict mathematical definition involves the accumulation of discrete trapezoidal areas.

Finally, the $RUI$ could be obtained by weighted summation of CDR and RAD:

$$ RUI = p \cdot CDR + q \cdot RAD(x) \quad (15) $$

where $p$ and $q$ are set to 10 and 5 in this paper, respectively. A sculptural visualization of the $RUI$'s contours is crystallized in Fig. 6 (b).

## 4.3 Quantitative Evaluations of Literature Reviews

We have carefully selected numerous reviews within the field of PAMI to ensure a representative sample for our evaluation.
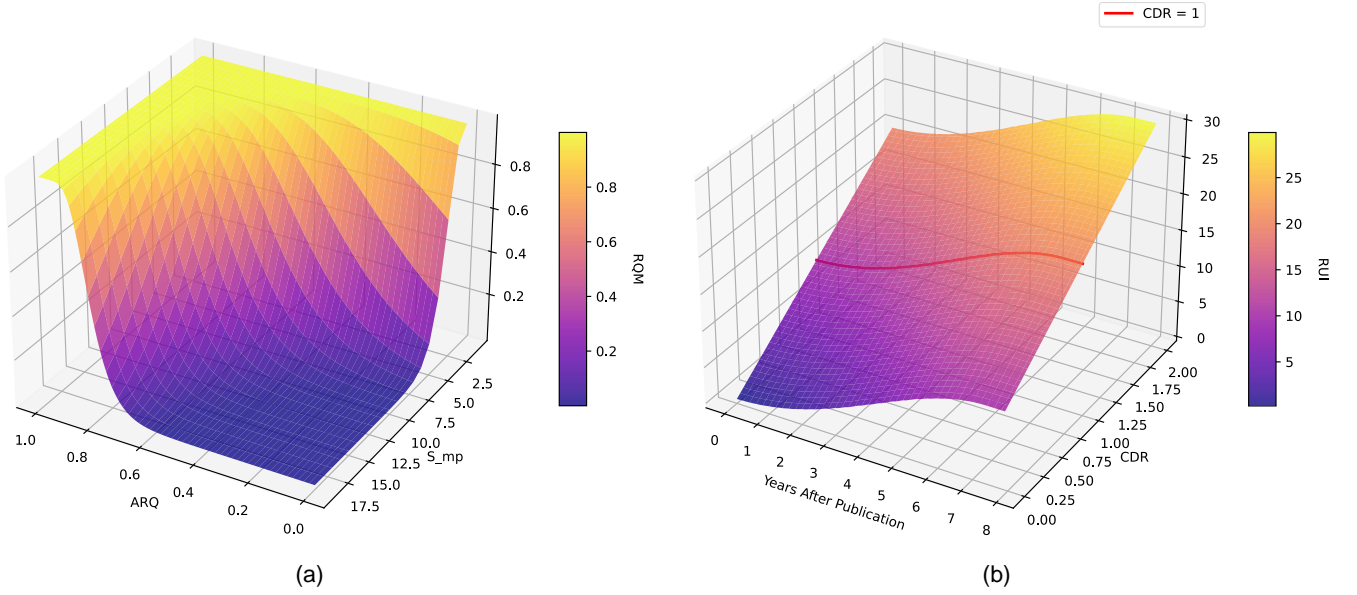
Fig. 6. 3D visualization of the proposed quality indicators. Panels (a) and (b) respectively depict the value landscapes of $RQM$ and $RUI$, shaped by two independent variables.

By employing the proposed quantitative indicators, we conduct a thorough assessment of these scholarly surveys, aiming to gauge their quality and potential influence within the field. In Fig. 7, we present a set of scatter plots depicting the references' quality of five randomly selected literature reviews published in two representative journals between 20015 and 2018. To prevent conflicts of interest, we use Journal A to represent a journal with an impact factor greater than 20. Journal B represents a journal with an impact factor below 5. The reference quality visualizations of journal A and journal B are presented in Fig. 7 (a), (b), respectively. The horizontal axis represents the reference age, and the vertical axis corresponds to the proposed $aTNCSI$. The color and size of the scatter points indicate the $IEI$ of the references. A positive $IEI$ value signifies a gradually increasing citation trend, resulting in warm-colored scatter points. The larger the size of the scatter point, the greater the value of $IEI$. Conversely, when the $IEI$ value is negative, the scatter points are cool-colored, and the size decreases as the value decreases. A closer examination of Fig. 7 reveals distinct distribution patterns among various journals depicted in their respective scatter plots. Journal with higher impact factors tends to exhibit clustered dots in the upper left corner, indicating a trend toward referencing more recent and influential sources. Conversely, journal with lower impact factors is prone to display more dispersed dots, with a notable distribution in the lower half of the graph. We argue that this discrepancy stems from authors' tendency to cite recent papers in their literature reviews. Given that these relatively recent papers only had limited received citations, their citation count might not adequately reflect the potential value of the literature. Therefore, those seasoned researchers submitting reviews to top-tier journals often render more precise judgments about the prospective significance of a certain paper, as evidenced in the figure.

More numerical results can be found in Tab. 4. As can be seen from the table, literature reviews [7], [18] of different topics that receive approximately the same number of citations exhibit notable discrepancies in their associated $aTNCSI$. This diver-

gence primarily stems from the varying levels of research interest across their respective topics. Surveys [101] in some emerging fields reveal a significant increase in their $IEI$, as expected. The majority of leading journals demonstrate high $RQM$ values, which reflect the effectiveness of high-standard peer review. In addition, even if the topics covered by the reviews [117], [119] are similar and their publication dates are relatively close, the proposed $RUI$ still provide some indication of the extent to which they need to be updated. Note that, given the time span of less than six months, the calculation of $RUI$ for certain papers is unavailable (denoted with "-").

## 5 HUMAN-AUTHORED VS. AI-GENERATED

### 5.1 Overview of AI-generated Literature Reviews

Traditionally, literature reviews have been manually conducted by researchers who analyze and synthesize scholarly sources to provide an overview of the current state of knowledge in a specific field. Recently, with the advancement of AI technologies, there's been a growing interest in leveraging artificial intelligence techniques, especially large language models, to automate or assist in the generation of the literature review. Typically, users are simply required to indicate their area of research interest, and the system will then automatically generate a literature review.

The automated creation of a literature review is a multidisciplinary endeavor that integrates knowledge from various fields. It relies not only on artificial intelligence technologies but also requires the merging of knowledge from other fields such as data science, bibliometrics, database engineering, etc. These components are instrumental in extracting, storing, and synthesizing relevant information from vast amounts of literature.

Early attempts in AI-generated literature reviews involve training language models, e.g. LLama [97], on a large corpus of academic papers, research articles, and other scholarly content. These models can then be used to generate coherent and contextually relevant text based on prompts or queries related to a specific

TABLE 4
Comparison of various reviews with the proposed indicators and metrics.

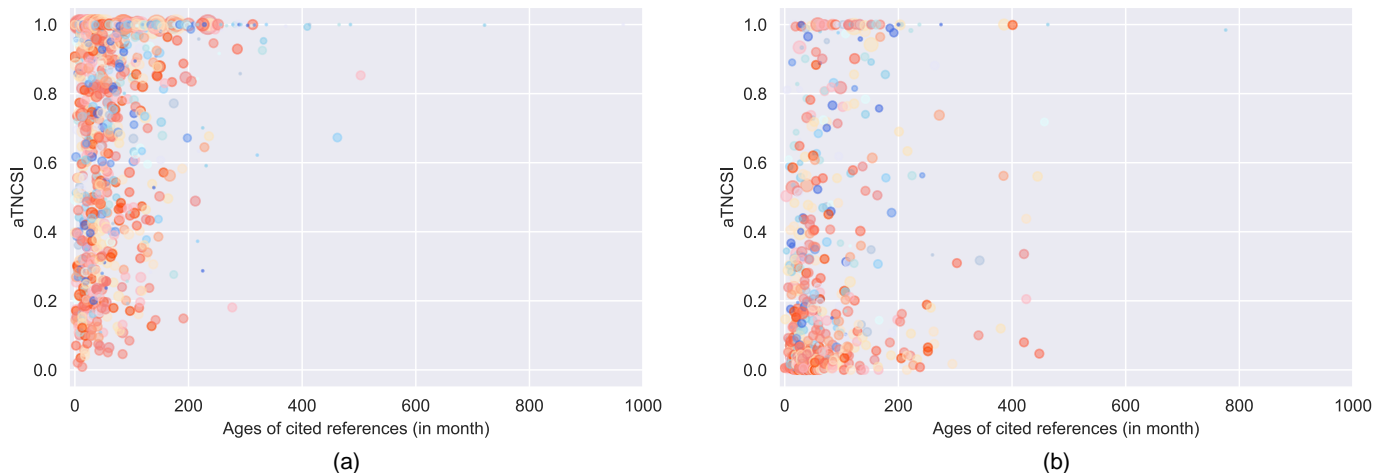| Title | Meta-Data | | | | Evaluation | | | |
|---|---|---|---|---|---|---|---|---|
| | Year | Citations | References | Topic | TNCSI↑ | IEI↑ | RQM↑ | RUI↓ |
| Object Detection with Deep Learning: A Review [122] | 2018 | 2980 | 254 | deep learning-based objection detection | 1.0 | −4.68 | 1.0 | 143.02 |
| Few-shot Object Detection: a Survey [6] | 2022 | 20 | 69 | few-shot object detection | 0.13 | −0.66 | 0.97 | 13.89 |
| Recent Few-shot Object Detection Algorithms: A Survey with Performance Comparison [53] | 2022 | 5 | 187 | few-shot object detection | 0.03 | 0.19 | 0.95 | 11.29 |
| Recent progresses on object detection: a brief review [118] | 2019 | 31 | 110 | object detection | 0.03 | 0.01 | 0.82 | 34.01 |
| A Survey of Modern Deep Learning Based Object Detection Models [116] | 2021 | 370 | 114 | object detection | 0.29 | −0.89 | 0.91 | 19.23 |
| A Survey on Curriculum Learning [102] | 2021 | 241 | 148 | curriculum learning | 0.77 | 1.8 | 0.42 | 9.38 |
| Review of Automatic Text Summarization Techniques & Methods [106] | 2020 | 99 | 117 | automatic text summarization | 0.73 | −0.41 | 0.2 | 17.37 |
| Automatic Text Summarization Methods: A Comprehensive Review [87] | 2022 | 15 | 102 | automatic text summarization | 0.18 | 0.01 | 0.02 | 3.7 |
| Graph Self-supervised Learning: A Survey [56] | 2021 | 256 | 184 | graph self-supervised learning | 0.88 | 3.3 | 1.0 | 111.17 |
| Self-supervised Learning on Graphs: Contrastive, Generative, or Predictive [107] | 2021 | 101 | 136 | graph self-supervised learning | 0.56 | −0.68 | 1.0 | 72.91 |
| A Review of Deep-learning-based Medical Image Segmentation Methods [54] | 2021 | 258 | 112 | medical image segmentation | 0.62 | −0.16 | 0.72 | 22.21 |
| Biomedical Image Segmentation: A Survey [5] | 2021 | 11 | 152 | biomedical image segmentation | 0.09 | −0.18 | 0.0 | 6.95 |
| A Survey of Methods, Datasets, and Evaluation Metrics for Visual Question Answering [88] | 2021 | 20 | 211 | visual question answering | 0.11 | 0.34 | 0.44 | 8.47 |
| From Image to Language: a Critical Analysis of Visual Question Answering (VQA) Approaches, Challenges, and Opportunities [42] | 2023 | 0 | 304 | visual question answering | 0.0 | - | 0.33 | 0.42 |
| A Survey on Vision Transformer [34] | 2020 | 767 | 328 | vision transformer | 1.0 | −0.27 | 1.0 | 748.26 |
| Transformers in Vision: A Survey [44] | 2021 | 1260 | 286 | transformers in computer vision | 1.0 | −4.98 | 1.0 | 312.82 |
| A Survey of Visual Transformers [58] | 2021 | 115 | 235 | visual transformers | 0.73 | 1.2 | 1.0 | 68.37 |
| A Survey on Efficient Vision Transformers: Algorithms, Techniques, and Performance Benchmarking [76] | 2023 | 0 | 81 | efficient vision transformers | 0.0 | - | 0.97 | 2.41 |
| A Survey on Large Language Model Based Autonomous Agents [101] | 2023 | 107 | 187 | LLM-based autonomous agents | 0.87 | 5.13 | 1.0 | 40.1 |
| A Systematic Review of Aspect-based Sentiment Analysis (ABSA): Domains, Methods, and Trends [25] | 2023 | 0 | 83 | aspect-based sentiment analysis | 0.0 | - | 0.2 | 0.39 |
| A Survey on Graph Diffusion Models: Generative AI in Science for Molecule, Protein and Material [119] | 2023 | 20 | 148 | graph diffusion models | 0.4 | 0.45 | 1.0 | 2.6 |
| A Survey on Audio Diffusion Mdels: Text to Speech Synthesis and Enhancement in Generative AI [117] | 2023 | 26 | 141 | audio diffusion models | 0.3 | 0.44 | 0.93 | 17.57 |
| Diffusion Models: a Comprehensive Survey of Methods and Applications [111] | 2022 | 355 | 393 | diffusion models | 0.64 | 0.85 | 0.98 | 16.38 |
| Geometric Deep Learning on Molecular Representations [7] | 2021 | 140 | 208 | molecular representations | 0.83 | 1.16 | 1.0 | 14.19 |
| Ensemble Deep Learning in Bioinformatics [18] | 2020 | 142 | 113 | ensemble bioinformatics | 0.67 | −0.21 | 0.82 | 23.43 |
| A Survey on Self-supervised Learning: Algorithms, Applications, and Future Trends [32] | 2023 | 1 | 272 | self-supervised learning | 0.0 | 0.21 | 0.78 | 5.18 |
| Self-Supervised Learning: Generative or Contrastive [55] | 2020 | 918 | 184 | self-supervised learning | 0.96 | 2.71 | 1.0 | 262.82 |

Fig. 7. Visualization of references quality of randomly selected reviews published in different journals. Panel (a) displays the citation quality distribution for a journal with an IF Score above 20, while Panel (b) shows the same for a journal with an IF Score below 5. This suggests that reviews published in journals with higher impact factors tend to reference sources of superior quality and greater timeliness.

research topic. However, given that a fully trained or fine-tuned language model has no access to the latest scholarly advances, these models would not generate a literature review that includes the latest scholarly materials. Thus, most existing advanced AI-generated literature review systems contain three main steps: knowledge retrieval, synthesis, and report.

The generation of literature review begins by gathering knowledge from various sources. There is no doubt that a high-quality knowledge retrieval procedure provides a richer context for the LLMs, resulting in the production of more accurate and decent responses. The system usually employs several predefined criteria such as keywords, publication date ranges, and citation numbers to locate, filter, and retrieve relevant resources to be synthesized. These resources are not limited to academic publications, but may also include blogs, official tutorials, and GitHub repositories. Once the relevant knowledge is retrieved, the AI system aims to synthesize the gathered information to create a coherent and comprehensive literature review. This step involves analyzing the retrieved content, identifying relationships between different sources, and organizing the information into meaningful sections. The system may employ LLMs and elaborate prompts to generate the text for each section or the entire survey at once. Finally, the system seamlessly transitions into the report generation stage, where a formatted review is crafted in a cohesive and well-structured style. This stage involves converting the synthesized content into a final report which is ready for presentation. One of the simplest and most common forms of presentation is plain text. For certain systems, the generated text is arranged in various sections based on the predefined setting. Such a multi-step procedure not only streamlines the literature review generation workflow but also ensures the production of acceptable content.

Though researchers are expecting to save significant time and effort in conducting comprehensive surveys by embracing AI-generated literature reviews, numerous concerns about ethical issues and fabrications in AI-generated academic content have been raised. A critical appraisal by Zybaczynska [124] highlights that current AI systems typically fall short in providing substantial, accurate information and critical discernment. Elali *et al.* [28] critically examines the profound challenges brought by the fabrication and falsification of AI-generated research, underscoring its signif-

icant impact on the scientific community. Furthermore, numerous academic publishers are cautious about content created by AI. For instance, journals such as "*Science*" and "*Nature*" prohibit the use of any text generated by ChatGPT or other automated tools in papers published in their issues. Nevertheless, we argue that the automated generation of literature reviews not intended for publication remains meaningful. It can assist researchers in staying rapidly abreast of the latest advancements in thriving fields.

## 5.2 Current State of AI-generated Literature Review

In this subsection, our investigation delves into various efforts that have been made in the field of employing AI techniques to generate literature reviews.

Most of the existing systems are only capable of generating literature reviews in plain text form and with a relatively fixed layout. A popular Github repository ChatPaper [43] retrieves papers based on a user-defined keyword with the use of the Google Scholar crawler. After a cosine similarity-based filtering of the retrieved papers, the selected papers are processed sequentially using ChatGPT. The output of Chatpaper is brief and concise, which usually contains a narrative description of the selected papers. Paper Digest [77] is an AI-based online platform that provides the service of automated review generation. Users can generate a literature review by setting specific keywords and further refining the cited source through constraints on publication dates. Similar to ChatPaper, the output of Paper Digest is also a narrative plain text description. Not all systems for generating literature reviews are fully automated. For example, Jenni AI [4] requires user interaction through "AI Command" during the review creation process, enabling the generation of highly customized literature reviews. Additionally, there are numerous online platforms [15], [86] and plugins available in the GPT Store that offer automated literature review generation services. Given the constraints on manuscript length, detailed introductions to these systems will not be provided.

Despite previous discussions indicating that certain publishers refuse to accept content generated by LLMs, scholars have still successfully published AI-generated literature reviews in some of academic journals. Aydın *et al.* investigate how well can LLMs

perform in the generation of the literature review. They employed well-known LLMs including ChatGPT and Google Bard to generate reviews on digital twin in healthcare [9] and meta-universes [8]. In their attempts, the Google Bard (or the ChatGPT) was first adopted to paraphrase the abstracts of papers within the last 3 years. Then, they designed several question prompts (such as "What is digital twin") to query the LLMs and rearranged the response in a formatted style.

While promising progress has been achieved in AI-generated literature reviews, there remains considerable room for improvement. The pursuit of automated literature reviews that rival the quality of those written by humans is an ongoing area of research, underscoring the need for continued exploration and development in this field. For example, most existing literature review generation systems are not capable of extracting information from tables or figures. As a result, tables and figures will not be included in the generated reviews. More details about the differences between AI and human-authored reviews will be presented in the next subsection.

### 5.3 Case Study: Comparing Human and AI in Conducting Literature Reviews

We crafted a case study aimed at elucidating the distinctive approaches employed by human researchers and their AI counterparts. To prevent potential information leakage and ensure a fair comparison due to the possibility of ChatGPT having been exposed to relevant reviews during training, we selected a newly emerging field, i.e. prompt learning for large visual models, which developed after September 2021 (the cutoff date for ChatGPT's training), as the focus of our investigation. Considering that literature reviews generated by AI lack academic influence, we only calculate $RQM$ (as introduced in Sec. 4.2) of them.

The data presented in Tab 5 reveals that, at present, most AI-generated reviews fail to match the quality of those crafted by human authors. There is still a gap in knowledge retrieval, synthesis, and reporting steps where AI systems lag behind human performance. Despite the high $RQM$ achieved by the PaperDigest, the overall quality of the generated review still falls short of the human-authored counterpart, especially given its reliance on only 5 references. Additionally, to our knowledge, few automated review systems can seamlessly integrate visual elements such as figures and tables. This landscape of predominantly plain text-based AI-generated reviews invites intriguing possibilities for future exploration and enhancement in the field.

It is noteworthy that despite the higher quality of the human-authored review [100], the timeline from draft to journal acceptance, as indicated by its preprint publication and acceptance dates, extends over six months. In contrast, AI systems can generate reviews within seconds. This delay implies that in rapidly advancing research areas, some of the newest findings might be overlooked by human authors.

## 6 CHALLENGES AND THE FUTURE OF THE LITERATURE REVIEW

Challenges and future opportunities for both human-authored and AI-generated reviews are discussed in this section.

### 6.1 Challenges

**High Rate of obsolescence** Knowledge in scientific research continually evolves and updates. Literature reviews require significant time and effort to compile. Due to the slower pace of review crafting, they may not reflect the latest research advancements in a timely manner, especially when these reviews are conducted by human authors and in need of a peer-review process.

**Information Overload** While our proposed indicators mitigate the problem of information overload to a certain degree, several challenges still await resolution. With the continuous growth of scientific research, the volume of literature is rapidly expanding. Collecting and screening a large number of publications within those blooming fields would cause inevitable incomplete searches and other similar issues. Both human and AI systems have to explore how to retrieve relevant literature more comprehensively and effectively.

**Bias** Literature reviews conducted in a specific language or region may unintentionally omit relevant studies published in other languages or regions. This language and geographical bias can limit the global perspective and generalizability of the review's findings. Such a bias appears not only in literature reviews written by humans but also in those generated by AI systems (even featuring multilingual capabilities).

### 6.2 Future

**AI Empowerment Dynamic Literature Review** In the rapidly advancing landscape of scientific research, there is a growing trend toward dynamic reviews generated in real time. While debates persist regarding the adherence of AI-generated content to academic ethics and its suitability for formal publication, the generated reviews for non-publishing purposes remain widely embraced. Given the labor- and time-intensive nature of conducting literature reviews, the exploration of AI technology for the automated generation of dynamic literature reviews stands as an area warranting further investigation.

**Automated Literature Appraisal** Literature appraisal is important for both human authors and AI systems. Although four quantitative indicators are proposed in this paper, this still makes it difficult to fully capture the academic impact and the quality of a certain review paper. In the future, with more powerful large multimodal models (LMMs) [3], [50], it will be possible to extract valuable information from non-textual modal content (e.g., images and tables) for automated literature appraisal.

**Open Science and Advanced Search Engines** The open science movement endorses the sharing and accessibility of literature data, thereby providing a greater number of resources for analysis and synthesis. It is recommended that the percentage of open-access papers should be further increased in the era of e-publishing. Furthermore, search engines that rely on semantic similarity, recommend systems, or co-citation networks rather than keyword matching are also encouraged. Such engines will primarily benefit both researchers and AI systems by enabling them to retrieve more relevant papers, thus enhancing the quality of any publications in all of the research fields.

## 7 CONCLUSION

This Analysis has provided a fresh perspective on the wealth of literature reviews in PAMI, introducing a systematic approach to categorize and evaluate them. In total, it presents four large language models-empowered quantitative evaluation indicators, a subjective evaluation that includes a typology for literature reviews, a comparison between human-authored and AI-generated

| Review | Reference | RQM↑ | Automation Level | Visual Elements | SALSA Analysis [31] |
|---|---|---|---|---|---|
| [100] by human authors | 160 | **0.99** | Manullay | ✓ | A typically narrative, method clustering-based review aims to offer valuable insights in visual prompt learning. The selection criteria for references are not specified, but it's clear that each reference has been thoroughly appraised. |
| Review by Jenni [4] | - | - | Semi-automated | ✗ | Automated searching and appraising references is not supported. The generation process relies on user interaction, and only narrative description is provided. |
| Review by ChatPaper [43] | 15 | 0 | Automated | ✗ | Retrieving relevant literature from arXiv based on multiple LLM-generated keywords with no appraisal step. Each section contains plain descriptions of the related reference. |
| Review by PaperDigest [77] | 5 | **0.99** | Automated | ✗ | Searching papers with the user-specified keyword. It seems to appraise the quality of references with private criteria. The generated content is more like a summary. |
| Review by askyourpdf [15] | 9 | 0.31 | Automated | ✗ | No official explanations are found for how to retrieve references and appraise the quality of the literature. The generated review includes a brief analysis and description of the related concept, current state of development, and gaps. |

TABLE 5
Comparisons between human-authored and AI-generated literature reviews.

reviews, and a meta-data database named RiPAMI, accompanied by a dataset of review topic key phrases.

The proposed evaluation indicators offer an innovative alternative to traditional bibliometric analysis, enabling a cost-efficient, field-normalized, and real-time assessment of literature reviews' quality and impact. These bibliometric indicators not only furnish clear numerical hints for human researchers but could also offer substantial assistance in the appraisal processes of AI-generated review systems.

Subject evaluations for reviews in PAMI are also provided. Numerous representative reviews of popular research topics are investigated to offer readers a brief overview of these exemplary reviews. In addition, a typology for reviews is introduced. It category the reviews into three major types which are method clustering-based, challenge-oriented, and hybrid literature review. Such typology enhances the understanding of how literature reviews are organized and synthesized.

The study highlights the differences between human-authored and AI-generated literature reviews, pinpointing significant gaps in retrieval, synthesis, and reporting capabilities of AI-generated reviews compared to those crafted by humans. The insights gained from this comparison not only help to understand the current state of AI-generated reviews in PAMI but also suggest how they might evolve with advancing technology.

By constructing the RiPAMI database, we analyze approximately 3,000 review samples in the field of PAMI, which provides statistical support for the proposed quantitative indicators. Furthermore, a dataset of topic keywords is manually annotated to validate the effectiveness of various prompts, enabling the selection of the optimal prompt for more accurate data retrieval.

To further support and expand upon this work, we release the code framework that encompasses functionalities including metadata retrieval, indicator calculation, and data analysis, etc. While it is designed for analysis of literature reviews in PAMI, the framework is universally adaptable, catering to the expansive needs of researchers from disparate fields.

# REFERENCES

[1] A. Abrishami and S. Aliakbary. Predicting citation counts based on deep neural network learning techniques. *Journal of Informetrics*, 13(2):485–499, 2019.

[2] T. Al-Moslmi, M. G. Ocaña, A. L. Opdahl, and C. Veres. Named entity extraction for knowledge graphs: A literature overview. *IEEE Access*, 8:32862–32881, 2020.

[3] J.-B. Alayrac, J. Donahue, P. Luc, A. Miech, I. Barr, Y. Hasson, K. Lenc, A. Mensch, K. Millican, M. Reynolds, et al. Flamingo: a visual language model for few-shot learning. *Advances in Neural Information Processing Systems*, 35:23716–23736, 2022.

[4] Altum Inc. Jenni ai - your ai research assistant. https://jenni.ai/, 2023. Accessed 25 December 2023.

[5] Y. Alzahrani and B. Boufama. Biomedical image segmentation: a survey. *SN Computer Science*, 2:1–22, 2021.

[6] S. Antonelli, D. Avola, L. Cinque, D. Crisostomi, G. L. Foresti, F. Galasso, M. R. Marini, A. Mecca, and D. Pannone. Few-shot object detection: A survey. *ACM Computing Surveys (CSUR)*, 54(11s):1–37, 2022.

[7] K. Atz, F. Grisoni, and G. Schneider. Geometric deep learning on molecular representations. *Nature Machine Intelligence*, 3(12):1023–1032, 2021.

[8] Ö. AYDIN. Google bard generated literature review: metaverse. *Journal of AI*, 7(1):1–14, 2023.

[9] Ö. Aydın and E. Karaarslan. Openai chatgpt generated literature review: Digital twin in healthcare. *Available at SSRN 4308687*, 2022.

[10] A. Bandini and J. Zariffa. Analysis of the hands in egocentric vision: A survey. *IEEE transactions on pattern analysis and machine intelligence*, 2020.

[11] J. Beel and B. Gipp. Google scholar's ranking algorithm: an introductory overview. In *Proceedings of the 12th international conference on scientometrics and informetrics (ISSI'09)*, volume 1, pages 230–241. Rio de Janeiro (Brazil), 2009.

[12] J. Beel and B. Gipp. Academic search engine spam and google scholar's resilience against it. *Journal of electronic publishing*, 13(3), 2010.

[13] J. Beel and B. Gipp. On the robustness of google scholar against spam. In *Proceedings of the 21st ACM Conference on Hypertext and Hypermedia*, pages 297–298, 2010.

[14] L. Blecher, G. Cucurull, T. Scialom, and R. Stojnic. Nougat: Neural optical understanding for academic documents, 2023.

[15] BlockTechnology OÜ. Ai-powered literature review generator. https://askyourpdf.com/tools/literature-review-writer, 2023. Accessed 25 December 2023.

[16] R. C. Bruce. Application of the gompertz function in studies of growth in dusky salamanders (plethodontidae: Desmognathus). *Copeia*, 104(1):94–100, 2016.

[17] M. Brzezinski. Power laws in citation distributions: evidence from scopus. *Scientometrics*, 103:213–228, 2015.

[18] Y. Cao, T. A. Geddes, J. Y. H. Yang, and P. Yang. Ensemble deep learning in bioinformatics. *Nature Machine Intelligence*, 2(9):500–508, 2020.

[19] M. A. Chandra and S. Bedi. Survey on svm and their application in image classification. *International Journal of Information Technology*, 13:1–11, 2021.

[20] L. Chen, S. Li, Q. Bai, J. Yang, S. Jiang, and Y. Miao. Review of image classification algorithms based on convolutional neural networks. *Remote Sensing*, 13(22):4712, 2021.

[21] G. Cheng, X. Yuan, X. Yao, K. Yan, Q. Zeng, X. Xie, and J. Han. Towards large-scale small object detection: Survey and benchmarks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.

[22] H. Cheng, P. Zhang, S. Wu, J. Zhang, Q. Zhu, Z. Xie, J. Li, K. Ding, and L. Jin. M6doc: A large-scale multi-format, multi-type, multi-layout, multi-language, multi-annotation category dataset for modern document layout analysis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15138–15147, 2023.

[23] H. M. Cooper. A taxonomy of literature reviews., 1985.

[24] H. M. Cooper. Organizing knowledge syntheses: A taxonomy of literature reviews. *Knowledge in society*, 1(1):104, 1988.

[25] P. Denny, K. Taskova, J. Wicker, et al. A systematic review of aspect-based sentiment analysis (absa): Domains, methods, and trends. *arXiv preprint arXiv:2311.10777*, 2023.

[26] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2020.

[27] L. Egghe et al. Citation age data and the obsolescence function: Fits and explanations. *Information Processing & Management*, 28(2):201–217, 1992.

[28] F. R. Elali and L. N. Rachid. Ai-generated research paper fabrication and plagiarism in the scientific community. *Patterns*, 4(3), 2023.

[29] A. S. for Cell Biology et al. San francisco declaration on research assessment (dora), 2012.

[30] FWCI. Field-weighted citation impact. https://libguides.usc.edu.au/researchmetrics/researchmetrics-field-weighted-citation-impact. Accessed 25 December 2023.

[31] M. J. Grant and A. Booth. A typology of reviews: an analysis of 14 review types and associated methodologies. *Health information & libraries journal*, 26(2):91–108, 2009.

[32] J. Gui, T. Chen, Q. Cao, Z. Sun, H. Luo, and D. Tao. A survey on self-supervised learning: Algorithms, applications, and future trends. *arXiv preprint arXiv:2301.05712*, 2023.

[33] Y. Guo, H. Wang, Q. Hu, H. Liu, L. Liu, and M. Bennamoun. Deep learning for 3d point clouds: A survey. *IEEE transactions on pattern analysis and machine intelligence*, 43(12):4338–4364, 2020.

[34] K. Han, Y. Wang, H. Chen, X. Chen, J. Guo, Z. Liu, Y. Tang, A. Xiao, C. Xu, Y. Xu, Z. Yang, Y. Zhang, and D. Tao. A survey on vision transformer. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(1):87–110, 2023.

[35] S. Hao, Y. Zhou, and Y. Guo. A brief survey on semantic segmentation with deep learning. *Neurocomputing*, 406:302–321, 2020.

[36] T. Hao, X. Li, Y. He, F. L. Wang, and Y. Qu. Recent progress in leveraging deep learning methods for question answering. *Neural Computing and Applications*, pages 1–19, 2022.

[37] J. E. Harmon. Evolution of the scientific paper. Technical report, Argonne National Lab., IL (United States), 1992.

[38] D. Hicks, P. Wouters, L. Waltman, S. De Rijcke, and I. Rafols. Bibliometrics: the leiden manifesto for research metrics. *Nature*, 520(7548):429–431, 2015.

[39] I. A. Huijben, W. Kool, M. B. Paulus, and R. J. Van Sloun. A review of the gumbel-max trick and its extensions for discrete stochasticity in machine learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(2):1353–1371, 2022.

[40] T. Hussain, K. Muhammad, W. Ding, J. Lloret, S. W. Baik, and V. H. C. de Albuquerque. A comprehensive survey of multi-view video summarization. *Pattern Recognition*, 109:107567, 2021.

[41] B. I. Hutchins, X. Yuan, J. M. Anderson, and G. M. Santangelo. Relative citation ratio (rcr): a new metric that uses citation rates to measure influence at the article level. *PLoS biology*, 14(9):e1002541, 2016.

[42] M. F. Ishmam, M. S. H. Shovon, M. Mridha, and N. Dey. From image to language: A critical analysis of visual question answering (vqa) approaches, challenges, and opportunities. *arXiv preprint arXiv:2311.00308*, 2023.

[43] kaixindelele. Chatpaper, 2023.

[44] S. Khan, M. Naseer, M. Hayat, S. W. Zamir, F. S. Khan, and M. Shah. Transformers in vision: A survey. *ACM computing surveys (CSUR)*, 54(10s):1–41, 2022.

[45] L. Langfeldt, I. Reymert, and D. W. Aksnes. The role of metrics in peer assessments. *Research Evaluation*, 30(1):112–126, 2021.

[46] J. Li, A. Sun, J. Han, and C. Li. A survey on deep learning for named entity recognition. *IEEE Transactions on Knowledge and Data Engineering*, 34(1):50–70, 2020.

[47] K. Li, G. Wan, G. Cheng, L. Meng, and J. Han. Object detection in optical remote sensing images: A survey and a new benchmark. *ISPRS journal of photogrammetry and remote sensing*, 159:296–307, 2020.

[48] Z. Li, F. Liu, W. Yang, S. Peng, and J. Zhou. A survey of convolutional neural networks: analysis, applications, and prospects. *IEEE transactions on neural networks and learning systems*, 2021.

[49] G. Litjens, T. Kooi, B. E. Bejnordi, A. A. A. Setio, F. Ciompi, M. Ghafoorian, J. A. Van Der Laak, B. Van Ginneken, and C. I. Sánchez. A survey on deep learning in medical image analysis. *Medical image analysis*, 42:60–88, 2017.

[50] H. Liu, C. Li, Q. Wu, and Y. J. Lee. Visual instruction tuning, 2023.

[51] L. Liu, W. Ouyang, X. Wang, P. Fieguth, J. Chen, X. Liu, and M. Pietikäinen. Deep learning for generic object detection: A survey. *International journal of computer vision*, 128:261–318, 2020.

[52] P. Liu, Y. Guo, F. Wang, and G. Li. Chinese named entity recognition: The state of the art. *Neurocomputing*, 473:37–53, 2022.

[53] T. Liu, L. Zhang, Y. Wang, J. Guan, Y. Fu, J. Zhao, and S. Zhou. Recent few-shot object detection algorithms: A survey with performance comparison. *ACM Transactions on Intelligent Systems and Technology*, 14(4):1–36, 2023.

[54] X. Liu, L. Song, S. Liu, and Y. Zhang. A review of deep-learning-based medical image segmentation methods. *Sustainability*, 13(3):1224, 2021.

[55] X. Liu, F. Zhang, Z. Hou, L. Mian, Z. Wang, J. Zhang, and J. Tang. Self-supervised learning: Generative or contrastive. *IEEE transactions on knowledge and data engineering*, 35(1):857–876, 2021.

[56] Y. Liu, M. Jin, S. Pan, C. Zhou, Y. Zheng, F. Xia, and P. S. Yu. Graph self-supervised learning: A survey. *IEEE Transactions on Knowledge and Data Engineering*, 35(6):5879–5900, 2023.

[57] Y. Liu, Y. Zhang, Y. Wang, F. Hou, J. Yuan, J. Tian, Y. Zhang, Z. Shi, J. Fan, and Z. He. A survey of visual transformers. *IEEE Transactions on Neural Networks and Learning Systems*, 2023.

[58] Y. Liu, Y. Zhang, Y. Wang, F. Hou, J. Yuan, J. Tian, Y. Zhang, Z. Shi, J. Fan, and Z. He. A survey of visual transformers. *IEEE Transactions on Neural Networks and Learning Systems*, pages 1–21, 2023.

[59] J. Lu, P. Gong, J. Ye, J. Zhang, and C. Zhang. A survey on machine learning from few samples. *Pattern Recognition*, 139:109480, 2023.

[60] J. Ma, J. Chen, M. Ng, R. Huang, Y. Li, C. Li, X. Yang, and A. L. Martel. Loss odyssey in medical image segmentation. *Medical Image Analysis*, 71:102035, 2021.

[61] G. R. Machado, E. Silva, and R. R. Goldschmidt. Adversarial machine learning in image classification: A survey toward the defender's perspective. *ACM Computing Surveys (CSUR)*, 55(1):1–38, 2021.

[62] Z. Mai, R. Li, J. Jeong, D. Quispe, H. Kim, and S. Sanner. Online continual learning in image classification: An empirical survey. *Neurocomputing*, 469:28–51, 2022.

[63] M. Malik, M. K. Malik, K. Mehmood, and I. Makhdoom. Automatic speech recognition: a survey. *Multimedia Tools and Applications*, 80:9411–9457, 2021.

[64] M. Masana, X. Liu, B. Twardowski, M. Menta, A. D. Bagdanov, and J. Van De Weijer. Class-incremental learning: survey and performance evaluation on image classification. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(5):5513–5533, 2022.

[65] N. Maslej, L. Fattorini, E. Brynjolfsson, J. Etchemendy, K. Ligett, T. Lyons, J. Manyika, H. Ngo, J. C. Niebles, V. Parli, Y. Shoham, R. Wald, J. Clark, and R. Perrault. The ai index 2023 annual report. Technical report, AI Index Steering Committee, Institute for Human-Centered AI, Stanford University, Stanford, CA, April 2023.

[66] E. Matricciani. The probability distribution of the age of references in engineering papers. *IEEE Transactions on Professional Communication*, 34(1):7–12, 1991.

[67] M. A. Mazurowski, H. Dong, H. Gu, J. Yang, N. Konz, and Y. Zhang. Segment anything model for medical image analysis: an experimental study. *Medical Image Analysis*, 89:102918, 2023.

[68] R. K. Merton. The matthew effect in science: The reward and communication systems of science are considered. *Science*, 159(3810):56–63, 1968.

[69] B. Min, H. Ross, E. Sulem, A. P. B. Veyseh, T. H. Nguyen, O. Sainz, E. Agirre, I. Heintz, and D. Roth. Recent advances in natural language processing via large pre-trained language models: A survey. *ACM Computing Surveys*, 2021.

[70] S. Minaee, Y. Boykov, F. Porikli, A. Plaza, N. Kehtarnavaz, and D. Terzopoulos. Image segmentation using deep learning: A survey. *IEEE transactions on pattern analysis and machine intelligence*, 44(7):3523–3542, 2021.

[71] H. F. Moed. Statistical relationships between downloads and citations at the level of individual documents within a single journal. *Journal*

*of the American Society for Information Science and Technology*, 56(10):1088–1097, 2005.

[72] H. F. Moed. Measuring contextual citation impact of scientific journals. *Journal of informetrics*, 4(3):265–277, 2010.

[73] D. Nadeau and S. Sekine. A survey of named entity recognition and classification. *Lingvisticae Investigationes*, 30:3–26, 2007.

[74] Z. Nasar, S. W. Jaffry, and M. K. Malik. Named entity recognition and relation extraction: State-of-the-art. *ACM Computing Surveys (CSUR)*, 54(1):1–39, 2021.

[75] OpenAI. Chatgpt: optimizing language models for dialogue. https://openai.com/blog/chatgpt/, 2022. Accessed 25 December 2023.

[76] L. Papa, P. Russo, I. Amerini, and L. Zhou. A survey on efficient vision transformers: algorithms, techniques, and performance benchmarking. *arXiv preprint arXiv:2309.02031*, 2023.

[77] Paper Digest. Paper digest ai-powered research platform. https://www.paperdigest.org/review/, 2023. Accessed 25 December 2023.

[78] G. Paré, M.-C. Trudel, M. Jaana, and S. Kitsiou. Synthesizing information systems knowledge: A typology of literature reviews. *Information & Management*, 52(2):183–199, 2015.

[79] PDFMiner. PDFMiner.six. [Online; accessed 14-Nov-2023].

[80] A. Purkayastha, E. Palmaro, H. J. Falk-Krzesinski, and J. Baas. Comparison of two article-level, field-independent citation metrics: Field-weighted citation impact (fwci) and relative citation ratio (rcr). *Journal of Informetrics*, 13(2):635–642, 2019.

[81] R. Qian, X. Lai, and X. Li. 3d object detection for autonomous driving: A survey. *Pattern Recognition*, 130:108796, 2022.

[82] X. Qu, Y. Gu, Q. Xia, Z. Li, Z. Wang, and B. Huai. A survey on arabic named entity recognition: Past, recent advances, and future trends. *arXiv preprint arXiv:2302.03512*, 2023.

[83] I. Qureshi, J. Yan, Q. Abbas, K. Shaheed, A. B. Riaz, A. Wahid, M. W. J. Khan, and P. Szczuko. Medical image segmentation using deep semantic-based methods: A review of techniques, applications and emerging trends. *Information Fusion*, 90:316–352, 2023.

[84] W. Rawat and Z. Wang. Deep convolutional neural networks for image classification: A comprehensive review. *Neural computation*, 29(9):2352–2449, 2017.

[85] L. Schmarje, M. Santarossa, S.-M. Schröder, and R. Koch. A survey on semi-, self-and unsupervised learning for image classification. *IEEE Access*, 9:82146–82168, 2021.

[86] Seamlees. Seamless - ai literature review tool for scientific research. https://seaml.es/, 2023. Accessed 25 December 2023.

[87] G. Sharma and D. Sharma. Automatic text summarization methods: A comprehensive review. *SN Computer Science*, 4(1):33, 2022.

[88] H. Sharma and A. S. Jalal. A survey of methods, datasets and evaluation metrics for visual question answering. *Image and Vision Computing*, 116:104327, 2021.

[89] W. Shen, Z. Peng, X. Wang, H. Wang, J. Cen, D. Jiang, L. Xie, X. Yang, and Q. Tian. A survey on label-efficient deep image segmentation: Bridging the gap between weak supervision and dense prediction. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.

[90] Z. Shen, S. Tong, F. Chen, and L. Yang. The utilization of paper-level classification system on the evaluation of journal impact. *arXiv e-prints*, pages arXiv–2006, 2020.

[91] N. Siddique, S. Paheding, C. P. Elkin, and V. Devabhaktuni. U-net and its variants for medical image segmentation: A review of theory and applications. *Ieee Access*, 9:82031–82057, 2021.

[92] H. Small. Co-citation in the scientific literature: A new measure of the relationship between two documents. *Journal of the American Society for information Science*, 24(4):265–269, 1973.

[93] K. B. Smith. Typologies, taxonomies, and the benefits of policy classification. *Policy studies journal*, 30(3):379–395, 2002.

[94] A. G. Stacey. Ages of cited references and growth of scientific knowledge: an explication of the gamma distribution in business and management disciplines. *Scientometrics*, 126(1):619–640, 2021.

[95] J. A. Teixeira da Silva. Citescore: Advances, evolution, applications, and limitations. *Publishing Research Quarterly*, 36(3):459–468, 2020.

[96] S. Tong, F. Chen, L. Yang, and Z. Shen. Novel utilization of a paper-level classification system for the evaluation of journal impact: An update of the cas journal ranking. *Quantitative Science Studies*, pages 1–16, 2023.

[97] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.

[98] N. S. Trueger, B. Thoma, C. H. Hsu, D. Sullivan, L. Peters, and M. Lin. The altmetric score: a new measure for article-level dissemination and impact. *Annals of emergency medicine*, 66(5):549–553, 2015.

[99] C. Vaghi, A. Rodallec, R. Fanciullino, J. Ciccolini, J. P. Mochel, M. Mastri, C. Poignard, J. M. Ebos, and S. Benzekry. Population modeling of tumor growth curves and the reduced gompertz model improve prediction of the age of experimental tumors. *PLoS computa-*

[100] J. Wang, Z. Liu, L. Zhao, Z. Wu, C. Ma, S. Yu, H. Dai, Q. Yang, Y. Liu, S. Zhang, et al. Review of large vision models and visual prompt engineering. *Meta-Radiology*, page 100047, 2023.

[101] L. Wang, C. Ma, X. Feng, Z. Zhang, H. Yang, J. Zhang, Z. Chen, J. Tang, X. Chen, Y. Lin, et al. A survey on large language model based autonomous agents. *arXiv preprint arXiv:2308.11432*, 2023.

[102] X. Wang, Y. Chen, and W. Zhu. A survey on curriculum learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(9):4555–4576, 2021.

[103] Y. Wang, C. Albrecht, N. A. A. Braham, L. Mou, and X. Zhu. Self-supervised learning in remote sensing: A review. *IEEE Geoscience and Remote Sensing Magazine*, 2022.

[104] M. Wankhade, A. C. S. Rao, and C. Kulkarni. A survey on sentiment analysis methods, applications, and challenges. *Artificial Intelligence Review*, 55(7):5731–5780, 2022.

[105] R. Weegar, A. Pérez, A. Casillas, and M. Oronoz. Recent advances in swedish and spanish medical entity recognition in clinical texts using deep neural approaches. *BMC medical informatics and decision making*, 19(7):1–14, 2019.

[106] A. P. Widyassari, S. Rustad, G. F. Shidik, E. Noersasongko, A. Syukur, A. Affandy, et al. Review of automatic text summarization techniques & methods. *Journal of King Saud University-Computer and Information Sciences*, 34(4):1029–1046, 2022.

[107] L. Wu, H. Lin, C. Tan, Z. Gao, and S. Z. Li. Self-supervised learning on graphs: Contrastive, generative, or predictive. *IEEE Transactions on Knowledge and Data Engineering*, 2021.

[108] S. Xun, D. Li, H. Zhu, M. Chen, J. Wang, J. Li, M. Chen, B. Wu, H. Zhang, X. Chai, et al. Generative adversarial networks in medical image segmentation: A review. *Computers in biology and medicine*, 140:105063, 2022.

[109] A. Yadav and D. K. Vishwakarma. Sentiment analysis using deep learning architectures: a review. *Artificial Intelligence Review*, 53(6):4335–4385, 2020.

[110] L. Yang, H. Jiang, Q. Song, and J. Guo. A survey on long-tailed visual recognition. *International Journal of Computer Vision*, 130(7):1837–1872, 2022.

[111] L. Yang, Z. Zhang, Y. Song, S. Hong, R. Xu, Y. Zhao, W. Zhang, B. Cui, and M.-H. Yang. Diffusion models: A comprehensive survey of methods and applications. *ACM Computing Surveys*, 56(4):1–39, 2023.

[112] J. Yu, H. Yin, X. Xia, T. Chen, J. Li, and Z. Huang. Self-supervised learning for recommender systems: A survey. *IEEE Transactions on Knowledge and Data Engineering*, 2023.

[113] Y. Yu, X. Si, C. Hu, and J. Zhang. A review of recurrent neural networks: Lstm cells and network architectures. *Neural computation*, 31(7):1235–1270, 2019.

[114] L. Yue, W. Chen, X. Li, W. Zuo, and M. Yin. A survey of sentiment analysis in social media. *Knowledge and Information Systems*, 60:617–663, 2019.

[115] L. Yujian and L. Bo. A normalized levenshtein distance metric. *IEEE transactions on pattern analysis and machine intelligence*, 29(6):1091–1095, 2007.

[116] S. S. A. Zaidi, M. S. Ansari, A. Aslam, N. Kanwal, M. Asghar, and B. Lee. A survey of modern deep learning based object detection models. *Digital Signal Processing*, 126:103514, 2022.

[117] C. Zhang, C. Zhang, S. Zheng, M. Zhang, M. Qamar, S.-H. Bae, and I. S. Kweon. A survey on audio diffusion models: Text to speech synthesis and enhancement in generative ai. *arXiv preprint arXiv:2303.13336*, 2, 2023.

[118] H. Zhang and X. Hong. Recent progresses on object detection: a brief review. *Multimedia Tools and Applications*, 78:27809–27847, 2019.

[119] M. Zhang, M. Qamar, T. Kang, Y. Jung, C. Zhang, S.-H. Bae, and C. Zhang. A survey on graph diffusion models: Generative ai in science for molecule, protein and material. *arXiv preprint arXiv:2304.01565*, 2023.

[120] Q. Zhao and X. Feng. Utilizing citation network structure to predict paper citation counts: A deep learning approach. *Journal of Informetrics*, 16(1):101235, 2022.

[121] W. X. Zhao, K. Zhou, J. Li, T. Tang, X. Wang, Y. Hou, Y. Min, B. Zhang, J. Zhang, Z. Dong, et al. A survey of large language models. *arXiv preprint arXiv:2303.18223*, 2023.

[122] Z.-Q. Zhao, P. Zheng, S.-t. Xu, and X. Wu. Object detection with deep learning: A review. *IEEE transactions on neural networks and learning systems*, 30(11):3212–3232, 2019.

[123] C. Zhou, Q. Li, C. Li, J. Yu, Y. Liu, G. Wang, K. Zhang, C. Ji, Q. Yan, L. He, et al. A comprehensive survey on pretrained foundation models: A history from bert to chatgpt. *arXiv preprint arXiv:2302.09419*, 2023.

[124] J. Zybaczynska, M. Norris, S. Modi, J. Brennan, P. Jhaveri, T. J. Craig, and T. Al-Shaikhly. Artificial intelligence-generated scientific literature-a critical appraisal. *The Journal of Allergy and Clinical Immunology:*

*tional biology*, 16(2):e1007178, 2020.

## APPENDIX A
## VISUALIZATION OF THE PROPOSED IMPACT INDICA-TORS.

To provide readers with a clear understanding of the pro-posed quality metrics, We render the graphical representations of $TNCSI$ and $IEI$ in Fig. 8 and Fig. 9, respectively.

As can be seen in Fig. 8, TNCSI equals the area under the probability density function curve, which is fitted with the use of maximum likelihood estimation. In Fig. 9, the horizontal axis labeled 0 to 5 inversely denotes the months prior to the current month, with 0 representing 6 months ago and 5 denoting the previous month. The corresponding $IEI$ is approximately $-0.76$, which indicates a slightly decreasing citation trend of the investigated review over the last six months.
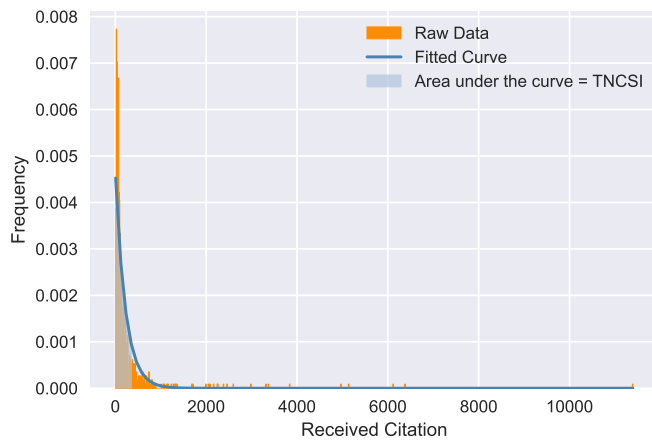


Fig. 8. Visualization of the proposed TNCSI. The TNCSI equals the area under the fitted probability density function curve.
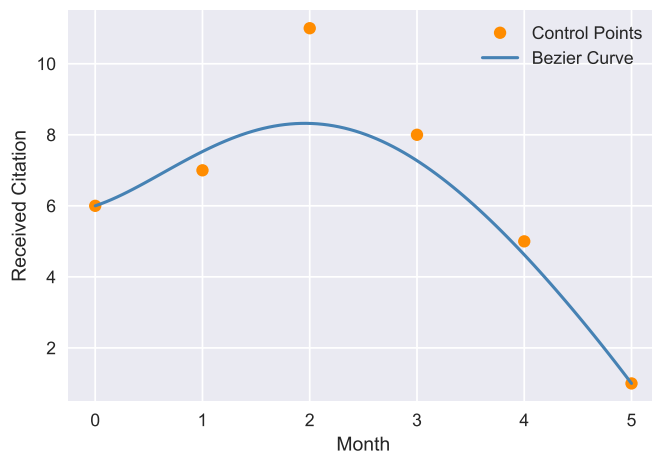


Fig. 9. Visualization of the proposed IEI. IEI is the average of the derivatives of each control point on the Bézier curve.