

# A Survey on Multimodal Recommender Systems: Recent Advances and Future Directions

Jinfeng Xu, Zheyu Chen, Shuo Yang, Jinze Li, Wei Wang,  
Xiping Hu, Steven Hoi, *Fellow, IEEE* and Edith Ngai, *Senior Member, IEEE*

**Abstract**—Acquiring valuable data from the rapidly expanding information on the internet has become a significant concern, and recommender systems have emerged as a widely used and effective tool for helping users discover items of interest. The essence of recommender systems lies in their ability to predict users' ratings or preferences for various items and subsequently recommend the most relevant ones based on historical interaction data and publicly available information. With the advent of diverse multimedia services, including text, images, video, and audio, humans can perceive the world through multiple modalities. Consequently, a recommender system capable of understanding and interpreting different modal data can more effectively refer to individual preferences. Multimodal Recommender Systems (MRS) not only capture implicit interaction information across multiple modalities but also have the potential to uncover hidden relationships between these modalities. The primary objective of this survey is to comprehensively review recent research advancements in MRS and to analyze the models from a technical perspective. Specifically, we aim to summarize the general process and main challenges of MRS from a technical perspective. We then introduce the existing MRS models by categorizing them into four key areas: Feature Extraction, Encoder, Multimodal Fusion, and Loss Function. Finally, we further discuss potential future directions for developing and enhancing MRS. This survey serves as a comprehensive guide for researchers and practitioners in MRS field, providing insights into the current state of MRS technology and identifying areas for future research. We hope to contribute to developing a more sophisticated and effective multimodal recommender system. To access more details of this paper, we open source a repository: <https://github.com/Jinfeng-Xu/Awesome-Multimodal-Recommender-Systems>.

**Index Terms**—Information systems, Data mining, Multimedia information systems, Multimodal recommender systems.

## I. INTRODUCTION

THE rapid expansion of the Internet has resulted in an overwhelming abundance of information, making it increasingly challenging for users to identify what is useful

and relevant. This phenomenon, referred to as information overload, arises due to the near impossibility of controlling the generation and dissemination of information in the digital age. Consequently, there is an urgent need for robust filtering mechanisms that prioritize pertinent content to facilitate efficient communication and decision-making processes. Recommender systems, which personalize content filters according to specific requirements across various domains, have demonstrated their efficacy in mitigating the adverse effects of information overload. These systems have proven particularly successful in commercial applications such as e-commerce, advertising, and social media, where personalization is crucial to user engagement and satisfaction [1]–[4].

The primary function of recommender systems is to predict users' ratings or preferences for various items and recommend the most likely and relevant items based on historical interaction data and publicly available information. However, traditional ID-based recommendation methods, which operate on the principle that users tend to select items akin to those they have previously liked, often strongly depend on enough user-item interactions. Despite their successes, recommender systems face two significant challenges: data sparsity and the cold start problem. Data sparsity arises from the natural sparsity of interaction data between users and products, making it difficult to accurately predict users' preferences. This sparsity can lead to unreliable recommendations, especially in systems with large item catalogs but relatively few user interactions. The cold start problem occurs because traditional recommender system models rely heavily on ID embeddings, which struggle to make satisfactory predictions for new users or products with little to no historical interaction data. This challenge is particularly pronounced in dynamic environments where new items and users are continuously introduced.

To alleviate these issues, multimodal information is increasingly being integrated into recommendation systems. MRS leverages auxiliary multimodal information, such as text, images, videos, and audio, to complement the historical interactions between users and items. This approach enhances recommendation performance by providing a richer and more comprehensive understanding of user preferences. The essential goal of recommender systems is to cater to people's preferences, and since human perception of the world is inherently multimodal, incorporating diverse modal information can capture preferences at a finer level of granularity. This leads to more accurate and personalized recommendations, thereby improving user satisfaction and engagement.

Research in multimodal recommendation is rapidly growing

Jinfeng Xu, Jinze Li, Shuo Yang, and Edith C. H. Ngai\* are with the Department of Electrical and Electronic Engineering, The University of Hong Kong, Hong Kong, China (e-mail: jinfeng, lijnize-hku, shuo.yang, @connect.hku.hk, chngai@eee.hku.hk).

Zheyu Chen is with the Department of Electrical and Electronic Engineering, The Hong Kong Polytechnic University, Hong Kong, China (e-mail: zheyu.chen@connect.polyu.hk).

Wei Wang and Xiping Hu are with the Department of Engineering, Shenzhen MSU-BIT University, Shenzhen, China, and also with the School of Medical Technology, Beijing Institute of Technology, Beijing, China (e-mail: ehomewang@ieee.org, huxp@bit.edu.cn).

Steven Hoi is with the School of Computing and Information Systems, Singapore Management University, Singapore (e-mail: chhoi@smu.edu.sg).

\*The corresponding author.

This paper was produced by the IEEE Publication Technology Group. They are in Piscataway, NJ.

Manuscript received XXXX; revised XXXX.

TABLE I  
RELATED RECOMMENDER SYSTEM SURVEY

Surveys	Key Contributions	Differences from Our Survey
Zhang et al. [6]	provide a panorama for advances in deep learning based recommender systems and provide a survey of future direction and challenges.	This work comprehensively demonstrates the advanced development in deep learning based recommender systems, which include MRS but lacks the fine-grained introduction to existing state-of-the-art technologies.
Guo et al. [7]	provide a fine-grained survey for the existing approaches utilizing the KG to improve the recommendation result and introduce some datasets used in different scenarios.	Our survey focuses on the taxonomy of the process for MRS and state-of-the-art technologies of the multimodal recommender system; the KG only discussed as a part of the techniques in our work.
Deldjoo et al. [8]	provide a comprehensive and coarse-grained survey and a coarse-grained categorization by the modalities.	The categorization in this work is coarse-grained and unreasonable to some extent, while our work provides a fine-grained categorization for techniques.
Jannach et al. [9]	explore the field of CRS and provide a taxonomic survey of existing techniques.	This work discusses the recent approaches in the CRS field but lacks a combination with multimodal information.
Deldjoo et al. [1]	discuss the state-of-the-art approaches to content-driven MRS and provide a survey of challenges and historical evolution.	This work focuses on the content-driven MRS rather than covering all perspectives of MRS.
Wu et al. [10]	provide a comprehensive survey for utilizing GNN techniques in the RS field and list several limitations and future directions.	Our survey focuses on the MRS field and provides a more fine-grained classification of GNN techniques in the MRS field.
Deldjoo et al. [2]	provide a comprehensive survey of RS in the fashion field according to the task in the market, and provide some vital evaluation goals in the fashion field.	This work focuses on the RS in the fashion field but lacks a general and fine-grained survey for RS.
Meng et al. [11]	provide a comprehensive analysis for personalized news recommendations via technologies and list several limitations and future directions.	Our survey focuses on the taxonomy of the process for general MRS and state-of-the-art technologies, which are also effective in the news field.
Zhou et al. [12]	summarize the main methods used in MRS and provide a common framework for commonly used MRS models.	This work comprehensively demonstrates the previous approaches in MRS but lacks a fine-grained process of MRS.
Liu et al. [5]	summarize the main methods used in MRS and provide a common framework for commonly used MRS models.	This work delineates the MRS approach from a process perspective, rather than a technology development perspective, and does not allow the reader to fully understand the direction of research in the field.

and evolving. To assist researchers in quickly understanding MRS and to support community development, a comprehensive review from a technical perspective is urgently needed. Existing work [5] attempts to categorize MRS from a technical standpoint; however, the rapid advancement of the field has rendered some of its categorizations outdated. Therefore, we aim to collect recent work and propose a more up-to-date categorization framework to help researchers grasp the latest progress in the MRS community. This review will provide a thorough overview of current MRS technologies, highlight emerging trends, and identify potential future directions for research and development in this dynamic field. By systematically examining the state-of-the-art (SOTA) works in MRS, we hope to contribute to the ongoing efforts to enhance the capabilities and applications of recommender systems in a multimodal digital world.

#### A. Search Strategy for Relevant Papers

We conducted a comprehensive survey on Multimodal Recommendation Systems (MRS) by systematically retrieving and analyzing articles from leading conferences and journals in the field. The conferences and journals included, but were not limited to, MM, KDD, WWW, SIGIR, AAAI, ICLR, IJCAI, CIKM, WSDM, TMM, TKDE, TPAMI, and INFFUS. This rigorous selection process ensured that our survey covered the most influential and cutting-edge research in MRS.

Our search approach was methodically divided into three distinct stages:

- **Collection of High-Quality Articles:** In the initial stage, we gathered articles from the aforementioned top conferences and journals. This selection was based on the reputation and impact factor of the sources, ensuring that only high-quality and peer-reviewed research was included in our survey.
- **Filtering and Post-Processing:** Following the collection phase, we meticulously filtered and post-processed the articles. This step involved removing duplicates, assessing the relevance of each article to the topic of MRS, and ensuring that only the most pertinent studies were retained. This rigorous filtering process was crucial for maintaining the focus and quality of our survey.
- **Technical Analysis and Synthesis:** In the final stage, we conducted a detailed analysis of the techniques employed in each article. This involved examining the methodologies, models, and algorithms used, as well as the motivations behind these approaches. We also reviewed related works cited within each article to provide a comprehensive understanding of the evolution and current trends in MRS. By synthesizing this information, we were able to summarize the key techniques and motivations driving the field.

Through this systematic approach, our survey offers a thorough and insightful overview of the SOTA works in Multimodal Recommendation Systems. It highlights the significant advancements, emerging trends, and potential future directions in the field, providing valuable guidance for researchers and

practitioners alike.

### B. Compared with Related Surveys

Several surveys have been published on recommender systems, addressing either general aspects or specific facets of these systems. However, none provide a comprehensive and reasonable taxonomy of the processes and detailed technologies utilized in recent SOTA MRS works, which is an emerging and crucial requirement in this field. The objective of MRS is to enhance the capability of extracting deeper and more accurate interactions between users and items by incorporating multimodal information into recommender systems. This paper discusses the main contributions and limitations of existing related surveys and highlights the unique contributions of our work, as summarized in Table I.

Zhang et al. [6] offer a panoramic view of advances in deep learning-based recommender systems, surveying future directions and challenges, including joint representation learning, explainability, deeper models, and machine reasoning. However, their work lacks a fine-grained introduction to existing SOTA technologies. Deldjoo et al. [8] provide a comprehensive survey and a coarse-grained categorization by modalities, including common features such as audio, visual, and textual, as well as special features like motion, metadata, and semantic orientation. Nonetheless, this categorization is somewhat coarse-grained and lacks precision.

Jannach et al. [9] explore the field of conversational recommender systems (CRS) and offer a taxonomic survey of existing techniques, but their work does not integrate multimodal information. Deldjoo et al. [1] discuss SOTA approaches to content-driven MRS, surveying challenges and historical evolution, including increasing recommendation diversity and novelty, providing transparency and explanations, achieving context-awareness, improving scalability and efficiency, and alleviating the cold start problem. However, their focus is primarily on content-driven MRS, rather than covering the general MRS landscape.

Previous works [7], [10] focus on graph structure in recommendation systems. Guo et al. [7] provide a fine-grained survey of approaches utilizing knowledge graphs (KG) to enhance recommendation results, categorizing methods into embedding-based, path-based, and unified approaches. Wu et al. [10] offer a comprehensive survey of graph neural network (GNN) techniques in recommender systems, identifying several limitations and future directions, including diverse and uncertain representation, scalability, dynamics, receptive fields, self-supervised learning, robustness, privacy-preserving methods, and fairness.

Deldjoo et al. [2] provide a comprehensive survey of recommender systems in the fashion domain, categorizing tasks in the market and outlining vital evaluation goals specific to fashion. Meng et al. [11] present a thorough analysis of personalized news recommendations, discussing technologies and listing several limitations and future directions, including privacy protection, fake news mitigation, and de-biasing.

Zhou et al. [12] summarize the main methods employed in MRS and propose a common framework for commonly

used MRS models. While their work offers a comprehensive overview of previous approaches in MRS, the pipeline proposed for MRS requires more detailed elaboration. More recently, Liu et al. [5] also summarize the main methods used in MRS and provide a common framework for MRS models. However, this work delineates MRS from a process perspective rather than focusing on technological developments, limiting the reader's ability to fully understand the research directions in the field.

In conclusion, our work aims to fill these gaps by providing a more detailed and up-to-date taxonomy of MRS processes and technologies, thereby advancing the understanding and development of this rapidly evolving domain.

Our survey focuses on a refined categorization of MRS from a technological perspective to provide researchers with insight into the technological development of MRS. Finally, we discuss potential future directions for developing and improving multimodal recommendations.

### C. The Outline of the Survey

The structure of our survey is organized according to the following:

- **Section I: Introduction**

We briefly outline the historical development of RS and underscore the significance of leveraging multimodal information to enhance recommendations. Subsequently, we detail the search strategy employed to ensure the quality of our work. Additionally, we offer a comparative analysis with previous surveys. Finally, we present the structure of this survey and highlight the main contributions of our research.

- **Section II: Technological Taxonomy**

We present a technological taxonomy for MRS summarize all recent SOTA works, and then overview the four key technologies respectively.

- **Section III: Feature Extraction**

We have comprehensively reviewed recent MRS works in feature extraction techniques for visual and textual modalities, identifying prevailing trends. This analysis has culminated in a curated standard setting, sparing researchers the complexity of navigating diverse methods.

- **Section IV: Encoder**

We categorize the MRS encoder types from a technical perspective and summarize all recent SOTA works. This comprehensive overview provides insights into potential research gaps and future directions, facilitating ongoing development and refinement in MRS technologies.

- **Section V: Multimodal Fusion**

We provide a comprehensive categorization for multimodal fusion in MRS, from both timing and strategy perspectives. We also summarize all recent works from these perspectives and reveal the interplays between these two perspectives. This comprehensive categorization is intended to help researchers select appropriate timing and strategies for modal fusion.

- **Section VI: Loss Function**

We provide a detailed introduction to loss functions

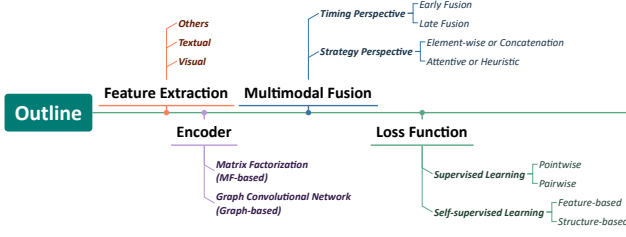


Fig. 1. The outline structure of our survey.

utilized in MRS, including both supervised and self-supervised learning frameworks. Specifically, we introduce and analyze several widely used loss functions, detailing their applications and effectiveness across different learning scenarios within MRS. This exploration aims to assist researchers in selecting and implementing the most suitable loss functions.

#### • Section VII: Future Direction

We delve into the potential future directions of the MRS field. We aim to stimulate and encourage further research, development, and innovation in this rapidly evolving area of research.

#### • Section VIII: Conclusion

We briefly summarize the contents and contributions of this survey.

Figure 1 illustrates a detailed outline structure of our survey (Including Section III - Section VI).

The main contributions of our survey are as follows:

- We provide a comprehensive review of MRS, summarize the huge number of SOTA works, and structure a general process of MRS to present how multimodal information is utilized in the RS field.
- We categorize and analyze the techniques and motivations for each main step of this general process of MRS, which can extremely guide the researcher to conduct further research.
- We introduce the commonly available datasets for MRS and provide a detailed characterization of them. And organized the datasets used in SOTA works to help researchers choose the suitable dataset.
- We discuss the existing challenges for MRS based on previous works and list some future directions, which is worthy of in-depth research.

Table II lists the abbreviations used throughout this paper.

## II. TECHNOLOGICAL TAXONOMY

Based on the current MRS work at SOTA to summarize and organize, we classify the technologies in MRS into four parts as shown in Figure 2. Specifically, there are four parts: **Feature Extraction**, **Encoder**, **Multimodal Fusion**, and **Loss Function**. We briefly overview these parts and discuss them in detail in the subsequent sections.

### A. Feature Extraction

Different application scenarios encompass varying types of modality information, leading to diverse datasets with distinct

TABLE II  
LIST OF ABBREVIATIONS USED THROUGHOUT THIS PAPER

Abbreviation	Term
RS	Recommender Systems
MRS	Multimodal Recommender Systems
CRS	Conversational Recommender Systems
CF	Collaborative Filtering
MF	Matrix Factorization
GCN	Graph Convolution Network
GNN	Graph Neural Network
KG	Knowledge Graphs
SOTA	State-of-the-art
CL	Contrastive Learning
SSL	Self-supervised Learning

multimodal features. Most datasets, however, include at least three primary modalities: interaction, visual, and textual. For instance, large platforms such as Amazon, Netflix, and TikTok provide datasets rich in image and textual information, thus covering both visual and textual modalities. Specifically, TikTok datasets often include additional modalities, such as audio and video [13]–[15]. Furthermore, datasets from specialized domains sometimes have rare modalities. For example, datasets in popular areas like fashion and healthcare often include a variety of specialized modalities.

Feature extraction is a critical process aimed at the representation of low-dimensional, interpretable channel features using embedding techniques. Different pre-extraction methods are employed for distinct modalities. For the visual modality, models such as ResNet [16] and ViT [17] are utilized to extract features. In the case of the textual modality, models like BERT [18] and Sentence-Transformer [19] are used to derive features. Audio features are typically extracted using models such as LSTM [20] and GRU [21].

A detailed introduction to feature extraction is provided in Section III, where we delve into the specifics of each modality and the corresponding extraction techniques.

### B. Encoder

The encoder utilizes features extracted from multimodal information and historical interaction data to infer user preference representations, which are subsequently used in predicting user-item interactions for making recommendations. Similar to traditional recommender systems, encoders for multimodal recommendation can be broadly classified into **Matrix Factorization (MF [22])-based** and **Graph Convolutional Network (Graph [23])-based** approaches. The MF-based approach is known for its simplicity and effectiveness, whereas the Graph-based approach leverages the bipartite graph inherent in user-item interactions to learn higher-order neighbor features.

With the rapid advancement of MRS, more sophisticated encoders have been proposed and utilized to fully exploit the rich multimodal information, thereby enhancing recommendation performance. These advanced encoders enable the integration of diverse multimodal data, leading to more accurate and personalized recommendations.

In Section III, we will provide a detailed introduction to the development and motivations behind both types of encoders.

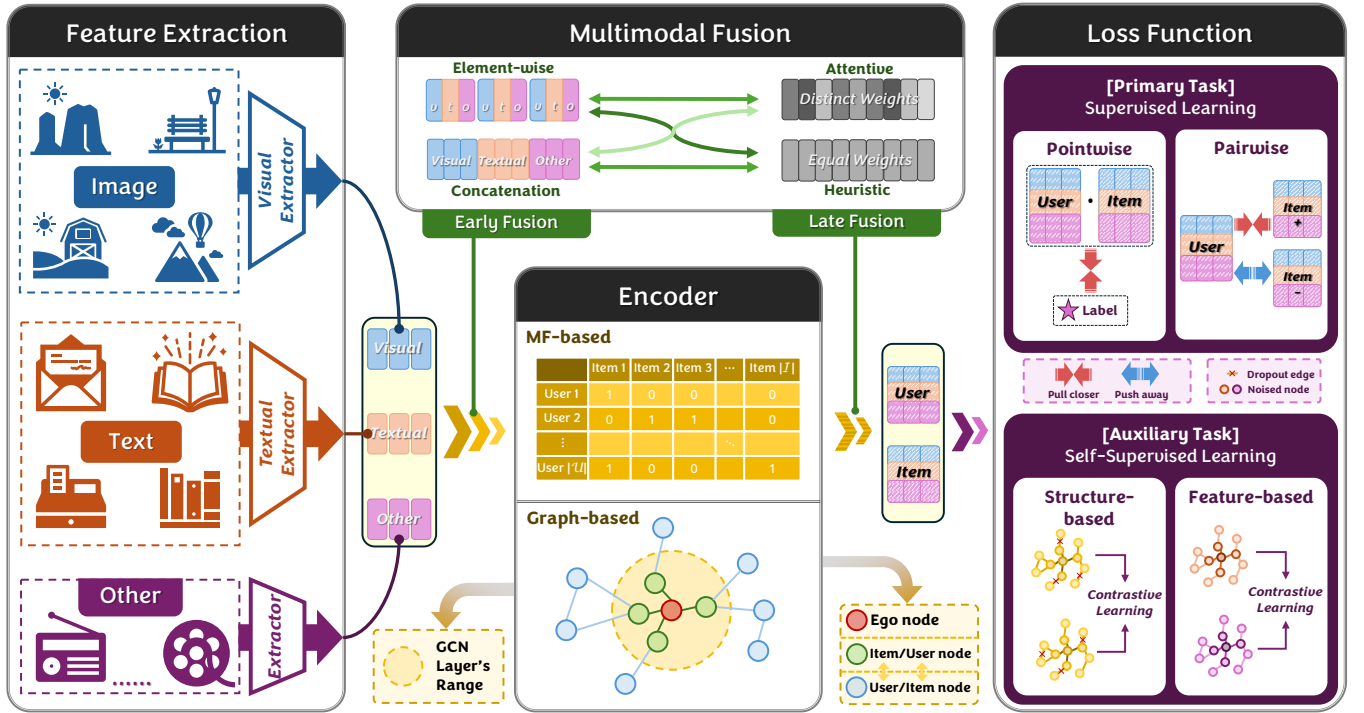


Fig. 2. Technologies in MRS from process pipeline perspectives.

This includes an exploration of how MF-based methods efficiently capture user-item interactions and how Graph-based methods extend this capability by incorporating complex graph structures. By examining these methodologies, we aim to elucidate the strengths and limitations of each approach, as well as their contributions to the MRS community.

### C. Multimodal Fusion

One of the key research focuses in MRS is the Multimodal Fusion. Recent studies have demonstrated that the timing of modality fusion can significantly impact the effectiveness of the recommendations. Multimodal fusion involves integrating information from different modalities at various stages, and this timing can be crucial for achieving optimal performance.

**Early Fusion:** Early fusion involves combining different modality features before they are processed by the encoder. This approach can effectively uncover hidden relationships between modalities, as the integrated multimodal features allow the encoder to learn richer and higher-quality representations. Early fusion can capture the intricate interactions between different types of data, such as text, images, and audio, leading to a more holistic understanding of user preferences. Techniques for early fusion often include concatenation, attention mechanisms, and neural network-based integration methods, which aim to create a unified representation of multimodal data.

**Late Fusion:** Late fusion combines the scores or predictions from each modality after the individual modality-specific encoders have processed them. This approach focuses on leveraging the strengths of each modality-specific model and then combining their outputs to make the final

recommendation. Late fusion can be particularly effective in scenarios where certain modalities are more informative or reliable than others. By deferring the fusion process until after the prediction phase, late fusion allows for more targeted and refined extraction of specific modality information, enhancing the overall recommendation accuracy.

In Section V, we will provide a detailed classification of existing work based on the timing of fusion, categorizing them into early fusion and late fusion approaches. This classification will offer a comprehensive understanding of how different fusion strategies impact the performance of MRS systems. We will explore various methodologies and techniques employed in both early and late fusion, analyzing their advantages, limitations, and application scenarios.

### D. Loss Function

MRS leverages loss functions that can be broadly divided into two components: primary tasks and auxiliary tasks. The primary tasks are supervised learning, which typically involves clearly defined labels to guide the model's learning process. These tasks ensure that the model learns to make accurate predictions based on labeled data. The auxiliary tasks are self-supervised learning (SSL) [24]. SSL generates supervision signals from the inherent structure or patterns within the data itself, rather than relying solely on externally labeled data. This approach allows recommender systems to utilize unlabeled data effectively, extracting meaningful representations and making accurate predictions even in data sparsity scenarios.

Supervised Learning can be further subdivided into Pointwise Loss and Pairwise Loss:

**Pointwise Loss:** This loss is calculated by comparing the predicted score for each individual item with its actual label. Common pointwise loss functions include Mean Squared Error (MSE) [25] and Cross-Entropy Loss (CE) [26], which are used to directly assess the accuracy of individual predictions.

**Pairwise Loss:** This loss focuses on the relative ranking of items. It evaluates the model's ability to correctly order each pair of items based on user preferences. Common pairwise loss functions include Bayesian Personalized Ranking (BPR) [27] and Hinge Loss [28], which aim to optimize the rank order of items rather than their absolute scores.

Self-supervised Learning can be categorized into feature-based and structure-based methods:

**Feature-based SSL:** This method involves creating auxiliary tasks that predict or reconstruct certain features of the data. For example, a model might be trained to predict missing features of an item or user based on the available data, thereby learning more robust representations.

**Structure-based SSL:** This approach leverages the structural properties of the data, such as the relationships and interactions between users and items. Graph-based methods, for instance, might use node similarity or subgraph patterns to generate supervision signals, enhancing the model's ability to capture complex dependencies and interactions.

In Section VI, we will provide a detailed introduction to these loss functions. We will explore the motivations behind each type of loss, their implementation details, and their impact on the performance of multimodal recommender systems. By examining both supervised and self-supervised learning strategies, we aim to offer a comprehensive understanding of how different loss functions contribute to the effectiveness of multimodal recommendation.

### III. FEATURE EXTRACTION

We have summarized feature extraction in visual and textual modalities in Table III for advanced MRS methods in recent years. Within the visual domain, early investigations predominantly utilized convolutional architectures such as CNNs, along with specific models like VGG [77], Inception [78], Caffe [79], and ResNet [16], which have demonstrated remarkable efficacy in handling various visual recognition tasks. These models are prized for their deep learning capabilities, which allow for the extraction of high-level, complex features from raw images. Moreover, in the textual domain, a diverse array of techniques has been employed for feature extraction. These include traditional methods like TF-IDF [80], as well as more sophisticated neural network approaches such as GRU [21], PV-DM (PV-DBOW) [81], and Glove [82]. The introduction of attention mechanisms, along with models like BERT [18], Word2Vec [83], Sentence-Transformer [19], and Sentence2Vec [84], has further revolutionized the ability to understand and process language by allowing for contextually enriched text representations.

As the field progresses, the MMRec<sup>1</sup> open-source framework exemplifies this trend by standardizing feature extraction methods for both visual and textual modalities, thereby

facilitating a more controlled and reproducible experimental settings. Furthermore, data quality poses challenges, such as some corrupted and missing images in the Amazon datasets. Recent approaches have thus shifted towards utilizing pre-provided features in the visual modality to circumvent the issues related to manual processing. Similarly, in the textual modality, reliance on pre-trained models like BERT [57], [58], [64], [65], [67], [73] or Sentence-Transformer [15], [53], [55], [56], [59], [60], [62], [63], [66]–[68], [70]–[72], [74], [75] has become commonplace. These models provide a robust foundation for feature extraction, leveraging vast pre-existing knowledge bases to enhance the accuracy and depth of textual analysis.

Besides, not all studies have engaged both visual and textual modalities comprehensively. Several works have predominantly focused on the visual aspect as a form of auxiliary information. For instance, VBPR [29] and DVBPR [33] underscore the significance of visual features in enhancing recommendation systems, while neglecting textual data. Similarly, VMCF [30] and ACF [31] incorporate visual information to refine the accuracy of their models, yet they do not integrate textual insights which could potentially enrich the contextual understanding of the data. On the adversarial front, AMR [38] leverages visual modality to bolster the robustness of their models against adversarial attacks, yet the textual modality remains unexplored. Conversely, ADDVAE [48] focuses exclusively on textual modality, thereby providing a nuanced understanding of textual data but omitting the rich, descriptive power of visual information.

### IV. ENCODER

In recommender systems, an encoder is typically used to extract a feature representation of a user or item. Encoders can be many types of models, from simple linear models to complex deep neural networks. The main purpose is to convert raw data (e.g., user behavioral data, item attributes, etc.) into a fixed-size embedding that captures the core features of the input data. However, in multimodal recommendation, the encoder plays the same role, but due to the greater variety of features, more diverse and specially designed encoders have been proposed to better utilize the multimodal information and user interaction data for more accurate recommendations. From a technical point of view, we roughly categorize all the encoders into MF-based encoders and Graph-based encoders. We depict these two types of encoders in Figure 3.

**Preliminary 1:** Let a set of users  $\mathcal{U} \in \mathbb{R}^{|\mathcal{U}| \times d}$  and a set of items  $\mathcal{I} \in \mathbb{R}^{|\mathcal{I}| \times d}$ , where  $d$  is the hidden dimension. We use  $\mathbf{R} = [r_{u,i}]^{|\mathcal{U}| \times |\mathcal{I}|}$  to denote the user-item interaction matrix. For Graph-based encoder  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  be a given graph with node set  $\mathcal{V}$  and edge set  $\mathcal{E}$ , where  $|\mathcal{V}| = |\mathcal{U}| + |\mathcal{I}|$ .  $\mathbf{U}$  and  $\mathbf{I}$  denote the hidden embeddings for users and items, respectively.  $\mathbf{E} = \{\mathbf{U}|\mathbf{I}\}$  denotes the embedding for  $\mathcal{V}$ .

For the MRS scenario, we need to restate a more comprehensive preliminary.

**Preliminary 2:** Restate item sets  $\mathcal{I}_m \in \mathbb{R}^{|\mathcal{I}| \times d_m}$ , where  $m \in \mathcal{M}$ ,  $\mathcal{M}$  is the set of modalities, and  $d_m$  is hidden dimension for modality  $m$ . Besides,  $\mathcal{G}_m = (\mathcal{V}_m, \mathcal{E})$  be a given

<sup>1</sup>MMRec: <https://github.com/enoeche/MMRec.git>

TABLE III  
SUMMARY OF MRS FEATURE EXTRACTION FOR BOTH VISUAL AND TEXTUAL MODALITIES.

Modality			Visual					Textual										
Method	Year	Publication	Provided	VGG	Inception	Caffe	ResNet	Provided	TF-IDF	GRU	PV-DM	PV-DBOW	Glove	Attention	BERT	Word2Vec	Sentence-transformer	Sentence 2Vec
(V)VBPR [29]	2016	AAAI	✓															
(V)VMCF [30]	2017	WWW	✓															
(V)ACF [31]	2017	SIGIR					✓											
JRL [32]	2017	CIKM	✓									✓						
(V)DVBPR [33]	2017	ICDM					✓											
GraphCAR [34]	2018	SIGIR	✓					✓										
VECF [35]	2019	SIGIR		✓						✓								
UVCAN [36]	2019	WWW			✓													
MAML [37]	2019	MM				✓					✓							
MMGCN [13]	2019	MM	✓					✓										
(V)AMR [38]	2019	TKDE					✓											
MGAT [39]	2020	I&M	✓					✓										
GRCN [40]	2020	MM	✓					✓										
MKGAT [41]	2020	CIKM					✓										✓	
IMRec [42]	2021	MM	✓											✓				
PMGT [43]	2021	MM			✓										✓			
LATTICE [15]	2021	MM	✓														✓	
HHFAN [44]	2021	TMM	✓					✓										
MVGAE [45]	2021	TMM					✓											✓
DualGNN [14]	2021	TMM	✓					✓										
PAMD [46]	2022	WWW		✓									✓					
MMGCL [47]	2022	SIGIR	✓					✓										
(T)ADDDVAE [48]	2022	KDD							✓									
EliMRec [49]	2022	MM					✓											✓
EgoGCN [50]	2022	MM	✓															✓
InvRL [51]	2022	MM					✓											✓
A2BM2GL [52]	2022	MM	✓					✓										
HCGCN [53]	2022	MM	✓														✓	
SLMRec [54]	2022	TMM	✓					✓										
DMRL [55]	2022	TMM					✓									✓		
BM3 [56]	2023	WWW	✓														✓	
MMSSL [57]	2023	WWW	✓												✓			
BCCL [58]	2023	MM	✓												✓			
FREEDOM [59]	2023	MM	✓														✓	
MGCN [60]	2023	MM	✓														✓	
PaInvRL [61]	2023	MM					✓											✓
DRAGON [62]	2023	ECAI	✓														✓	
LGMRec [63]	2024	AAAI	✓														✓	
LLMRec [64]	2024	WSDM	✓												✓			
PromptMM [65]	2024	WWW	✓												✓			
MCDRec [66]	2024	WWW	✓														✓	
DiffMM [67]	2024	MM	✓												✓			
SOIL [68]	2024	MM	✓														✓	
CKD [69]	2024	MM	✓														✓	
GUME [70]	2024	CIKM	✓														✓	
POWERec [71]	2024	INFFUS	✓														✓	
DGVAE [72]	2024	TMM	✓														✓	
VMoSE [73]	2024	TMM	✓												✓			
SAND [74]	2024	arXiv	✓														✓	
MENTOR [75]	2025	AAAI	✓														✓	
SMORE [76]	2025	WSDM	✓														✓	

graph with node set  $\mathcal{V}_m$  and edge set  $\mathcal{E}$ , where  $|\mathcal{V}_m| = |\mathcal{U}| + |\mathcal{I}_m|$ .  $\mathbf{I}_m$  denotes the items hidden embeddings for modality  $m$ .  $\mathbf{E}_m = \{\mathbf{U}|\mathbf{I}_m\}$  denotes the embedding for  $\mathcal{V}_m$ .

#### A. MF-based Encoder

The core idea of MF-based encoders is to decompose the user-item rating matrix  $\mathbf{R}$  into two low-rank hidden embed-

dings  $\mathbf{U}$  and  $\mathbf{I}$ . The approximation of the rating matrix  $\mathbf{R}$  can be expressed as:

$$\mathbf{R} \approx \hat{\mathbf{R}} = \mathbf{U}\mathbf{I}^T, \quad (1)$$

where  $^T$  means the transpose operation for the matrix. The loss function can be defined as:

$$\min_{\mathbf{U}, \mathbf{I}} \|\mathbf{R} - \hat{\mathbf{R}}\|_F^2 + \lambda(\|\mathbf{E}\|_F^2), \quad (2)$$



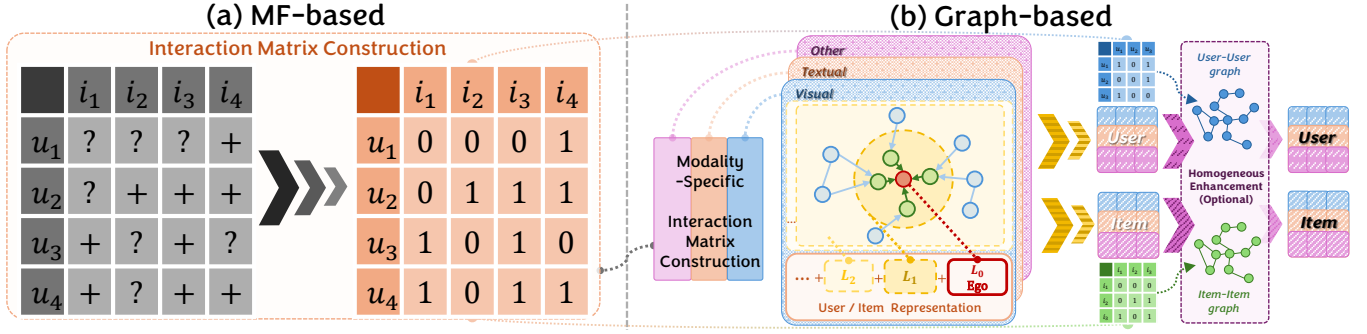


Fig. 3. The illustration of two types of encoders.

where  $\|\cdot\|_F$  means Frobenius norm and  $\lambda$  is the regularization parameter to control overfitting.

Due to the diversified features of the item, we categorize the strategies into a unified MF-based encoder and multiple MF-based encoders. The design of both relies on the choice of multimodal fusion, which we will discuss in detail in Section V. Roughly speaking, the approximation of the rating matrix  $\mathbf{R}$  for MRS can be expressed as:

$$\mathbf{R} \approx \hat{\mathbf{R}} = \mathbf{U}\mathbf{I}^T, \quad \mathbf{I} = \text{Aggr}(\mathbf{I}_m), \quad (3)$$

$$\mathbf{R} \approx \hat{\mathbf{R}} = \text{Aggr}(\mathbf{U}\mathbf{I}_m^T), \quad (4)$$

where  $\text{Aggr}(\cdot)$  denotes multimodal fusion. The loss function can be defined as:

$$\min_{\mathbf{U}, \mathbf{I}} \|\mathbf{R} - \hat{\mathbf{R}}\|_F^2 + \lambda \left( \sum_{m \in M} \|\mathbf{E}_m\|_F^2 \right). \quad (5)$$

### B. Graph-based Encoder

The core idea of Graph-based encoders is to use the features of nodes and the structural information of the graph to learn the representation of nodes. We introduce the commonly used graph-based encoder paradigm Graph Convolution Network (GCN) [85]. For a graph  $\mathcal{G}$  and its adjacency matrix  $\mathbf{A}$ , a layer propagation of GCN can be expressed as:

$$\mathbf{E}^{(l)} = \sigma(\tilde{\mathbf{D}}^{-\frac{1}{2}} \tilde{\mathbf{A}} \tilde{\mathbf{D}}^{-\frac{1}{2}} \mathbf{E}^{(l-1)} \mathbf{W}^{(l-1)}), \quad (6)$$

where  $\mathbf{E}^{(l)}$  is the  $l$  layer hidden embeddings for  $\mathcal{V}$ ,  $\tilde{\mathbf{A}} = \mathbf{A} + \mathbf{I}$ ,  $\tilde{\mathbf{D}}$  is the degree matrix of  $\tilde{\mathbf{A}}$ , where  $\tilde{\mathbf{D}}_{ii} = \sum_j \tilde{\mathbf{A}}_{ij}$ .  $\mathbf{W}^{(l-1)}$  is the weight matrix for  $l-1$  layer and  $\sigma(\cdot)$  is the active function. In the recommendation domain, LightGCN [86] proves that weight matrices and active functions in GCNs are useless and even increase training difficulty. To this end, a widely-used simplified GCN can be expressed as:

$$\mathbf{E}^{(l)} = \tilde{\mathbf{D}}^{-\frac{1}{2}} \tilde{\mathbf{A}} \tilde{\mathbf{D}}^{-\frac{1}{2}} \mathbf{E}^{(l-1)}, \quad (7)$$

Then it stacks multiple layers by  $\bar{\mathbf{E}} = \text{Stack}_{l \in L}(\mathbf{E}^{(l)})$ , where  $L$  is the total layer number of GCN. We simply define this entire Graph-based Encoder as  $\bar{\mathbf{E}} = \text{GCN}(\mathbf{U}, \mathbf{I})$ . The entire representation  $\bar{\mathbf{E}}$  can be split into user and item parts by  $\bar{\mathbf{U}}, \bar{\mathbf{I}} = \text{Sp}(\bar{\mathbf{E}})$ . To better mine user and item representations, two homogeneous type graphs, user-user and item-item, are proposed. A portion of Graph-based encoders will use either

one or all of them to better learn the representations. Homogeneous graphs first retain the top- $k$  items/users by similarity:

$$\mathbf{S}_{i,i'}^I = \begin{cases} 1, & \mathbf{S}_{i,i'}^I \in \text{top-}k(\mathbf{S}_{i,i'}^I) \\ 0, & \text{otherwise} \end{cases}, \quad (8)$$

$$\mathbf{S}_{u,u'}^U = \begin{cases} 1, & \mathbf{S}_{u,u'}^U \in \text{top-}k(\mathbf{S}_{u,u'}^U) \\ 0, & \text{otherwise} \end{cases}. \quad (9)$$

To mitigate the issue of significant disparities in the feature distributions between vertices with high degrees and those with low degrees within graph-structured data, it is a common practice to apply symmetric normalization to the adjacency matrix, denoted as  $\mathbf{S}^U = (\mathbf{D}^U)^{-1/2} \mathbf{S}^U (\mathbf{D}^U)^{-1/2}$  and  $\mathbf{S}^I = (\mathbf{D}^I)^{-1/2} \mathbf{S}^I (\mathbf{D}^I)^{-1/2}$ , where  $\mathbf{D}^U$  and  $\mathbf{D}^I$  represent the diagonal degree matrix of  $\mathbf{S}^U$  and  $\mathbf{S}^I$ , respectively. This normalization process is crucial as it adjusts the influence of each vertex based on its connectivity, thereby preventing vertices with a higher degree from disproportionately dominating the feature representations. Then propagate  $\bar{\mathbf{U}}/\bar{\mathbf{I}}$  through:

$$\hat{\mathbf{U}} = (\mathbf{S}^U)^{L_u} \bar{\mathbf{U}}, \quad \hat{\mathbf{I}} = (\mathbf{S}^I)^{L_i} \bar{\mathbf{I}}, \quad (10)$$

where  $L_u$  and  $L_i$  are the layer number of the user-user graph and item-item graph, respectively. These two representations can optionally enhance user and item representations by  $\tilde{\mathbf{U}} = \hat{\mathbf{U}} + \bar{\mathbf{U}}$  and  $\tilde{\mathbf{I}} = \hat{\mathbf{I}} + \bar{\mathbf{I}}$ .

For the sake of simplicity, we define the composite Graph-based encoder including optional user-user and item-item graphs as:

$$\tilde{\mathbf{U}}, \tilde{\mathbf{I}} = \text{C-GCN}(\mathbf{U}, \mathbf{I}) \quad (11)$$

The approximation of the rating matrix  $\mathbf{R}$  for MRS can be expressed as:

$$\mathbf{R} \approx \hat{\mathbf{R}} = \tilde{\mathbf{U}}\tilde{\mathbf{I}}^T, \quad \tilde{\mathbf{U}}, \tilde{\mathbf{I}} = \text{C-GCN}(\mathbf{U}, \mathbf{I}). \quad (12)$$

The loss function can be defined as Eq 1. Same as MF-based encoders, Graph-based encoders can also be categorized into a unified Graph-based encoder and multiple Graph-based encoders with different fusion strategies, detailed in Section V. Roughly speaking, the approximation of the rating matrix  $\mathbf{R}$  for MRS can be expressed as:

$$\mathbf{R} \approx \hat{\mathbf{R}} = \tilde{\mathbf{U}}\tilde{\mathbf{I}}^T, \quad \tilde{\mathbf{U}}, \tilde{\mathbf{I}} = \text{C-GCN}(\mathbf{U}, \text{Aggr}(\mathbf{I}_m)), \quad (13)$$

$$\mathbf{R} \approx \hat{\mathbf{R}} = \text{Aggr}(\tilde{\mathbf{U}}\tilde{\mathbf{I}}_m^T), \quad \tilde{\mathbf{U}}, \tilde{\mathbf{I}}_m = \text{C-GCN}(\mathbf{U}, \mathbf{I}_m). \quad (14)$$

The loss function can be defined as Eq 5.



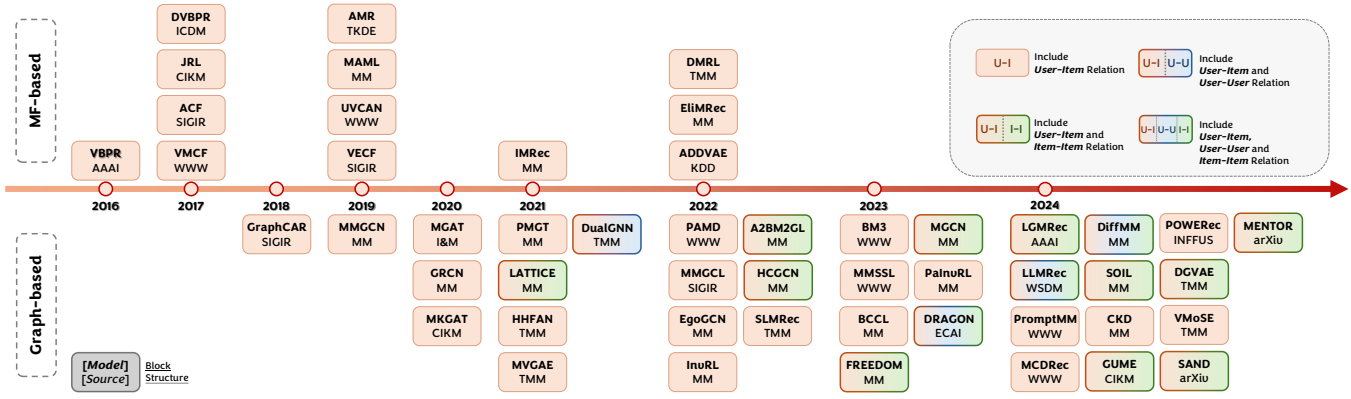


Fig. 4. Taxonomy of Encoders.

### C. Taxonomy

Following this analysis, we categorize recent works in MRS into MF-based and Graph-based encoders, as summarized in Figure 4. For Graph-based encoders, we delineate the types of relationships, including user-item (U-I), user-user (U-U), and item-item (I-I) relations.

**MF-based Encoder:** VBPR [29] integrates visual features extracted by a pre-trained deep CNN into MF for enhanced preference prediction. VMCf [30] builds a product-affinity network that incorporates visual appearance and inter-item relationships. ACF [31] introduces item- and component-level attention mechanisms to better handle implicit feedback, marking a first in applying attention in CF. JRL [32] and its extension, eJRL, employ multi-view machine learning to merge diverse information sources, improving top-N recommendations without needing to retrain for new data. DVBPR [33] and UVCAN [36] enhance recommendations by integrating visual signals and employing co-attention mechanisms, respectively. MAML [37] and ADDVAE [48] focus on modeling user preferences using concatenated textual and visual inputs through neural networks, with ADDVAE also exploring disentangled representations. AMR [38] addresses vulnerabilities in MRS using adversarial learning to create more robust models. IMRec aligns recommendation processes with user reading habits by focusing on local news details. ELIMRec [49] and DMRL [55] utilize causal inference and disentangled representations, respectively, to reduce biases and better capture independent modal factors.

**Graph-based Encoder:** GraphCAR [34] combines multimedia content with the traditional CF method. MMGCN [13] employs GCNs to learn representations for each modality, which are then fused with ID embeddings to form the final item representations. Building on the MMGCN framework, MGAT [39] uses standard GCN aggregation and a similar fusion method for combining results. GRCN [40] refines user-item interaction graphs by identifying and cutting noise edges. MKGAT [41] constructs a multimodal knowledge graph with an entity-based approach, using specialized encoders for different data types and an attention layer for effective information aggregation from neighboring entities. PMGT [43]

is a pre-trained model that utilizes fused multimodal features and interactions, employing an attention mechanism to derive multimodal embeddings. These embeddings are then enhanced with position and role-based embeddings to initialize node embeddings for pre-training and subsequent downstream tasks. DualGNN [14] introduces a user-user graph to uncover hidden preference patterns among users. LATTICE [15] develops an item-item graph to detect semantically correlated signals among items. HHFAN [44] develops a heterogeneous graph incorporating user, item, and multimodal information, and uses random walks to sample neighbors based on node type. It employs a Fully Connected (FC) layer to unify various node vectors into a single space and uses LSTM for aggregating embeddings of the same node type within the intra-type feature aggregation network. MVGAE [45] is a multimodal variational graph auto-encoder model that utilizes the modality-specific variational encoder to learn the node representation. PAMD [46] employs a disentangled encoder to separate common and unique characteristics of objects into respective representations, and uses contrastive learning for cross-modality alignment of these representations. MMGCL [47] integrates self-supervised learning with graph-based approaches for micro-video recommendation, featuring an innovative negative sampling method that highlights inter-modality relationships. EgoGCN [50] introduces an effective graph fusion method, which is not confined to unimodal graph information propagation but aggregates informative inter-modal messages from neighboring nodes. InvRL [51] addresses spurious correlations in multimedia recommendations by learning stable item representations across diverse environments. A2BM2GL [52] integrates collaborative and semantic representation learning to effectively model nodes and video features. An anti-bottleneck module with attention mechanisms enhances node relationship expressiveness. Additionally, an adaptive recommendation loss dynamically adjusts to user preference variations, improving item recommendation accuracy. HCGCN [53] extends the LATTICE framework by utilizing co-clustering and item-clustering losses to refine user-item preference feedback and adjust modality importance. SLMRec [54] proposes a self-supervised learning framework for multimodal recommendations, establishing a node self-discrimination task to reveal

hidden multimodal patterns of items. BM3 [56] simplifies SLMRec by replacing the random negative example sampling mechanism with a dropout strategy. MMSSL [57] designs a modality-aware interactive structure learning paradigm via adversarial perturbations, and proposes a cross-modal comparative learning method to disentangle the common and specific features among modalities. BCCL [58] integrates a bias constraint module for data augmentation, a modal awareness module, and a sparse enhancement module to collaboratively produce high-quality samples. FREEDOM [59] refines LAT-TICE by freezing the item-item graph and reducing noise in the user-item graph. MGCN [60] purifies modal features using item behavior information to reduce noise contamination and models modal preferences based on user behavior. PaInvRL [61] adaptively balances ERM (Empirical Risk Minimization) and IRM (Invariant Risk Minimization) losses to achieve Pareto-optimal solutions, effectively enhancing model performance by minimizing losses for Pareto optimality. DRAGON [62] learns the dual representations of users and items by constructing homogeneous and heterogeneous graphs. LGM-Rec [63] integrates local embeddings, which capture local topological nuances, with global embeddings, which consider hypergraph dependencies. LLMRec [64] employs three simple yet effective LLM-based graph augmentation strategies to enhance recommendation performance. PromptMM [65] enhances knowledge distillation by disentangling collaborative relationships to enable augmented distillation. MCDRec [66] integrates modal perceptual features with collaborative information to improve item representation and employs diffusion-aware representation to denoise the user-item interaction graph. DiffMM [67] introduces a well-designed modality-aware graph diffusion model to improve modality-aware user representation learning. SOIL [68] exploits candidate items from the perspective of constructing interest-aware graphs. CKD [69] aims to solve the modal imbalance problem and make the best use of all modalities. MENTOR [75] leverage aligned modalities while preserving interaction information with multi-level cross-modal alignment. GUME [70] achieves outstanding performance in long-tail scenarios by combining MGCN and MENTOR. POWERec [71] leverages prompt learning to model modality-specific interests. DGVAE [72] utilizes GCNs to encode and disentangle ratings and multimodal information, learning item representations from an item-item graph and enhancing interpretability by projecting multimodal data into text. VMoSE [73] enhances robustness by adaptive sampling and fusing noisy multi-modal signals based on uncertainty estimates. SAND [74] aligns modal-generic representations efficiently without negative sampling and distinctly separates modal-unique representations to preserve modality-independent information.

## V. MULTIMODAL FUSION

Multimodal fusion represents a critical investigation area in MRS. The effectiveness of multimodal integration is significantly influenced by the timing of when different modalities are fused, as well as the strategies employed to execute this fusion. Therefore, it becomes essential to systematically

analyze and categorize existing works based on these two perspectives: timing and strategy.

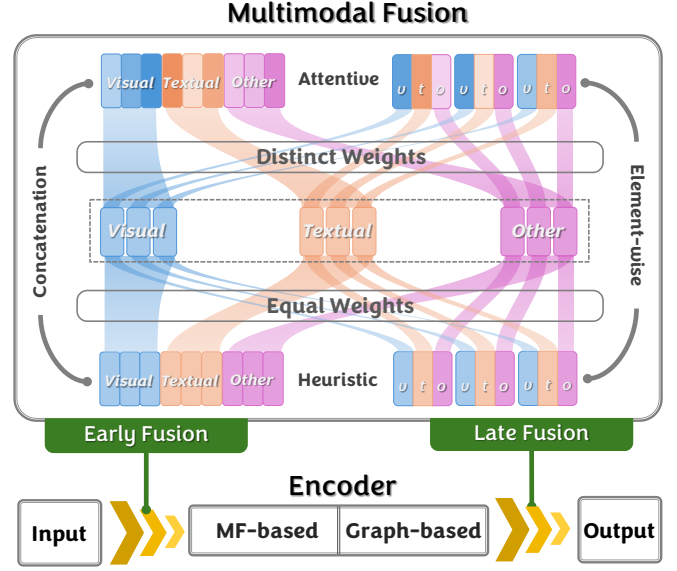


Fig. 5. The illustration of Multimodal Fusion.

### A. Timing Perspective

We first categorize all recent works in MRS from the timing perspective. Specifically, it can be decided into **Early fusion** and **Late fusion** strategies, respectively. Formally,

$$\text{Early fusion: } \bar{\mathbf{E}} = \text{Encoder}(\mathbf{U}, \text{Aggr}(\mathbf{I}_m)), \quad (15)$$

$$\text{Late fusion: } \bar{\mathbf{E}} = \text{Aggr}(\text{Encoder}(\mathbf{U}, \mathbf{I}_m)), \quad (16)$$

where we utilize the same preliminaries as detailed in Section IV. Encoder denotes various encoders. Early fusion denotes fusing all modalities before graph message propagation and aggregation. Late fusion denotes fusing all modalities after graph message propagation and aggregation. Although both early and late fusion timings demonstrate commendable performance in various scenarios, each approach exhibits distinct limitations. Specifically, early fusion tends to integrate modalities at the early stage, which can result in modal-specific features not being fully exploited due to premature integration. On the other hand, the late fusion strategy, which combines modalities at a later stage, faces challenges in fully capturing and leveraging the correlations among different modalities. It might limit the model's ability to extract the correlation among different modalities, potentially reducing the overall effectiveness of leveraging multimodal information.

### B. Strategy Perspective

From the strategy perspective, all recent works can be fine-grained categories from two dimensions, **Element-wise or Concatenation** and **Attentive or Heuristic**. Many works in MRS have introduced subtle adjustments to the fusion strategy. Despite these variations, it is commonly feasible to categorize these approaches in a broad manner based on the underlying motivations driving their use, as delineated in Table IV.

TABLE IV

SIMPLIFIED SUMMARY OF FUSION STRATEGY. WE SIMPLIFIED EXPRESS THEM WITH ONLY TWO MODALITIES VISUAL AND TEXTUAL.  $E_v$  AND  $E_t$  DENOTE REPRESENTATIONS FOR VISUAL AND TEXTUAL MODALITIES, RESPECTIVELY.

Fusion Strategy	Concatenation	Element-wise
Heuristic	$E_v    E_t$	$E_v + E_t$
Attentive	$\alpha_v E_v    \alpha_t E_t$	$\alpha_v E_v + \alpha_t E_t$

$||$  denotes concatenation operation.  $\alpha_v$  and  $\alpha_t$  denote learnable weights for each modality.

The method of element-wise fusion provides a more profound integration of different modalities compared to the concatenation approach. However, this deeper integration might inadvertently amplify the inherent noise present within the modality data. On the other hand, the attentive approach, as opposed to heuristic methods, offers a more dynamic allocation of modality weights, allowing for a more adaptive and responsive handling of input features based on their relevance to the specific task. Despite its advantages, the attentive mechanism incurs significantly higher computational costs and increases the complexity of the training phase.

The fusion strategy often interacts synergistically with the timing of fusion, thereby influencing the model performance significantly. Recognizing this interplay, we conduct a comprehensive summary of the recent strategies and timing of modality fusion within MRS, as shown in Figure 5. Our objective is to provide researchers with a clearer perspective on how these perspectives coalesce to influence model performance.

**Early Fusion:** VBPR [29] and GraphCAR [34] directly utilize visual features within MF frameworks. In an extension to this approach, VMCF [30] introduces visual matrix co-factorization. DVBPR [33] leverages a deeper neural network to enhance the representation of visual data. Attention mechanisms have been increasingly applied to dynamically allocate weights across different modalities. Works such as ACF [31], VECF [35], UVCAN [36], MAML [37], GRCN [40], IMRec [42], PMGT [43], BCCL [58], and MCDRec [66] utilize attention mechanisms to effectively allocate weights for different modalities. MKGAT [41] constructs a multimodal knowledge graph that fuses modal features to enhance recommendation systems. More recently, LATTICE [15] exploits raw features to construct item-item graphs for each modality and dynamically fuses weights using an attention mechanism to form modality-specific graphs. Moreover, HCGCN [53] employs a multimodal item-item graph to enhance the user-item graph, enabling the discovery of user preferences within similar items and facilitating simultaneous clustering.

**Late Fusion:** Models such as JRL [32], MMGCN [13], DualGNN [14], MMGCL [47], EgoGCN [50], FREEDOM [59], MGCN [60], DRAGON [62], LGMRec [63], DiffMM [67], SOIL [68], and POWERec [71] first independently processing different modalities and then combining them effectively to enhance prediction accuracy. Moreover, HHFAN [44] employs self-attention-aware neural networks to integrate information from all modalities, enhancing the model's ability to focus on more informative features. Attention mechanisms are applied in PAMD [46], A2BM2GL [52], and LLMRec [64], which dynamically adjust the influence of each modality on the final prediction, providing a flexible and context-sensitive fusion strategy. Similarly, MVGAE [45] uses a product-of-experts framework to harmonize modality-specific distributions, ensuring an effective fusion of modal inputs. Furthermore, some methods focus on leveraging auxiliary modality alignment tasks to better integrate modalities. ADDVAE [48], MMSSL [57], PromptMM [65], GUME [70], DGVAE [72], SAND [74], VMoSE [73], and MENTOR [75] all incorporate auxiliary tasks to help modality fusion. Additionally, EliMRec [49] integrate modalities from a causal perspective, and CKD [69] applies the Average Treatment Effect (ATE) strategy for modality fusion, aiming to quantify the impact of each modality on the recommendation outcome effectively. SMORE [76] projects the multi-modal features into the frequency domain and leverages the spectral space for fusion.

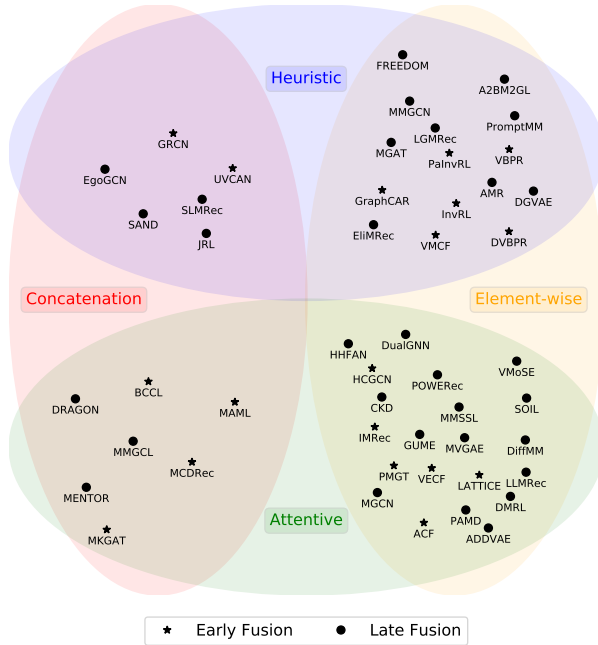


Fig. 6. Taxonomy of Multimodal Fusion. five-pointed star and circle symbols denote early fusion and late fusion, respectively.

### C. Taxonomy

Following the above analysis, we categorize recent works from both fusion timing and fusion strategy perspectives, as summarized in Figure 6. We further provide a specific analysis of all recent MRS works from a fusion timing perspective.

Since the goal of each fusion strategy is relatively fixed, we will not go into too much detail from the perspective of the fusion strategy.

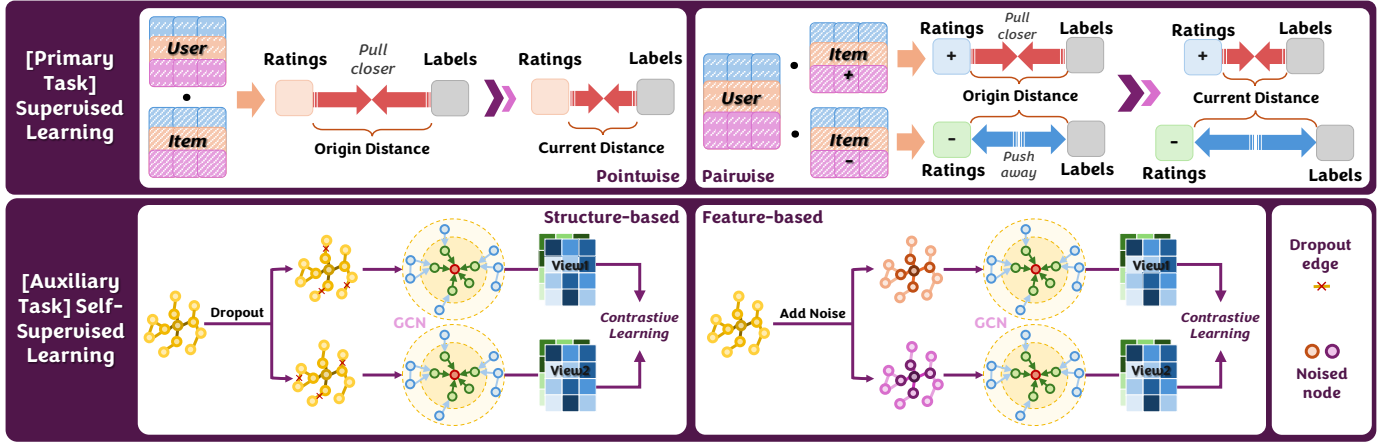


Fig. 7. The illustration of Loss Functions.

## VI. LOSS FUNCTION

MRS utilizes loss functions consisting of two key types: primary and auxiliary tasks, as shown in Figure 7. The primary tasks involve supervised learning, where clearly defined labels guide the model's learning, ensuring accurate predictions from labeled data. Conversely, the auxiliary tasks employ self-supervised learning (SSL) [24], which generates supervisory signals from the data's inherent structure, independent of external labels. This method enables effective use of unlabeled data, allowing the system to extract meaningful representations and maintain accuracy, even when labeled data is sparse. We introduce these two types of loss functions in detail.

### A. Supervised Learning

**Preliminary:** In the RS scenario, we simplified define the recommendation model as  $f(\cdot)$ , it can predict the score between user  $u$  and item  $i$  by:

$$y_{u,i} = f(\mathbf{R}, u, i), \quad (17)$$

where  $\mathbf{R} = [r_{u,i}]^{|U| \times |I|}$  represents the user-item interaction matrix. In the MRS scenario, the multimodal recommendation model will be expanded as  $f_m(\cdot)$ , it can predict the score between user  $u$  and item  $i$  with multimodal information by:

$$y_{u,i} = f_m(\mathbf{R}, u, i, \mathbf{M}), \quad (18)$$

where  $\mathbf{M}$  incorporates external item attributes with different modalities.

The supervised learning task is to force the predicted score between any user  $u$  and item  $i$  closer to the ground label. It can be divided into **Pointwise Loss** and **Pairwise Loss**.

For Pointwise Loss, we introduce two common loss functions that are also widely used in many machine learning downstream tasks, such as Computer Vision and Natural Language Processing. 1) Mean Squared Error (MSE) [25] and 2) Cross-Entropy Loss (CE) [26].

MSE can be defined as:

$$\frac{1}{|\mathcal{D}|} \sum_{(u,i) \in \mathcal{D}} (r_{u,i} - y_{u,i})^2, \quad (19)$$

where  $\mathcal{D}$  denotes the training set. This loss function pulls the predicted score for interacted user-item pairs more approximately to 1, and non-interacted user-item pairs more approximately to 0.

CE can be defined as:

$$-\frac{1}{|\mathcal{D}|} \sum_{(u,i) \in \mathcal{D}} [r_{u,i} \log(y_{u,i}) + (1 - r_{u,i}) \log(1 - y_{u,i})]. \quad (20)$$

Utilizing CE in the RS (MRS) fields treats the prediction function as a two-fold classification task, classifying user-item pairs into interacted and non-interacted categories. This improves classification accuracy by minimizing the discrepancy between the predicted probability distribution and the true label distribution, allowing the model to better learn the probability distribution of classified labels. However, due to the natural sparsity of user-item interaction data, we need to value the scarce labeled data more, and thus further propose the Pairwise Loss.

For Pointwise Loss, the learning objective is to pull closer the positive pairs, while pushing away the negative pairs. We introduce two common loss functions that are also widely used loss function. 1) Bayesian Personalized Ranking (BPR) [27] and 2) Hinge Loss [28].

BPR can be defined as:

$$\sum_{(u,i^+,i^-) \in \mathcal{D}} -\log(\sigma(y_{u,i^+} - y_{u,i^-})), \quad (21)$$

where  $y_{u,i^+}$  and  $y_{u,i^-}$  are the ratings of user  $u$  to the positive item  $i^+$  and negative item  $i^-$ .  $\sigma$  is the active function.

Hinge Loss can be defined as:

$$\max(0, 1 - r_{u,i} \cdot y_{u,i}). \quad (22)$$

Hinge aims to not only correctly classify data points, but also to maximize the spacing of classification decision boundaries.

### B. Self-supervised Learning

**Preliminary:** In the RS scenario, self-supervised learning first creates two different views for contrastive learning, we simply define this view creator as:

$$w = \mathcal{C}(\mathbf{R}, \mathbf{E}), \quad (23)$$



where  $\mathcal{C}(\cdot)$  is a view creator.  $\mathbf{R}$  and  $\mathbf{E}$  are use the same definition as Section IV. In the MRS scenario, the view creator can be further defined as:

$$w = \mathcal{C}(\mathbf{R}, \mathbf{E}, \mathbf{M}). \quad (24)$$

The view creator can be divided into two types: 1) **Feature-based SSL** and 2) **Structure-based SSL**.

Feature-based SSL is designed to generate multiple views of the same data instance by perturbing the features. This approach is strategically employed to enhance the robustness of the learned representations. By introducing variations in the input features, the model is compelled to focus on the invariant aspects across these modifications, thereby acquiring a deeper, more generalized understanding of the user-item interaction. It can be simplify expressed as:

$$\omega = \mathcal{C}_{feature}(\mathbf{R}, \mathbf{E}, \mathbf{M}) = (\mathbf{R}, \text{Perturb}(\mathbf{E}), \mathbf{M}), \quad (25)$$

where  $\text{Perturb}(\cdot)$  can be an MLP, feature dropout, adding random noise, etc. It is worth noting that different modalities can be naturally considered as two feature-based views. Applying feature-based SSL for two different modalities representation can be seen as modality alignments [75].

Structure-based SSL is designed to generate multiple views of the same data instance by perturbing the graph structures. This approach is employed to intricately capture the complex dependencies and interactions inherent within graph structures, which are pivotal in enhancing the robustness and performance of graph-based learning models. Such manipulations enable the learning algorithms to discern and generalize from the essential features of the data, potentially leading to more effective and insightful representations in domains where graph-based data is prevalent. It can be simplify expressed as:

$$\omega = \mathcal{C}_{structure}(\mathbf{R}, \mathbf{E}, \mathbf{M}) = (\text{Perturb}(\mathbf{R}), \mathbf{E}, \mathbf{M}), \quad (26)$$

where  $\text{Perturb}(\cdot)$  can be an MLP, node/edge dropout, adding random noise, etc.

The self-supervised learning task is to force the generated views closer to enhance the representation ability of the MRS model. We introduce two widely used self-supervised loss functions. 1) InfoNCE [87] and 2) Jensen-Shannon divergence (JS) [88].

InfoNCE is a variant of Noise Contrastive Estimation [89], which gained wide adoption as a self-supervised learning loss function in RS. It can be expressed as:

$$\mathbb{E}[-\log \frac{\exp(f(\omega'_i, \omega''_i))}{\sum_{i,j} \exp(f(\omega'_i, \omega''_j))}], \quad (27)$$

where  $f(\cdot)$  represents a critic function that calculates a score indicating the similarity between two views. The term  $\exp f(\omega'_i, \omega''_i)$  corresponds to the score of positive pairs, while the term  $\sum \exp(f(\omega'_i, \omega''_j))$  encompasses both the numerator and the scores of all negative pairs.

In addition to using InfoNCE estimation for mutual information, the lower bound can also be estimated using the Jensen-Shannon (JS) divergence. The derived learning objective is akin to combining InfoNCE with a standard binary cross-entropy loss [90], applied to positive pairs and negative pairs.

$$\mathbb{E}[-\log \sigma(f(\omega'_i, \omega''_i))] - \mathbb{E}[\log(1 - \sigma(f(\omega'_i, \omega''_j)))], \quad (28)$$

where  $\sigma$  represents the sigmoid function used to normalize the output of the critic function. The main idea behind this optimization is to assign the label 1 to positive pairs and 0 to negative pairs, thereby increasing the predicted value for positive pairs and enhancing the similarity between them.

### C. Taxonomy

To provide researchers with a clearer framework for selecting the most effective loss function for their work, we have systematically categorized recent works, as Figure 8. Most existing MRs methods leverage pairwise loss function, only VMCF [30], VECF [35], UVCAN [36], IMRec [42], PMGT [43], ADDVAE [48], A2BM2GL [52], BM3 [56], BCCL [58], DGVAE [72], and VMoSE [73] adopt point wise loss function. Furthermore, we provide a specific analysis for self-supervised learning of all recent MRS works.

**Feature-based SSL:** PAMD [46] uses self-supervised signals to assist in learning disentangled representation from the feature level. A2BM2GL [52] proposes aligning disentangled factors learned from ratings and textual content based on regularization and compositional de-attention mechanism. BCCL [58] further introduces a bias-constrained data augmentation method to ensure the quality of augmentation samples in contrastive learning. MGCN [60], LGMRec [63], SOIL [68], GUME [70] use contrastive learning from the feature level to improve the representation quality. PromptMM [65] introduces a learnable prompt module that dynamically bridges the semantic gap between the multi-modal context encoding in the teacher model and the collaborative relation modeling in the student model. MCDRec [66], DiffMM [67], DGVAE [72], and VMoSE [73] optimize KL divergence to compels posterior distribution closer to the prior distribution. SAND [74] uses self-supervised signals to distinguish modal-unique and modal-generic representations.

**Structure-based SSL:** MMGCL [47] uses modality masking and modality edge dropout to enhance the modal consistency of representations through self-supervision from a structural perspective. A2BM2GL [52] utilizes the attention mechanism to dynamically learn the importance weights of short-range and long-range neighboring nodes jointly to obtain more expressive representations. HCGCN [53] uses contrastive loss from a structural perspective to coordinate multimodal features. SLMRec [54] uses contrastive learning from a variety of common structural perspectives to improve representation quality.

**Mixed SSL:** BM3 [56] uses multiple common self-supervised signals at the feature level and structure level to enhance representation. MMSSL [57] and MENTOR [75] align different modal representations at the feature level and enhance representations using common contrastive learning from the feature perspective.

## VII. FUTURE DIRECTION

In this section, we explore potential future directions for the field of MRS. Our goal is to stimulate and foster further research, development, and innovation within this rapidly evolving domain. By identifying and discussing emerging

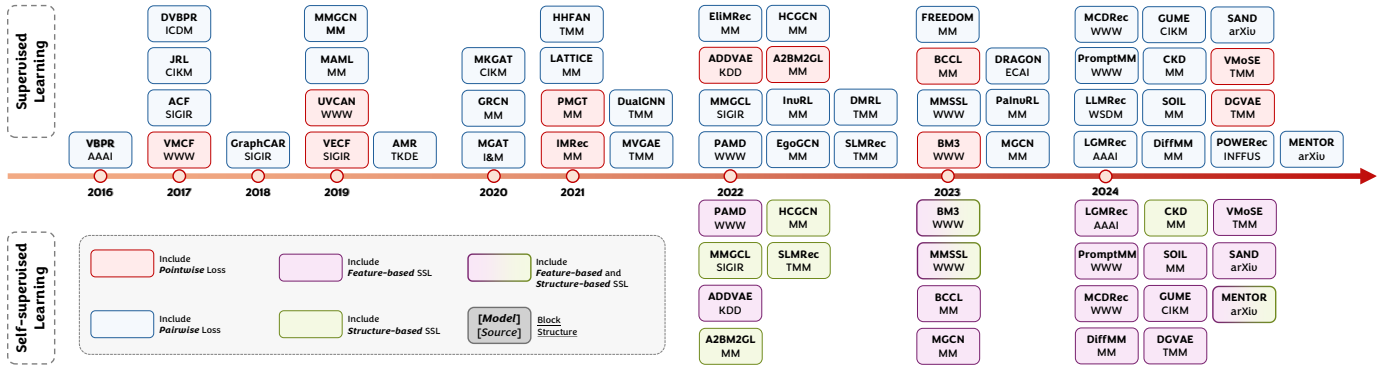


Fig. 8. Taxonomy of Loss Functions.

trends and unresolved challenges, we hope to inspire continued academic inquiry and technological advancement that will drive the next wave of breakthroughs in MRS.

#### A. Towards Unified MRS Model

Existing models in the MRS field typically segregate the feature extraction and representation encoding into two distinct processes. The former often leverages existing pre-trained models, while the latter receives more focused attention. However, this segregated process results in inherent multi-modal noise, leading to a disconnect between the extracted features and their subsequent encoding. Consequently, there is an urgent need for a unified model that integrates these processes more cohesively. Such a unified model will enhance the relevance and efficiency of the multimodal data, thereby improving the overall accuracy of the recommendation systems in leveraging complex multimodal information.

#### B. Resolves Cold-start Problem

Existing RS models typically operate under the assumption of a fixed number of users and items during the training phase, which poses challenges when adapting to the continuous influx of new data. In practical settings, these models are deployed within dynamic environments where new user-item interactions, as well as new users and items, are frequently introduced—a phenomenon often referred to as the cold-start problem. To address this issue, it is beneficial to leverage multimodal information in cold-start scenarios. The integration of diverse data modalities—such as textual descriptions, images, and metadata associated with users and items—can significantly enhance the model’s ability to understand and predict the preferences of new users and the attributes of new items effectively. By exploiting multimodal information, RS models can generate more accurate and reliable recommendations even when confronted with limited interaction data, thereby improving their adaptability and performance in dynamically evolving environments.

#### C. Towards Richer Variety of Modalities

Existing MRS models demonstrate considerable effectiveness in utilizing textual and visual modalities. However, the

burgeoning richness of information available on the Internet presents opportunities for leveraging a broader variety of modalities. This expansion includes auditory, olfactory, and kinesthetic data, among others, which can enrich the understanding and personalization capabilities of RS models.

In response to this evolving landscape, there is an urgent need for MRS models that can effectively integrate and synthesize information from these varied data sources. By harnessing a wider array of modal information, MRS models can achieve a more holistic understanding of user preferences and item characteristics. This comprehensive approach not only enhances the accuracy and relevance of recommendations but also significantly improves user engagement and satisfaction by catering to diverse sensory preferences and interaction styles. Thus, developing MRS models capable of effectively processing and integrating multiple modalities is crucial for advancing the state of the art in RS.

### VIII. CONCLUSION

The primary goal of this survey is to thoroughly examine the recent advancements in MRS and provide a technical analysis of various models. Our discussion categorizes existing MRS models into four critical aspects: Feature Extraction, Encoder, Multimodal Fusion, and Loss Function. Moreover, we review the technical contributions of existing works and explore potential future avenues for advancing and refining MRS technologies. Our contributions extend beyond mere summarization. We tailored technological taxonomy and proposed potential directions for future research. This survey serves as a valuable resource for researchers in the field, offering insights and guidance into the evolving landscape of multimedia recommendations.

### REFERENCES

- [1] Y. Deldjoo, M. Schedl, and P. Knees, “Content-driven music recommendation: Evolution, state of the art, and challenges,” *arXiv preprint arXiv:2107.11803*, 2021.
- [2] Y. Deldjoo, F. Nazary, A. Ramisa, J. McAuley, G. Pellegrini, A. Bellogin, and T. Di Noia, “A review of modern fashion recommender systems,” *arXiv preprint arXiv:2202.02757*, 2022.
- [3] J. Xu, Z. Chen, J. Li, S. Yang, W. Wang, X. Hu, and E. C.-H. Ngai, “Fourierkan-gcf: Fourier kolmogorov-arnold network—an effective and efficient feature transformation for graph collaborative filtering,” *arXiv preprint arXiv:2406.01034*, 2024.

- [4] J. Xu, Z. Chen, J. Li, S. Yang, H. Wang, and E. C. Ngai, "Aligngroup: Learning and aligning group consensus with member preferences for group recommendation," in *Proceedings of the 33rd ACM International Conference on Information and Knowledge Management*, 2024, pp. 2682–2691.
- [5] Q. Liu, J. Hu, Y. Xiao, J. Gao, and X. Zhao, "Multimodal recommender systems: A survey," *arXiv preprint arXiv:2302.03883*, 2023.
- [6] S. Zhang, L. Yao, A. Sun, and Y. Tay, "Deep learning based recommender system: A survey and new perspectives," pp. 1–38, 2019.
- [7] Q. Guo, F. Zhuang, C. Qin, H. Zhu, X. Xie, H. Xiong, and Q. He, "A survey on knowledge graph-based recommender systems," *IEEE Transactions on Knowledge and Data Engineering*, vol. 34, no. 8, pp. 3549–3568, 2020.
- [8] Y. Deldjoo, M. Schedl, P. Cremonesi, and G. Pasi, "Recommender systems leveraging multimedia content," *ACM Computing Surveys (CSUR)*, vol. 53, no. 5, pp. 1–38, 2020.
- [9] D. Jannach, A. Manzoor, W. Cai, and L. Chen, "A survey on conversational recommender systems," *ACM Computing Surveys (CSUR)*, vol. 54, no. 5, pp. 1–36, 2021.
- [10] S. Wu, F. Sun, W. Zhang, X. Xie, and B. Cui, "Graph neural networks in recommender systems: a survey," *ACM Computing Surveys*, vol. 55, no. 5, pp. 1–37, 2022.
- [11] X. Meng, H. Huo, X. Zhang, W. Wang, and J. Zhu, "A survey of personalized news recommendation," *Data Science and Engineering*, pp. 1–21, 2023.
- [12] H. Zhou, X. Zhou, Z. Zeng, L. Zhang, and Z. Shen, "A comprehensive survey on multimodal recommender systems: Taxonomy, evaluation, and future directions," *arXiv preprint arXiv:2302.04473*, 2023.
- [13] Y. Wei, X. Wang, L. Nie, X. He, R. Hong, and T.-S. Chua, "Mmgcn: Multi-modal graph convolution network for personalized recommendation of micro-video," in *Proceedings of the 27th ACM international conference on multimedia*, 2019, pp. 1437–1445.
- [14] Q. Wang, Y. Wei, J. Yin, J. Wu, X. Song, and L. Nie, "Dualgnn: Dual graph neural network for multimedia recommendation," *IEEE Transactions on Multimedia*, 2021.
- [15] J. Zhang, Y. Zhu, Q. Liu, S. Wu, S. Wang, and L. Wang, "Mining latent structures for multimedia recommendation," in *Proceedings of the 29th ACM International Conference on Multimedia*, 2021, pp. 3872–3880.
- [16] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [17] 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.
- [18] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," *arXiv preprint arXiv:1810.04805*, 2018.
- [19] N. Reimers and I. Gurevych, "Sentence-bert: Sentence embeddings using siamese bert-networks," *arXiv preprint arXiv:1908.10084*, 2019.
- [20] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [21] K. Cho, B. Van Merriënboer, D. Bahdanau, and Y. Bengio, "On the properties of neural machine translation: Encoder-decoder approaches," *arXiv preprint arXiv:1409.1259*, 2014.
- [22] A. Mnih and R. R. Salakhutdinov, "Probabilistic matrix factorization," *Advances in neural information processing systems*, vol. 20, 2007.
- [23] F. Wu, A. Souza, T. Zhang, C. Fifty, T. Yu, and K. Weinberger, "Simplifying graph convolutional networks," in *International conference on machine learning*. PMLR, 2019, pp. 6861–6871.
- [24] 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*, vol. 35, no. 1, pp. 857–876, 2021.
- [25] H. Marmolin, "Subjective mse measures," *IEEE transactions on systems, man, and cybernetics*, vol. 16, no. 3, pp. 486–489, 1986.
- [26] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *nature*, vol. 521, no. 7553, pp. 436–444, 2015.
- [27] S. Rendle, C. Freudenthaler, Z. Gantner, and L. Schmidt-Thieme, "Bpr: Bayesian personalized ranking from implicit feedback," *arXiv preprint arXiv:1205.2618*, 2012.
- [28] C. Cortes and V. Vapnik, "Support-vector networks," *Machine learning*, vol. 20, pp. 273–297, 1995.
- [29] R. He and J. McAuley, "Vbpr: visual bayesian personalized ranking from implicit feedback," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 30, no. 1, 2016.
- [30] C. Park, D. Kim, J. Oh, and H. Yu, "Do" also-viewed" products help user rating prediction?" in *Proceedings of the 26th International Conference on World Wide Web*, 2017, pp. 1113–1122.
- [31] J. Chen, H. Zhang, X. He, L. Nie, W. Liu, and T.-S. Chua, "Attentive collaborative filtering: Multimedia recommendation with item-and component-level attention," in *Proceedings of the 40th International ACM SIGIR conference on Research and Development in Information Retrieval*, 2017, pp. 335–344.
- [32] Y. Zhang, Q. Ai, X. Chen, and W. B. Croft, "Joint representation learning for top-n recommendation with heterogeneous information sources," in *Proceedings of the 2017 ACM on Conference on Information and Knowledge Management*, 2017, pp. 1449–1458.
- [33] W.-C. Kang, C. Fang, Z. Wang, and J. McAuley, "Visually-aware fashion recommendation and design with generative image models," in *2017 IEEE international conference on data mining (ICDM)*. IEEE, 2017, pp. 207–216.
- [34] Q. Xu, F. Shen, L. Liu, and H. T. Shen, "Graphcar: Content-aware multimedia recommendation with graph autoencoder," in *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval*, 2018, pp. 981–984.
- [35] X. Chen, H. Chen, H. Xu, Y. Zhang, Y. Cao, Z. Qin, and H. Zha, "Personalized fashion recommendation with visual explanations based on multimodal attention network: Towards visually explainable recommendation," in *Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2019, pp. 765–774.
- [36] S. Liu, Z. Chen, H. Liu, and X. Hu, "User-video co-attention network for personalized micro-video recommendation," in *The world wide web conference*, 2019, pp. 3020–3026.
- [37] F. Liu, Z. Cheng, C. Sun, Y. Wang, L. Nie, and M. Kankanhalli, "User diverse preference modeling by multimodal attentive metric learning," in *Proceedings of the 27th ACM international conference on multimedia*, 2019, pp. 1526–1534.
- [38] J. Tang, X. Du, X. He, F. Yuan, Q. Tian, and T.-S. Chua, "Adversarial training towards robust multimedia recommender system," *IEEE Transactions on Knowledge and Data Engineering*, vol. 32, no. 5, pp. 855–867, 2019.
- [39] Z. Tao, Y. Wei, X. Wang, X. He, X. Huang, and T.-S. Chua, "Mgat: Multimodal graph attention network for recommendation," *Information Processing & Management*, vol. 57, no. 5, p. 102277, 2020.
- [40] Y. Wei, X. Wang, L. Nie, X. He, and T.-S. Chua, "Graph-refined convolutional network for multimedia recommendation with implicit feedback," in *Proceedings of the 28th ACM international conference on multimedia*, 2020, pp. 3541–3549.
- [41] R. Sun, X. Cao, Y. Zhao, J. Wan, K. Zhou, F. Zhang, Z. Wang, and K. Zheng, "Multi-modal knowledge graphs for recommender systems," in *Proceedings of the 29th ACM international conference on information & knowledge management*, 2020, pp. 1405–1414.
- [42] J. Xun, S. Zhang, Z. Zhao, J. Zhu, Q. Zhang, J. Li, X. He, X. He, T.-S. Chua, and F. Wu, "Why do we click: visual impression-aware news recommendation," in *Proceedings of the 29th ACM International Conference on Multimedia*, 2021, pp. 3881–3890.
- [43] Y. Liu, S. Yang, C. Lei, G. Wang, H. Tang, J. Zhang, A. Sun, and C. Miao, "Pre-training graph transformer with multimodal side information for recommendation," in *Proceedings of the 29th ACM International Conference on Multimedia*, 2021, pp. 2853–2861.
- [44] D. Cai, S. Qian, Q. Fang, and C. Xu, "Heterogeneous hierarchical feature aggregation network for personalized micro-video recommendation," *IEEE Transactions on Multimedia*, vol. 24, pp. 805–818, 2021.
- [45] J. Yi and Z. Chen, "Multi-modal variational graph auto-encoder for recommendation systems," *IEEE Transactions on Multimedia*, vol. 24, pp. 1067–1079, 2021.
- [46] T. Han, P. Wang, S. Niu, and C. Li, "Modality matches modality: Pretraining modality-disentangled item representations for recommendation," in *Proceedings of the ACM Web Conference 2022*, 2022, pp. 2058–2066.
- [47] Z. Yi, X. Wang, I. Ounis, and C. Macdonald, "Multi-modal graph contrastive learning for micro-video recommendation," in *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2022, pp. 1807–1811.
- [48] N.-T. Tran and H. W. Lauw, "Aligning dual disentangled user representations from ratings and textual content," in *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 2022, pp. 1798–1806.
- [49] X. Liu, Z. Tao, J. Shao, L. Yang, and X. Huang, "Elimrec: Eliminating single-modal bias in multimedia recommendation," in *Proceedings of the 30th ACM International Conference on Multimedia*, 2022, pp. 687–695.



- [50] F. Chen, J. Wang, Y. Wei, H.-T. Zheng, and J. Shao, "Breaking isolation: Multimodal graph fusion for multimedia recommendation by edge-wise modulation," in *Proceedings of the 30th ACM International Conference on Multimedia*, 2022, pp. 385–394.
- [51] X. Du, Z. Wu, F. Feng, X. He, and J. Tang, "Invariant representation learning for multimedia recommendation," in *Proceedings of the 30th ACM International Conference on Multimedia*, 2022, pp. 619–628.
- [52] D. Cai, S. Qian, Q. Fang, J. Hu, and C. Xu, "Adaptive anti-bottleneck multi-modal graph learning network for personalized micro-video recommendation," in *Proceedings of the 30th ACM International Conference on Multimedia*, 2022, pp. 581–590.
- [53] Z. Mu, Y. Zhuang, J. Tan, J. Xiao, and S. Tang, "Learning hybrid behavior patterns for multimedia recommendation," in *Proceedings of the 30th ACM International Conference on Multimedia*, 2022, pp. 376–384.
- [54] Z. Tao, X. Liu, Y. Xia, X. Wang, L. Yang, X. Huang, and T.-S. Chua, "Self-supervised learning for multimedia recommendation," *IEEE Transactions on Multimedia*, 2022.
- [55] F. Liu, H. Chen, Z. Cheng, A. Liu, L. Nie, and M. Kankanhalli, "Disentangled multimodal representation learning for recommendation," *IEEE Transactions on Multimedia*, 2022.
- [56] X. Zhou, H. Zhou, Y. Liu, Z. Zeng, C. Miao, P. Wang, Y. You, and F. Jiang, "Bootstrap latent representations for multi-modal recommendation," in *Proceedings of the ACM Web Conference 2023*, 2023, pp. 845–854.
- [57] W. Wei, C. Huang, L. Xia, and C. Zhang, "Multi-modal self-supervised learning for recommendation," in *Proceedings of the ACM Web Conference 2023*, 2023, pp. 790–800.
- [58] W. Yang, Z. Fang, T. Zhang, S. Wu, and C. Lu, "Modal-aware bias constrained contrastive learning for multimodal recommendation," in *Proceedings of the 31st ACM International Conference on Multimedia*, 2023, pp. 6369–6378.
- [59] X. Zhou and Z. Shen, "A tale of two graphs: Freezing and denoising graph structures for multimodal recommendation," in *Proceedings of the 31st ACM International Conference on Multimedia*, 2023, pp. 935–943.
- [60] P. Yu, Z. Tan, G. Lu, and B.-K. Bao, "Multi-view graph convolutional network for multimedia recommendation," in *Proceedings of the 31st ACM International Conference on Multimedia*, 2023, pp. 6576–6585.
- [61] S. Huang, H. Li, Q. Li, C. Zheng, and L. Liu, "Pareto invariant representation learning for multimedia recommendation," in *Proceedings of the 31st ACM International Conference on Multimedia*, 2023, pp. 6410–6419.
- [62] H. Zhou, X. Zhou, L. Zhang, and Z. Shen, "Enhancing dyadic relations with homogeneous graphs for multimodal recommendation," in *ECAI 2023*. IOS Press, 2023, pp. 3123–3130.
- [63] Z. Guo, J. Li, G. Li, C. Wang, S. Shi, and B. Ruan, "Lgmrec: Local and global graph learning for multimodal recommendation," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 38, no. 8, 2024, pp. 8454–8462.
- [64] W. Wei, X. Ren, J. Tang, Q. Wang, L. Su, S. Cheng, J. Wang, D. Yin, and C. Huang, "Llmrec: Large language models with graph augmentation for recommendation," in *Proceedings of the 17th ACM International Conference on Web Search and Data Mining*, 2024, pp. 806–815.
- [65] W. Wei, J. Tang, L. Xia, Y. Jiang, and C. Huang, "Promptmm: Multi-modal knowledge distillation for recommendation with prompt-tuning," in *Proceedings of the ACM on Web Conference 2024*, 2024, pp. 3217–3228.
- [66] H. Ma, Y. Yang, L. Meng, R. Xie, and X. Meng, "Multimodal conditioned diffusion model for recommendation," in *Companion Proceedings of the ACM on Web Conference 2024*, 2024, pp. 1733–1740.
- [67] Y. Jiang, L. Xia, W. Wei, D. Luo, K. Lin, and C. Huang, "Diffmm: Multi-modal diffusion model for recommendation," in *Proceedings of the 32nd ACM International Conference on Multimedia*, 2024.
- [68] H. Su, J. Li, F. Li, K. Lu, and L. Zhu, "Soil: Contrastive second-order interest learning for multimodal recommendation," in *Proceedings of the 32nd ACM International Conference on Multimedia*, 2024.
- [69] J. Zhang, G. Liu, Q. Liu, S. Wu, and L. Wang, "Modality-balanced learning for multimedia recommendation," in *Proceedings of the 32nd ACM International Conference on Multimedia*, 2024.
- [70] G. Lin, Z. Meng, D. Wang, Q. Long, Y. Zhou, and M. Xiao, "Gume: Graphs and user modalities enhancement for long-tail multimodal recommendation," in *Proceedings of the 33th ACM international conference on information & knowledge management*, 2024.
- [71] X. Dong, X. Song, M. Tian, and L. Hu, "Prompt-based and weak-modality enhanced multimodal recommendation," *Information Fusion*, vol. 101, p. 101989, 2024.
- [72] X. Zhou and C. Miao, "Disentangled graph variational auto-encoder for multimodal recommendation with interpretability," *IEEE Transactions on Multimedia*, 2024.
- [73] J. Yi and Z. Chen, "Variational mixture of stochastic experts auto-encoder for multi-modal recommendation," *IEEE Transactions on Multimedia*, 2024.
- [74] Z. He, Z. Wang, Y. Yang, H. Bai, and L. Wu, "Boosting multimedia recommendation via separate generic and unique awareness," *arXiv preprint arXiv:2406.08270*, 2024.
- [75] J. Xu, Z. Chen, S. Yang, J. Li, H. Wang, and E. C.-H. Ngai, "Mentor: Multi-level self-supervised learning for multimodal recommendation," *arXiv preprint arXiv:2402.19407*, 2024.
- [76] R. K. Ong and A. W. Khong, "Spectrum-based modality representation fusion graph convolutional network for multimodal recommendation," *arXiv preprint arXiv:2412.14978*, 2024.
- [77] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [78] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 1–9.
- [79] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell, "Caffe: Convolutional architecture for fast feature embedding," in *Proceedings of the 22nd ACM international conference on Multimedia*, 2014, pp. 675–678.
- [80] G. Salton and C. Buckley, "Term-weighting approaches in automatic text retrieval," *Information processing & management*, vol. 24, no. 5, pp. 513–523, 1988.
- [81] Y. Koren, "The bellkor solution to the netflix grand prize," *Netflix prize documentation*, vol. 81, no. 2009, pp. 1–10, 2009.
- [82] J. Pennington, R. Socher, and C. D. Manning, "Glove: Global vectors for word representation," in *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, 2014, pp. 1532–1543.
- [83] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean, "Distributed representations of words and phrases and their compositionality," *Advances in neural information processing systems*, vol. 26, 2013.
- [84] S. Arora, Y. Liang, and T. Ma, "A simple but tough-to-beat baseline for sentence embeddings," in *International conference on learning representations*, 2017.
- [85] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," in *International Conference on Learning Representations*, 2017.
- [86] X. He, K. Deng, X. Wang, Y. Li, Y. Zhang, and M. Wang, "Lightgcn: Simplifying and powering graph convolution network for recommendation," in *Proceedings of the 43rd International ACM SIGIR conference on research and development in Information Retrieval*, 2020, pp. 639–648.
- [87] A. v. d. Oord, Y. Li, and O. Vinyals, "Representation learning with contrastive predictive coding," *arXiv preprint arXiv:1807.03748*, 2018.
- [88] R. D. Hjelm, A. Fedorov, S. Lavoie-Marchildon, K. Grewal, P. Bachman, A. Trischler, and Y. Bengio, "Learning deep representations by mutual information estimation and maximization," *arXiv preprint arXiv:1808.06670*, 2018.
- [89] M. Gutmann and A. Hyvärinen, "Noise-contrastive estimation: A new estimation principle for unnormalized statistical models," in *Proceedings of the thirteenth international conference on artificial intelligence and statistics. JMLR Workshop and Conference Proceedings*, 2010, pp. 297–304.
- [90] X. Xia, H. Yin, J. Yu, Q. Wang, L. Cui, and X. Zhang, "Self-supervised hypergraph convolutional networks for session-based recommendation," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 35, no. 5, 2021, pp. 4503–4511.