



Multimodal Lab: Exploring Innovation and Challenges

Ayush Pote, Haridass Bankar

Sharadchandra Pawar College of Engineering University Of Pune

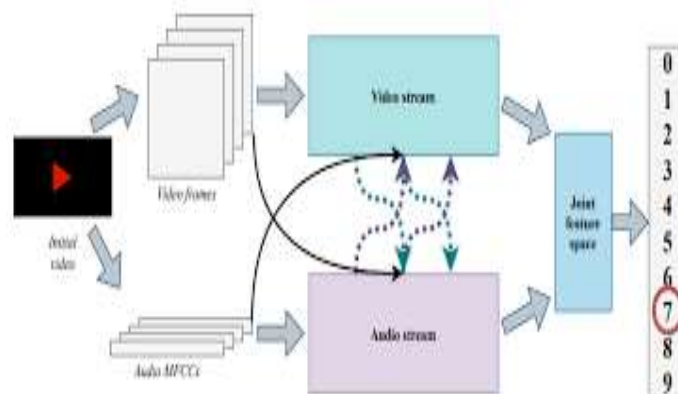
ABSTRACT

The Multimodal machine learning (MML) is a tempting multidisciplinary research area in which heterogeneous facts from a couple of modalities and system studying (ML) are blended to solve critical issues. commonly, studies works use statistics from a unmarried modality, along with photographs, audio, text, and indicators. however, actual-world troubles have become critical now, and managing them using a couple of modalities of records in place of a single modality can considerably effect finding solutions. ML algorithms play an critical role in tuning parameters in growing MML fashions. This paper critiques latest advancements inside the challenges of MML, namely: representation, translation, alignment, fusion and co-getting to know, and affords the gaps and challenges. a systematic literature evaluation (SLR) become applied to outline the development and trends on the ones challenges inside the MML domain. In overall, 1032 articles have been tested on this overview to extract capabilities like source, domain, utility, modality, and so forth. This studies article will help researchers recognize the consistent state of MML and navigate the choice of destiny research directions.

Index Terms-Multimodal machine learning, systematic literature review, representation, translation, alignment, fusion, co-learning.

I. INTRODUCTION

synthetic intelligence (AI) has progressed unexpectedly in the last few a long time. It affects human livelihood, fitness care, technology and era. With the glide of progress, AI development in its techniques to address needs more important realworld issues. ML is an utility of AI that gives a system the capability to learn and improve from its revel in mechanically. The present day fashion of ML is high and entails numerous issues to offer solutions. it's miles wealthy in algorithms and makes use of them to construct fashions that could manner one of a kind styles of records. information are ubiquitous and keep records inclusive of legit reports, clinical or economic statistics and so forth. The significance of facts is increasing with the progress of AI. It incorporates statistics in numerous bureaucracy, together with numeric, text, indicators and so forth. information can come from a awesome range of modalities wherein modality manner how some thing is experienced or occurs. visible, auditory, haptic, physiological alerts, and so on., are examples of modality. statistics may be defined as multimodal when a couple of modalities are concerned together. Speech popularity is a multimodal instance in which audio and visible statistics are blended to understand what a person is saying . The blend of multimodal statistics and ML frames the belief of MML, which focuses on building fashions to method multimodal data from more than one modalities. facts from various modalities are continually heterogeneous. Heterogeneous information are constantly ambiguous, and gaining knowledge of from these multimodal facts gives an possibility to understand the relationships among modalities . 5 middle technical challenges blanketed MML: representation, translation, alignment, fusion, and co-mastering . figure 1 depicts the classification diagram of MML. representation is the primary project, because of this providing statistics the usage of statistics from modalities. Representing a couple of modalities is vital due to the fact facts come from heterogeneous assets, incorporate noise, and can have missing.



The last challenge is synchronous learning, where one model transfers knowledge to another model. Multimodal collaborative learning comes from a heterogeneous background, is noisy, and may not have a parallel, unequal, and hybrid distribution. In parallel, modalities refer to a group, but in non-parallel, modalities refer to ideas or categories rather than examples. In mixed mode, two different modalities are connected by a shared mode. This SLR investigates five major recent adoption issues to answer the research questions.

The rest of this research is organized as follows. The motivation and contributions of this study are presented in Section 2. Section 3 concludes with the analysis of MML and its problems. Section 4 explains the details of the research methodology. Section 5 presents numerical results. The analysis performed is explained in Section 6. Finally, the article is concluded in Section 7.

II MOTIVATION AND CONTRIBUTION

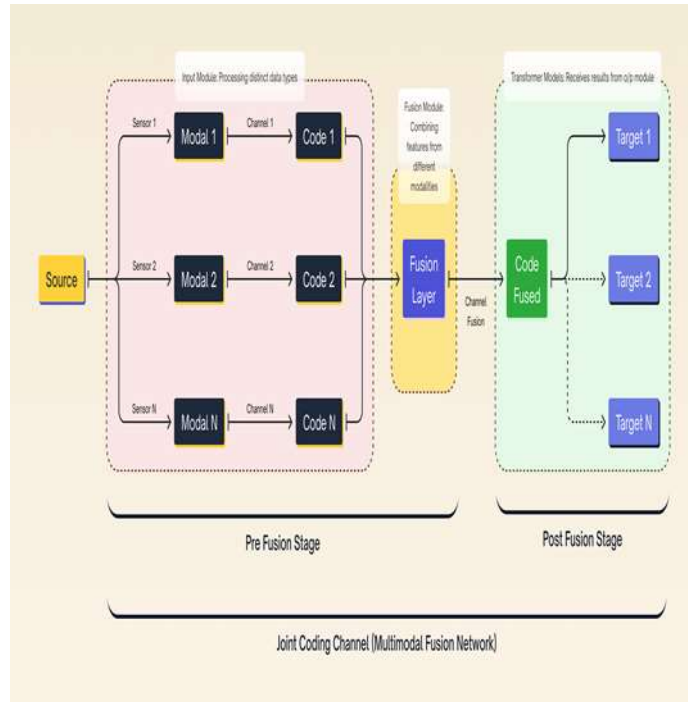
There are many studies on MML and its problems. Most studies have limitations in covering the samples and topics. For example, many studies are specific to one area only. There is also a lack of research with different models, not completely competitive. This limitation supports the conduct of this study. In the related research section, many publications are listed and compared with this research. The main goal is to find information about MML and models and to find insights that can guide future directions. In summary, the contributions of this paper are as follows:

- 374 articles on MML and its issues were identified.
- The results are explained to understand the current research.
- Algorithm for building MML models.

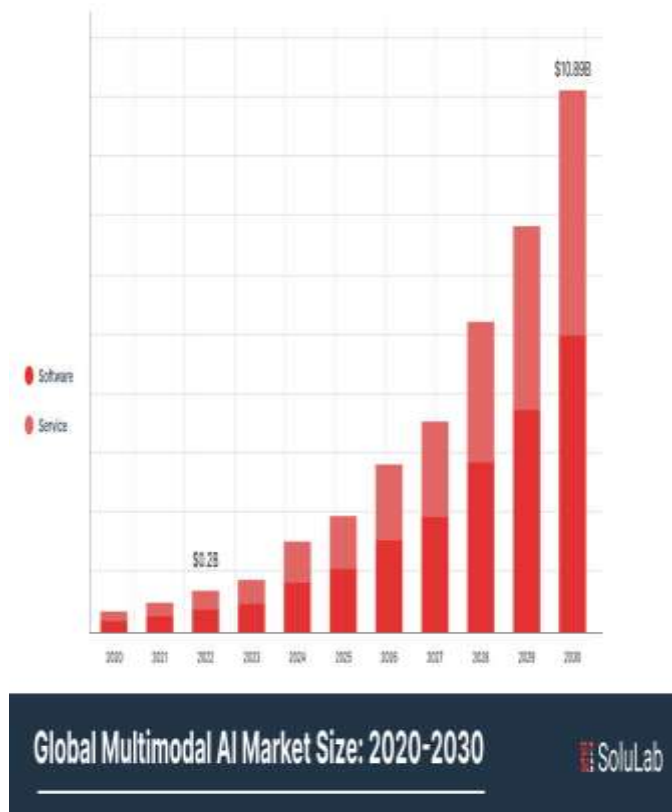
III. RELATED STUDIES

This section presents existing research or review studies on MML and its challenges. We highlight the differences and compare them with our research. MML is an advanced research program and is rapidly developing. A research team at Carnegie Mellon University has done a good analysis of MML by dividing all the problems of MML into different areas. The survey focused on audio, video and text. Another paper reviews methods and applications in deep learning, where the author focuses on some deep learning (DL) methods and applications. In addition to this paper, we also came across several studies discussing how MML can solve different problems with models. In, they investigated meme classification, sentiment analysis and content understanding using MML. Visual and speech analysis is a fascinating research area and new possibilities with MML are rapidly increasing. Review papers show the progress and evolution in using MML technology for computer vision, language and image analysis. Each tournament has a specific role in MML.

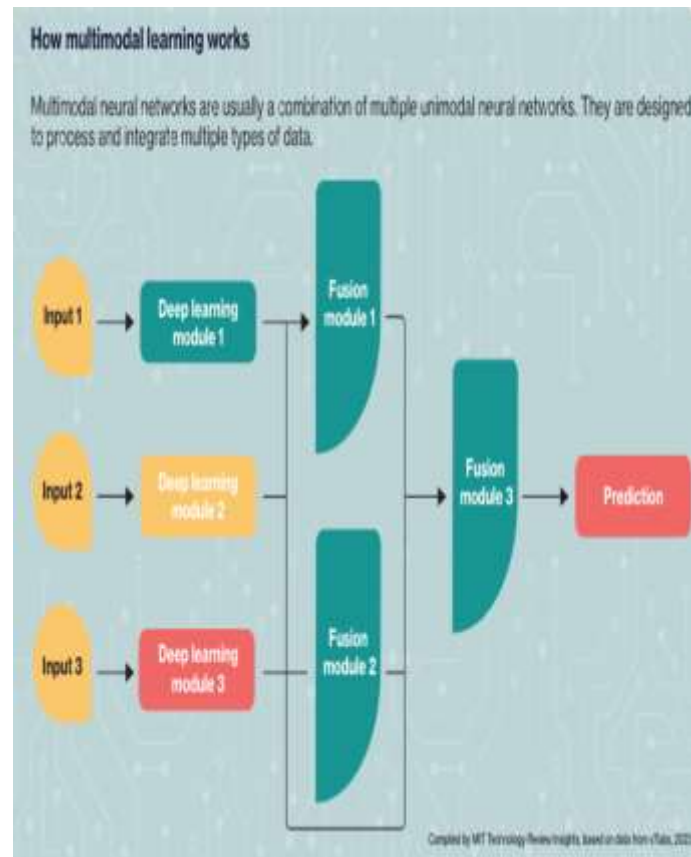
The first challenge of MML is representation, which means representing and aggregating data to show the integration and integration in the model. Two recent research works and provide an overview of the learning representation and planning process. The paper reviews the development of representation learning in unsupervised learning and deep learning. In, two groups introduced multiple representation learning and investigated many simple applications. performed a representation analysis to evaluate its process in human cognition. A pattern. The interpretation method is broad and often focuses on specific models. To date, no research paper has focused entirely on the translation of various technologies and advancements. A recent study revealed the range of multimodal interpretation of visual and verbal information. As in translation, there are not enough research papers focusing on different ways of competition. The paper introduces the manifold alignment kernel method (MAM) that can match the number of sources without similar pairs. In they discuss multiple relational representation and its applications as a class of multiple representations.



Fusion is the fourth and most studied tournament in MML. Connecting data from two or more variables is the main goal of multivariate integration . References are three recent studies that work on multiple composite systems that divide the computational process into different models. Sections summarize the challenges and prospects of the fusion approach. Without focusing on the fusion process, there are some studies on the use of fusion in various fields. Healthcare is a necessary field that can make a significant impact with joint efforts. In they conducted a study on medical fusion to promote smart medical care. Fusion is well known for combining images, especially medical images that have been evaluated in various ways, such as . Operations research and analysis systems include data from multiple sources, and fusion technology helps integrate all these changes. In they presented a research paper on cognitive processing using fusion techniques. References present studies on biometric systems that use fusion techniques to combine multiple data. Fusion is often used for audiovisual language fusion. References are good examples of research on audiovisual data fusion.



Except, above discussed related works mainly focused on one specific challenge and its application. Instead of focusing on a particular challenge, this survey included all and discussed their current advances, gaps and challenges. Although article focused on three modalities, this paper contained all the possible modalities. Article discussed the current use of ML methods and applications in MML, but they limited their review by selecting typical ML methods and applications. On the contrary, this study presented all ML algorithms, domains, and applications available in the search range. To compare included related surveys and this work, an analysis of the associated surveys are presented in Table 1. The table clearly distinguishes between the inclusion of modalities and the challenges of MML in the study. the trend between the published year and the number of citations of included survey papers. From the figure, it is visible that after the year 2018, researchers are interested in working on MML and its challenges. However, the number of citations is lower than in previous years, but it will increase over time.



IV. RESEARCH METHODOLOGY

To achieve the objectives of this study, a qualitative literature review (SLR) was conducted according to the guidelines of the article. Also known as a qualitative review, SLR aims to identify, analyze, and interpret all studies related to the research area or question. three steps Participating in SLR: planning, implementing and reporting reviews. Figure 4 shows each step and its steps. This section explains each step in detail.

A. Review the plan The planning phase is the first phase related to designing and creating the work process. It includes determining the importance of SLR in a particular area, defining the research questions that the SLR will address, and developing a review plan that describes the review process.

1) NECESSITY OF SLR

Before starting, it is necessary to check the importance of such a review. Recently, researchers have focused on using MML techniques to solve multimodal problems. However, there is still no word about complex technologies and the methods used do not always provide solutions. In addition, the use of the model requires rapid payment. An understanding of the difficulties and standards of MML.

2) RESEARCH QUESTIONS (RQs)

Teaching RQ is the most important in SLR. Analyzing the past performance of the competition and understanding MML is the main purpose of this SLR. It includes a definition of MML, approaches to problems, decision models, machine learning models used, and gaps for future research. In order to facilitate this study, the following four questions were asked, including one mini-question.

RQ1: What is the meaning of multimodal machine learning?

RQ2: What are the issues in building multimodal machine learning? Competing decision models in multimodal machine learning?

3) DEVELOPING SLR PROTOCOL

The analytical methods of SLR are used to initiate specific analyses of certain fields. In this study, a strategy has been developed to achieve this goal. First, we look . For preliminary research from well-known domain names. Secondly, we leave the flowers to the selection process. Data extraction and qualitative research analysis were carried out in the third and fourth steps respectively. The second will include data analysis. In this paper, we propose a new data replication strategy to support big data analysis.

B. CONDUCTING THE SLR

It is important to have an SLR that starts as soon as scientists agree on a policy . All the steps listed in the process to achieve the learning objectives are carried out in this section. It is divided into five sections and is discussed below.

1. Research and analysis

In order to find answers to all research questions, we use some keywords to search for basic research in the famous online repository of SLR. We search four major online databases to exclude biased results when searching for articles. A rich journal and conference library is the first reason to explore these four important repositories. Table 2 presents the names of the four repository locations. This process is divided into two parts: determining the search terms that include the MML domain and its five matches, and defining the query by placing the terms of the Boolean operator AND , the keywords, and the introduction. During the search, we also included information from evidence found in previous studies.

2) Examine the selection process

This study followed the PRISMA (Preferred Guidelines for Reviews and Meta-analyses) steps in to identify research articles. The steps of PRISMA include identification, analysis, qualification, and ranking. A total of 1009 articles were initially collected from archives and 23 from other sources using the content. The total number is 800 after removing duplicates. Initially, 593 articles were included in the eligibility section, of which 57 were excluded because they were not published in any journal or conference. 545 articles were then selected because of their research or studies, and 48 articles related to research or analysis were excluded. Of the 545 articles, 338 met the MML requirements and challenges, while 207 articles were removed.

3) Effective learning assessment.

This section evaluates each selected sentence according to the specified query. All questions were prepared according to the instructions in items , the total score is still between 0 points (bad) and 6 points (very good). 4 points or higher. Table 6 provides an example of the evaluation process using five sentences.

4) EXTRACTION OF DATA

The author has made a comprehensive analysis of the collected data from various perspectives. The features of the metadata give an idea about where the article is published. Then the products are divided into groups according to the registration and application status. The article is divided into algorithms, data and models in the study. This algorithm feature shows the different ML algorithms used to solve MML problems. The modal attribute provides a clear view of where the element appears. All categories and studies conducted in this review were done to understand and document current trends in MML.

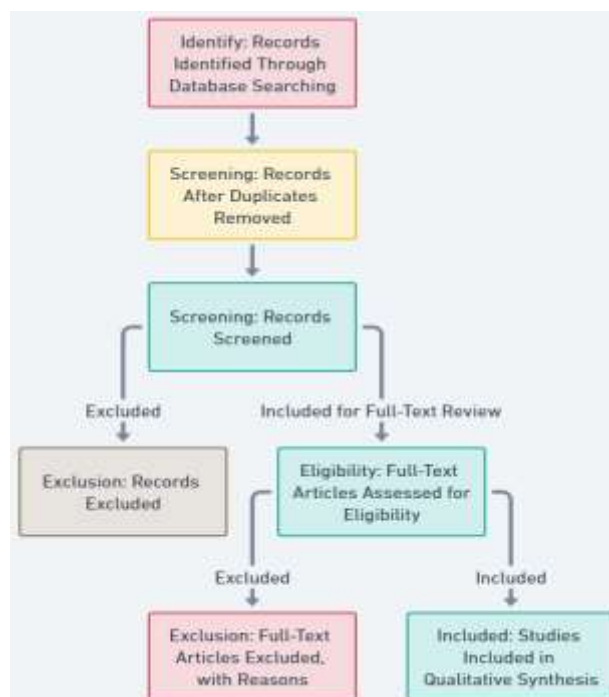


Fig 5. using flow diagram of prisma

#	Study Quality Assessment Question	Yes	Partially	No
Q1	Are the aims and objectives of the study clearly specified?	34 (64.1%)	19 (35.8%)	0 (0.0%)
Q2	Is the context of the study clearly stated?	34 (64.1%)	18 (33.9%)	1 (1.8%)
Q3	Does the research design support the aims of the study?	39 (73.5%)	14 (26.4%)	0 (0.0%)
Q4	Has the study an adequate description of the technique for visualization?	40 (75.4%)	13 (24.5%)	0 (0.0%)
Q5	Is there a clear statement of findings by applying visualization technique in software architecture?	27 (50.9%)	26 (49%)	0 (0.0%)
Q6	Do the researchers critically examine their potential bias and influence to the study?	2 (3.7%)	2 (3.7%)	49 (92.4%)
Q7	Are limitations of the study discussed explicitly?	4 (7.5%)	23 (43.3%)	26 (49%)

Table 5. quality assessment

V. RESULT

A. METADATA

This section presents all the metadata extracted from the text. According to the inclusion and exclusion criteria, all the rumors were found to be related to MML and its issues. Out of the 374 selected articles, 268 were published in journals and 106 were presented at conferences. A total of 28 publications were identified from the selected database. with the most newsletters and forums. Most of the related articles are published in IEEE Xplorer. ELSEVIER is another popular site that publishes similar journals after IEEE Xplorer. This article analyzes the number of articles from different countries. A total of 51 countries were collected, including China, the USA, Germany, and Great Britain. The total number of articles and the number of articles written from 2009 to 2022 . There has been rapid growth since 2018. In the image, all the names are on the right side, with arrows attached to each relevant application. in the domain

B. TASK

This section represents the information about the algorithm and the information used by the selected text. Each article uses a different type of machine learning algorithm to do the job. This research extracted a total of 79 algorithms from the literature. All the collection algorithms are grouped according to their types and methods. The four types of machine learning algorithms are supervised, semi-supervised, unsupervised, and supported. Firstly, the extraction algorithms are divided into supervised, semi-supervised and unsupervised. Data collection shows that the number of supervised, semi-supervised and unsupervised algorithms are 63, 10 and 6, respectively. Then, the extraction process is divided into 11 models: neural network (NN), support vector machine (SVM), mixed model (EM), nearest neighbor model (NNM), tree model (TB), Bayesian model (BM), linear model (LM), genetic algorithm (GA), graph-based model (GBM), encoder-decoder (ED) and k-language. Each model is associated with supervised, semi-supervised and unsupervised. For example, the algorithms related to neural networks can be supervised, semi-supervised and unsupervised, so there is a connection between them. The total number of encountered NN, SVM, EM, NNM, TB, BM, LM, GA, GBM, ED and k-word algorithms are 49, 10, 4, 2, 2, 4, 3, 1, 1, 2 and 1 respectively. Finally, the collected algorithms were divided into 79 unique types, each of which was associated with 11 pre-classified models. The list of 79 algorithms is Artificial neural network (ANN), AlexNet, bidirectional recurrent neural network (BRNN) and convolutional neural network.

Extreme Learning Machine (ELM), Fast Shaped Based Network (FS-Net), Feed Forward Neural Network (FNN), Fully Connected Neural Network (FCNN), FCNNCRF, iMageNet, Inception, LSTM, Liquid State Machine (LSM), MobileNet, multi-cluster network constraint CNN (MGNC-CNN), multi-view CNN short-term memory (MICNN-L), recurrent neural network (RNN), residual neural network (ResNet), volume-based region Convolutional Neural Network (RCNN), RNN-LSTM, Extrusion and excitation based ResNet (SE-ResNet), SE-CNN, TextNet, Text-CNN, Temporal CNN (T-CNN), Traffic Sign Recognition CNN (TSR-CNN), U-Net, Visual Geometry Group (VGG), VolNet, Bidirectional Long-Term Memory (Bi-LSTM), Bidirectional Convolution.

VI. DISCUSSION

Initially, the first research using the topic showed that there was interest in MML research. However, one or two aspects of MML were known by researchers more than a decade ago, but it started to gain interest as an important research topic in late 2016. All the research questions mentioned in Chapter 4 will be discussed and discussed in this chapter. Figure 14 shows the various support methods of this study. RQ1 Analysis

The first research question was determined to examine the meaning of MML. MML represents a design model that can process and describe data in various formats. To more accurately determine the meaning of MML, it is necessary to understand the relationship between multimodality and ML. Multimodality or multimodality is used in many ways in the curriculum. According to, multimodality means containing various modes such as image, text and sound. A recent paper describes MML as a machine learning that takes data from multiple models and processes them. They are further divided into three types: human-centric, machine-centric, and task-centric. Human-centric is about the transmission of knowledge through humans. Machine-centric means that the product is coded using machine learning before it is processed. Task-centric is concerned with the task that the learning machine has to perform, and accordingly, the input and output are represented differently. The human and machine-centered definition means to be more general, independent of the work. The task-centered approach tries to find the contribution of each input according to its relationship with the task. Therefore, based on the above points, the article [411] believes that the purpose of MML is not only to build machine learning models to realize various models, but

also to focus on the relationship between reform and work. As discussed above, two explanations can be made for MML. First, in MML, ML models are built as data; Second, these information come from different types and are related to the work.

B. RQ2 Analysis

Number of publications on five challenges between 2009 and 2022. Five challenges in the analysis conducted using MML in the research. Each MML problem solves a specific task; for example, concatenation is used to combine words from two or more variables, while translation transfers one variable to another to translate words. The first challenge of MML is representation; this is important because it is difficult to represent different objects in a meaningful way. As can be seen in recent papers such as, a good data representation is necessary to support the performance of machine learning models. Paper shows the representation of learning and extracting objects based on it. The use of neural networks to represent vision, sound and information is increasing. In cases where probabilistic graphical models are used, graphical models are used to represent latent random variables, such as Deep Boltzmann Machines (DBM). Network representation uses network models where hidden states are considered to represent data, such as Recurrent Neural Networks (RNN). Among the three subtypes of common representation, neural networks are mainly used due to their efficiency, but they cannot handle incomplete data. However, graphical models can handle incomplete data and complete models. Sequential models are used to represent sequences of data. Similarities and models are two types of collective representation. A similar formula applies to distance.

Combination of the two types. Structural models are often used for cross-mix analysis, such as canonical correlation analysis (CCA). In contrast, joint representation is best for all variables, while joint representation is suitable when only one model is available at the time of testing. In joint representation, more than two models are used, but the joint representation is limited to two models. There is currently a growing research on this topic. Alma is a simple example-based multimodal translation that tries to find the closest model in the dictionary to generate translation results. It is well known for regression models such as CNN and Kernel Canonical Correlation Analysis (KCCA). In contrast, there are compositional models that combine parts of the dictionary to generate a good translation. Grammar models model translations in specific domains using grammar. Encoder and decoder models first encode the original modality into a latent representation and then decode it to generate the target modality, for which CNNs are used in most cases. Continuum models always have results such as translation to interval. Short time-span (LSTM) RNNs often use continuous models such as. Modeling is simple in comparison but makes the model heavy because the model itself acts as a dictionary and sometimes produces incorrect definitions. Generative models are important for construction because the model must be able to generate symbols in sequence. This is why many researchers choose structural models to provide solutions.

VII. Conclusion

This paper has conducted a qualitative literature review on MML and its issues to provide an overview of recent trends. The study was conducted using the PRISMA method and the selection process is reported in detail. A total of 374 articles were selected from the initial collection of 1032 articles based on their relevance to the four research questions. The results of this review show that MML depends not only on the model but also on the ML algorithm and task. A survey of algorithms and data shows that neural network-based algorithms and image data are the most widely used. This review also highlights various aspects of multimodal representation, translation, integration, integration and collaborative learning and their differences. However, this SLR provides an overview of the work done in MML and should be expanded here to provide opportunities for future work for researchers interested in the subject.els to provide solutions.