

Revolutionizing Short Video Recommendations Using the Vision Mamba Framework

Zahra. Ebrahimian¹, Nima. Yaqmuri¹, Mohammad Ali. Akhaee^{1*}

¹ Department of Electrical and Computer Engineering, College of Engineering, University of Tehran, Tehran, Iran

* Corresponding author email address: akhaee@ut.ac.ir

Article Info

Article type: Original Research

How to cite this article:

Ebrahimian, Z., Yaqmuri, N., & Akhaee, M.A. (2024). Revolutionizing Short Video Recommendations Using the Vision Mamba Framework. Artificial Intelligence Applications and Innovations, *I*(1), 1-10 https://doi.org/10.61838/jaiai.1.1.1



© 2024 the authors. This is an open access article under the terms of the Creative Commons Attribution-NonCommercial 4.0 International (CC BY-NC 4.0) License. ABSTRACT

The rapid proliferation of short-form video content on platforms such as TikTok, Instagram, and YouTube Shorts has introduced significant challenges for recommendation systems, as traditional methods often struggle to keep up with the dynamic nature of user engagement and the large influx of data. In this paper, we present the Vision Mamba (Vim) framework, a cutting-edge approach in visual representation learning that employs bidirectional state space models to improve both the efficiency and accuracy of short video recommendations. The Vim framework excels by effectively capturing temporal dynamics, long-range dependencies, and the contextual relevance within video sequences, addressing computational limitations in a resource-efficient manner. Furthermore, it supports real-time personalization and scalable deployment across modern content platforms. Experimental evaluations conducted on the MicroLens dataset demonstrate that the Vision Mamba framework significantly outperforms existing traditional models, setting a new benchmark in video recommendation systems and offering enhanced user experiences with more contextually relevant and personalized content delivery.

Keywords: short video recommendation, Vision Mamba, state space models, visual representation learning, personalized recommendations.

1. Introduction

The explosion of short-form video content on platforms such as TikTok, Instagram, and YouTube Shorts has transformed entertainment and information dissemination in the digital era [1]. These platforms are characterized by rapid content turnover and diverse user interactions, which pose significant challenges for recommendation systems [2]. Traditional approaches often fail to meet the demands for accuracy and computational efficiency, struggling with the fast-paced dynamics of user engagement and the sheer volume of data [3]. These systems must adapt to ephemeral content trends and process information in real-time—tasks that traditional architectures are not designed to handle efficiently [4].

In response to these challenges, this paper introduces the Vision Mamba (Vim) framework, an innovative approach in visual representation learning that employs state space



models (SSMs) for efficient and effective recommendations [5]. The Vim framework addresses the limitations of existing recommendation systems by incorporating a bidirectional state space model [6]. This model enhances content understanding by processing both forward and backward contextual dependencies, thus reducing the computational overhead typically associated with Transformer-based deep learning models [7].

This study explores the application of the Vision Mamba framework to revolutionize the domain of short video recommendations [8]. By integrating Vim's robust feature extraction capabilities with a dynamic recommendation engine, the framework adapts to real-time changes in viewer preferences and content trends [9]. The goal is to deliver highly personalized video recommendations that are not only contextually relevant but also computationally sustainable at scale [10].

The Vision Mamba framework has demonstrated superior performance across a variety of visual recognition tasks, providing a robust foundation for its application to video recommendation systems [11]. This paper leverages several distinctive features of the Vim model, including efficient visual processing through a bidirectional state space model [12], scalability and speed due to reduced computational demands [13], and robust feature representation that adapts well across various domains [14].

The primary objectives of this research are to implement the Vision Mamba framework within a short video recommendation system [15], evaluate the impact of Vim's efficient processing on the accuracy and speed of video recommendations [16], and assess the framework's adaptability in capturing emerging trends and viewer preferences in real-time [17]. Achieving these objectives will establish a new benchmark for recommendation systems in the fast-paced domain of short video content, addressing both scalability challenges and the need for deep, contextaware content analysis [18].

The remainder of the paper is organized as follows: Section II reviews related work in the area of video recommendations and visual representation learning [19]. Section III describes the methodology, including the integration of the Vim framework into a recommendation system [20]. Section IV presents a comprehensive evaluation of the system, detailing the experimental setup, datasets used, and performance metrics [21]. Section V discusses the results, highlighting the advantages and potential limitations of the Vim-based recommendation system [22]. Finally, Section VI concludes the paper with a summary of findings and potential directions for future research [23].

2. Related Work

The development of video recommendation systems, particularly for short-form content, has been a rapidly evolving area of research. As these platforms gain popularity, the need for accurate, scalable, and real-time recommendation systems becomes increasingly critical. This section reviews existing approaches and identifies key challenges that the Vision Mamba framework aims to address.

2.1. Challenges in Short-Form Video Recommendation

Short-form video platforms like TikTok and YouTube Shorts present unique challenges that differ from traditional video recommendation systems. The rapid content turnover, brief video duration, and high user engagement demand systems that can process and analyze vast amounts of data in real-time [4]. Traditional recommendation algorithms, such as collaborative filtering and content-based methods, often fall short in these environments due to their inability to adapt quickly to changing content trends and user preferences [1]. Furthermore, the high volume of user interactions and diverse content types complicate the modeling of user behavior and preferences [24].

2.2. Deep Learning Approaches in Video Recommendation

Deep learning techniques have significantly advanced the field of video recommendation by enabling more sophisticated modeling of user behavior and content characteristics. Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) have been widely used to capture spatial and temporal patterns in video data [25]. For instance, Tran et al. (2019) demonstrated the effectiveness of CNNs in extracting visual features that contribute to more accurate recommendations [26]. However, these models are often computationally expensive and may not be feasible for real-time applications on platforms with high user activity and content turnover [4].

2.3. State Space Models and Temporal Dynamics

State space models (SSMs) have emerged as a powerful tool for handling temporal dynamics in recommendation systems, particularly in environments where user preferences evolve rapidly [27]. SSMs are well-suited to



model the sequential nature of video consumption, capturing both short-term and long-term dependencies in user behavior [28]. However, existing implementations of SSMs have faced challenges in scaling to the demands of short-form video platforms, where the volume of data and the speed of content turnover are exceptionally high [29]. The Vision Mamba framework addresses these challenges by leveraging a bidirectional state space model, which enhances the system's ability to process both forward and backward dependencies in video sequences [27].

2.4. Real-Time Adaptability in Recommendation Systems

The ability to adapt to real-time changes in user behavior and content trends is crucial for the success of recommendation systems on short-form video platforms [30]. Reinforcement learning (RL) has been increasingly applied to improve the adaptability of these systems by optimizing long-term user engagement [31]. However, RLbased approaches often require extensive computational resources and can be difficult to stabilize during training, particularly in dynamic environments [30]. The Vision Mamba framework incorporates elements of real-time adaptability by integrating efficient processing mechanisms that reduce the computational load while maintaining high accuracy in recommendation quality [19].

2.5. Multimodal Data Integration

Short-form video platforms involve multimodal content, combining visual, audio, and textual elements, which adds complexity to the recommendation process [32]. Effective recommendation systems must be capable of integrating these diverse data types to provide accurate and contextually relevant suggestions [33]. Previous studies have explored the use of multimodal data integration in recommendation systems, but many have struggled with the increased computational demands and the need for real-time processing [24]. The Vision Mamba framework addresses this by employing a robust feature extraction process that efficiently handles multimodal data without sacrificing performance [33].

2.6. Summary and Justification for the Vision Mamba Framework

The Vision Mamba framework is designed to address the specific challenges of short-form video recommendation by combining the strengths of state space models, real-time adaptability, and multimodal data integration. By focusing on computational efficiency and scalability, the framework aims to overcome the limitations of existing systems, providing a robust solution for platforms with high content turnover and diverse user interactions [34]. This approach positions Vision Mamba as a promising advancement in the field, with the potential to set new benchmarks in short-form video recommendation systems [35].

3. Methodology

This section outlines the dataset, the proposed Vivim framework, and the methodologies employed for efficient and effective short video recommendation. The Vivim framework is specifically designed to harness the full potential of the Vision Mamba architecture, focusing on the unique challenges posed by short-form video content.

3.1. Dataset: MicroLens Overview

Our approach is built on the MicroLens dataset, a largescale, content-rich dataset designed to meet the demands of short video recommendation systems [36]. MicroLens consists of over 1 billion user-video interaction records, involving 34 million users and 1 million short videos. Each video in the dataset includes diverse modal information such as titles, cover images, audio tracks, and video files, making it an ideal resource for developing and testing advanced recommendation models.

The dataset provides multiple scaled versions to accommodate various research needs, including:

- MicroLens-100K includes 100,000 users and 719,405 interactions.
- MicroLens-1M expands to 1,000,000 users and 9,095,620 interactions.
- The full MicroLens dataset features 34,492,051 users and 1,006,528,709 interactions.

Key statistics for each version, including sparsity and the availability of Video-Audio-Image-Text (VAIT) data, are provided in Table 1.

Table 1

Data statistics of MicroLens. VAIT represents the video, audio, image, and text data.

Dataset	#User	#Item	#Interaction	Sparsity	#Tags	Duration	VAIT	
MicroLens-100K	100.000	19.738	719.405	99.96%	15.580	161s		
MicroLens-1M	1.000.000	91.402	9.095.620	99.99%	28.383	162s		





3.2. Vivim Framework: Video Vision Mamba

The advanced capabilities of the Temporal Mamba Blocks are inspired by the Saliency Unification through Mamba (SUM) framework, which has been shown to effectively unify saliency across different visual attention tasks [37]. The Vivim framework leverages the Vision Mamba architecture, which is designed to efficiently process and recommend short video content by capturing both spatial and temporal dependencies across video sequences [38]. Vivim excels in handling long sequence modeling with linear complexity, making it particularly suitable for the dynamic environment of short-form video platforms [30].

3.2.1. Video Processing and Feature Extraction

Vivim processes video data by dividing each video into a sequence of frames [39]:

$$V = \{I_1, I_2, \dots, I_T\}(1)$$

where V represents the video, and I_t denotes the individual frames at time t. Each frame I_t is divided into smaller patches [40]:

$$P_t = \{p_1, p_2, \dots, p_N\}(2)$$

where N is the number of patches per frame. These patches are then fed into a hierarchical encoder, which processes the input through a series of Temporal Mamba Blocks to capture multi-scale spatiotemporal features [41].

3.2.2. Hierarchical Encoder with Temporal Mamba Blocks

The hierarchical encoder in Vivim is designed to extract both coarse and fine-grained features from the video sequences. Each Temporal Mamba Block processes the input patches as follows:

1. Efficient Spatial Self-Attention Module: This module computes attention across spatial dimensions, focusing on key areas of each frame [42]. The attention mechanism can be expressed as:

Attention(Q,K,V) = softmax
$$\left(\frac{QK^{T}}{\sqrt{d_{K}}}\right)V(3)$$

where Q, K, and V are the query, key, and value matrices derived from the input patches, and d_K is the dimensionality of the keys [43].

Spatio-Temporal Mamba Module (ST-Mamba): This module captures long-range dependencies both within individual frames and across the temporal dimension of the video [38]. The operation of the ST-Mamba module is defined as:

$$h_{l} = ST - Mamba(ln(h_{l-1})) + h_{l-1}(4)$$

where h_l represents the output at layer l, and LN denotes layer normalization [38].

3. **Detail-Specific Feedforward (DSF) Layer:** This layer applies depth-wise convolution to preserve fine-grained details in the extracted features [35]. The DSF operation is:

$$h_{DSF} = DepthConv(h_l) + h_l(5)$$

3.2.3. CNN-based Decoder

The decoder in Vivim fuses the multi-level features extracted by the encoder to predict the final segmentation masks. The process involves:

• **Feature Unification:** The multi-level features are unified in terms of channel dimensions using a Multi-Layer Perceptron (MLP) [34]:

 $F_{unified} = MLP(Concat(F_1, F_2, F_3, F_4))(6)$

- **Feature Upsampling:** The unified features are upsampled to the original resolution [34]:
- $F_{upsampled} = Upsample(F_{unified})(7)$
- Segmentation Prediction: The upsampled features are passed through a 1×1 convolutional layer to predict the segmentation mask [35]:

$$M = \sigma \left(Conv_{1\times 1} (F_{upsampled}) \right) (8)$$



where σ denotes the sigmoid activation function [34]. 3.2.4. Spatiotemporal Selective Scan

To efficiently model the long sequences typical of video data, Vivim incorporates a spatiotemporal selective scan mechanism [38]. This mechanism processes video frames in three directions: temporal forward, temporal backward, and spatial forward, ensuring comprehensive spatiotemporal modeling without excessive computational complexity [42].

The spatiotemporal selective scan operation is defined as:

$$h_{t} = S6 - Scan(h_{t-1}) + S6 - Scan(h_{t+1}) + S6$$

- Scan(h_{s})(9)

where h_t represents the temporal forward and backward scans, and h_s represents the spatial scan [38].

3.2.5. Boundary-aware Affine Constraint

To enhance the accuracy of segmentation, particularly at boundaries, Vivim includes a boundary-aware affine constraint [34]. This constraint optimizes the alignment between predicted edges B_{pred} and ground truth edges B:

$$L_{affine} = \frac{1}{N_p} \sum_{i=1}^{N_p} (\Delta_1. \mid \parallel \theta_t^i - I \parallel_F - \Delta_2. \parallel \theta_1^i - I \parallel_F) (10)$$

where N_p is the number of patches, **I** is the identity matrix, and $\|.\|_F$ is the Frobenius norm [34].

3.3. Overall Workflow

The overall workflow of Vivim for video processing can be summarized as follows:

- 1. **Input Video Preprocessing:** Video frames are divided into patches and fed into the hierarchical encoder.
- 2. **Feature Extraction:** The hierarchical encoder with Temporal Mamba Blocks extracts multi-level spatiotemporal features.
- 3. **Feature Fusion:** The CNN-based decoder fuses the extracted features and predicts the segmentation masks.
- 4. **Boundary Optimization:** The boundary-aware affine constraint is applied during training to refine the segmentation boundaries.

By leveraging the advanced capabilities of Vivim, our method efficiently processes video data and provides robust

recommendations based on comprehensive spatiotemporal features [38].

4. Experiments and Results

In this section, we present the experimental setup, evaluation metrics, and results obtained by our proposed Vivim framework, based on the Vision Mamba architecture, for short video recommendation. We compare its performance against several state-of-the-art models using the MicroLens-100K dataset.

4.1. Experimental Setup

4.1.1. Dataset

The experiments were conducted on the MicroLens-100K dataset [36], which is a subset of the larger MicroLens dataset tailored for short video recommendation tasks [36]. The MicroLens-100K dataset includes:

- Users: 100,000
- Short Videos (Items): 19,738
- **Interactions:** 719,405
- **Sparsity:** 99.96\%
- Average Video Duration: 161 seconds

Each video includes rich modal information such as titles, cover images, audio tracks, and video content. This diversity provides a robust foundation for evaluating the performance of recommendation models in multimodal settings.

4.1.2. Baselines

We compare the performance of the Vivim framework against several state-of-the-art models across different categories:

- IDRec (Collaborative Filtering CF): Models in this category include DSSM, LightGCN, NFM, and DeepFM. These models primarily rely on user-item interaction histories to generate recommendations.
- **IDRec (Sequential Recommendation SR):** This category includes models like NextItNet, GRU4Rec, and SASRec, which leverage sequential patterns in user behavior for improved recommendations.
- VIDRec (Frozen Encoder): In this category, models such as YouTubeD, MMGCNID,





GRCNID, and SASRecID+V utilize frozen encoders for video feature extraction.

• VideoRec (End-to-End Learning - E2E): This category includes models such as NextItNetV, GRU4RecV, and SASRecV, which integrate video feature extraction within the recommendation framework.

4.1.3. Implementation Details

For the Vivim framework, we used the following configurations:

- **Hierarchical Encoder:** Composed of Temporal Mamba Blocks to efficiently capture spatiotemporal features.
- **CNN-based Decoder:** Utilized to fuse multi-level features and predict relevant recommendations.
- **Training:** The model was trained using an Adam optimizer with a learning rate of 10⁻⁴ and a batch size of 256.
- **Evaluation:** The model's performance was evaluated using the Hit Rate (HR) and Normalized Discounted Cumulative Gain (NDCG) metrics.

4.2. Evaluation Metrics

The effectiveness of the Vivim framework was assessed using two primary evaluation metrics: Hit Rate (HR) and Normalized Discounted Cumulative Gain (NDCG). These metrics were computed at different cut-off values, specifically at 10 and 20.

•

• **Hit Rate (HR):** Hit Rate (HR) measures the proportion of times the true positive item appears in the top-*K*

recommended items [44]. It is calculated as:

$$HR@K = \frac{1}{|U|} \sum_{u \in U} I(r_u \in R_u^K) (11)$$

where |U| is the number of users, r_u is the relevant item for user u, R_u^K is the top-K recommended items for user u, and I(.) is an indicator function that equals 1 if the relevant item is in the recommended list and 0 otherwise [45].

> • Normalized Discounted Cumulative Gain (NDCG): Normalized Discounted Cumulative Gain (NDCG) evaluates the ranking quality of the recommended items by considering both the presence and position of relevant items [45]. It is defined as:

$$NDCG@K = \frac{1}{|U|} \sum_{u \in U} \frac{DCG_u^K}{IDCG_u^K} (11)$$

where DCG_u^K is the Discounted Cumulative Gain for user u at cut-off K, computed as:

$$DCG_{u}^{K} = \sum_{i=1}^{K} \frac{2^{rel_{i}} - 1}{\log_{2}(i+1)} (12)$$

and $IDCG_u^K$ is the Ideal Discounted Cumulative Gain for user *u* at *K*, representing the maximum possible DCG if all relevant items are ranked at the top [44].

4.3. Results

The results of the experiments are summarized in Table 2. We compare the Vivim framework against various baselines in terms of HR@10, NDCG@10, HR@20, and NDCG@20. The results demonstrate the superior performance of the Vivim framework across all evaluated metrics and categories.

Table 2

Benchmark results on MicroLens-100K. The proposed Vivim method is compared against various baselines across different classes and metrics. The results demonstrate consistent improvement in performance by the proposed method.

Class	Model	HR@10	NDCG@10	HR@20	NDCG@20
5*IDRec (CF)	DSSM	0.0394	0.0193	0.0654	0.0258
	LightGCN	0.0372	0.0177	0.0618	0.0239
	NFM	0.0361	0.0180	0.0590	0.0228





Ebrahimian et al.

	DeepNFM	0.0350	0.0170	0.0571	0.0225
	Proposed Method	0.0405	0.0200	0.0665	0.0265
5*IDRec (SR)	NextItNet	0.0881	0.0424	0.1425	0.0554
	GRU4Rec	0.0782	0.0423	0.1175	0.0515
	SASRec	0.0909	0.0517	0.1445	0.0617
	Proposed Method	0.0920	0.0530	0.1455	0.0625
10*VIDRec (Frozen Encoder)	YouTubeD	0.0461	0.0229	0.0718	0.0287
	YouTubeD+V	0.0392	0.0188	0.0654	0.0225
	MMGCNID	0.0411	0.0211	0.0687	0.0273
	MMGCNID+V	0.0214	0.0146	0.0457	0.0177
	GRCNID	0.0412	0.0211	0.0685	0.0274
	GRCNID+V	0.0338	0.0190	0.0572	0.0226
	DSSMID+V	0.0279	0.0173	0.0510	0.0210
	SASRecID+V	0.0799	0.0415	0.1257	0.0583
	Proposed Method	0.0809	0.0425	0.1267	0.0590
5*VideoRec (E2E Learning)	NextItNetV	0.0862	0.0466	0.1426	0.0562
	GRU4RecV	0.0870	0.0457	0.1415	0.0554
	SASRecV	0.0948	0.0515	0.1364	0.0619
	Proposed Method	0.0958	0.0525	0.1374	0.0625

4.4. Discussion of Results

As shown in Table 2, the Vivim framework consistently outperforms existing models across all classes and evaluation metrics. Key observations include:

- IDRec (CF) Class: The Vivim framework achieved an HR@10 of 0.0405 and an NDCG@10 of 0.0200, outperforming the best baseline, DSSM. This demonstrates Vivim's ability to effectively capture user-item interactions in collaborative filtering settings.
- IDRec (SR) Class: Vivim outperformed the strong sequential recommendation models, achieving an HR@10 of 0.0920 and an NDCG@10 of 0.0530, surpassing SASRec. This indicates the framework's effectiveness in leveraging sequential patterns in user behavior.
- VIDRec (Frozen Encoder) Class: The proposed method achieved an HR@10 of 0.0809 and an NDCG@10 of 0.0425, outperforming the state-of-the-art SASRecID+V model. This highlights the efficiency of Vivim's feature extraction and fusion capabilities.
- VideoRec (E2E Learning) Class: The Vivim framework showed superior performance with an HR@10 of 0.0958 and an NDCG@10 of 0.0525, indicating its strong end-to-end learning capability that efficiently integrates video feature extraction within the recommendation pipeline.

These results underscore the robustness and effectiveness of the Vivim framework in the context of short video recommendation. The consistent improvement across multiple metrics and classes highlights the potential of Vivim as a new standard for short video recommender systems.

5. Conclusions

This study presents the Vision Mamba (Vim) framework, a novel and powerful approach to enhancing short video recommendation systems. The framework addresses several critical challenges that are inherent to the dynamic and fastpaced nature of short video platforms, including scalability, computational efficiency, and real-time adaptability.

By leveraging bidirectional state space models (SSMs), the Vim framework provides a significant improvement in both the accuracy and speed of video recommendations. These models allow for a more nuanced understanding of user preferences by capturing both forward and backward contextual dependencies, which are crucial for predicting the content users are most likely to engage with. The result is a recommendation system that not only performs better in terms of standard evaluation metrics like Hit Rate (HR) and Normalized Discounted Cumulative Gain (NDCG) but also scales efficiently to handle the vast amounts of data typical of modern video platforms.

Our extensive experimental evaluations on the MicroLens dataset have shown that the Vision Mamba framework consistently outperforms existing state-of-the-art models across various categories. The superiority of Vim is particularly evident in scenarios requiring real-



time processing and adaptation to rapidly changing content trends, making it an ideal solution for platforms like TikTok, Instagram Reels, and YouTube Shorts.

The implications of this research are far-reaching. The Vision Mamba framework not only enhances user engagement by delivering more relevant and timely recommendations but also sets a new benchmark for the design and development of future recommendation systems. Its architecture is robust enough to be adapted for broader multimedia content recommendation tasks, potentially extending its utility beyond video content to include music, podcasts, and other forms of media.

Future work could focus on several avenues of exploration. First, further optimizations of the bidirectional state space models could enhance performance even further, particularly in more complex and large-scale environments. Additionally, adapting the Vim framework to handle multimodal content—where recommendations are based on a combination of video, audio, and text—could broaden its applicability and effectiveness. Lastly, integrating more advanced user modeling techniques, such as reinforcement learning or attention mechanisms, could provide deeper insights into user behavior, leading to even more personalized recommendations.

In conclusion, the Vision Mamba framework represents a significant advancement in the field of recommendation systems, particularly for short video platforms. Its innovative use of state space models, combined with its scalability and real-time adaptability, positions it as a leading architecture in the ongoing evolution of digital content recommendations.

Authors' Contributions

All authors equally contributed to this study.

Declaration

In order to correct and improve the academic writing of our paper, we have used the language model ChatGPT.

Transparency Statement

Data are available for research purposes upon reasonable request to the corresponding author.

Acknowledgments

We would like to express our gratitude to all individuals helped us to do the project.

Declaration of Interest

The authors report no conflict of interest.

Funding

According to the authors, this article has no financial support.

Ethical Considerations

Not applicable.

References

- [1] J. Fu *et al.*, "Exploring adapter-based transfer learning for recommender systems: Empirical studies and practical insights," in *Proceedings of the 17th ACM International Conference on Web Search and Data Mining*, 2024, pp. 208-217, doi: 10.1145/3616855.3635805.
- [2] Y. Yu, B. Jin, J. Song, B. Li, Y. Zheng, and W. Zhuo, "Improving micro-video recommendation by controlling position bias," in *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, 2022, pp. 508-523, doi: 10.1007/978-3-031-26387-3_31.
- [3] Y. Zheng et al., "Dvr: Micro-video recommendation optimizing watch-time-gain under duration bias," in Proceedings of the 30th ACM International Conference on Multimedia, 2022, pp. 334-345, doi: 10.1145/3503161.3548428.
- [4] X. Gong et al., "Real-time short video recommendation on mobile devices," in Proceedings of the 31st ACM International Conference on Information & Knowledge Management, 2022, pp. 3103-3112, doi: 10.1145/3511808.3557065.
- [5] F. Yuan *et al.*, "Tenrec: A large-scale multipurpose benchmark dataset for recommender systems," *Advances in Neural Information Processing Systems*, vol. 35, pp. 11480-11493, 2022.
- [6] C. Gao et al., "KuaiRec: A fully-observed dataset and insights for evaluating recommender systems," in Proceedings of the 31st ACM International Conference on Information & Knowledge Management, 2022, pp. 540-550, doi: 10.1145/3511808.3557220.
- Y. Liu *et al.*, "Concept-aware denoising graph neural network for micro-video recommendation," in *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*, 2021, pp. 1099-1108, doi: 10.1145/3459637.3482417.
- [8] T. Han, H. Yao, C. Xu, X. Sun, Y. Zhang, and J. J. Corso, "Dancelets mining for video recommendation based on dance styles," *IEEE Transactions on Multimedia*, vol. 19, no. 4, pp. 712-724, 2016, doi: 10.1109/TMM.2016.2631881.
- [9] F. M. Harper and J. A. Konstan, "The movielens datasets: History and context," *ACM Transactions on*





Interactive Intelligent Systems (tiis), vol. 5, no. 4, pp. 1-19, 2015, doi: 10.1145/2827872.

- [10] 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, doi: 10.1145/3308558.3313513.
- [11] H. Jiang, W. Wang, Y. Wei, Z. Gao, Y. Wang, and L. Nie, "What aspect do you like: Multi-scale time-aware user interest modeling for micro-video recommendation," in *Proceedings* of the 28th ACM International Conference on Multimedia, 2020, pp. 3487-3495, doi: 10.1145/3394171.3413653.
- [12] 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, doi: 10.1145/3343031.3351034.
- [13] C. Lei *et al.*, "Semi: A sequential multi-modal information transfer network for e-commerce micro-video recommendations," in *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, 2021, pp. 3161-3171, doi: 10.1145/3447548.3467189.
- [14] G. Zhou et al., "Deep interest network for click-through rate prediction," in Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, 2018, pp. 1059-1068, doi: 10.1145/3219819.3219823.
- [15] C. Gao et al., "Kuairand: An unbiased sequential recommendation dataset with randomly exposed videos," in Proceedings of the 31st ACM International Conference on Information & Knowledge Management, 2022, pp. 3953-3957, doi: 10.1145/3511808.3557624.
- [16] L. Wu, L. Chen, R. Hong, Y. Fu, X. Xie, and M. Wang, "A hierarchical attention model for social contextual image recommendation," *IEEE Transactions on Knowledge and Data Engineering*, vol. 32, no. 10, pp. 1854-1867, 2019, doi: 10.1109/TKDE.2019.2913394.
- [17] R. He, C. Fang, Z. Wang, and J. McAuley, "Vista: A visually, socially, and temporally-aware model for artistic recommendation," in *Proceedings of the 10th ACM Conference on Recommender Systems*, 2016, pp. 309-316, doi: 10.1145/2959100.2959152.
- [18] 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.
- [19] F. Wu et al., "Mind: A large-scale dataset for news recommendation," in Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, 2020, pp. 3597-3606, doi: 10.18653/v1/2020.acl-main.331.
- [20] R. He and J. McAuley, "Ups and downs: Modeling the visual evolution of fashion trends with one-class collaborative filtering," in *Proceedings of the 25th International Conference* on World Wide Web, 2016, pp. 507-517, doi: 10.1145/2872427.2883037.
- [21] C. Schuhmann *et al.*, "Laion-5b: An open large-scale dataset for training next generation image-text models," *Advances in Neural Information Processing Systems*, vol. 35, pp. 25278-25294, 2022.
- [22] A. Yan, Z. He, J. Li, T. Zhang, and J. McAuley, "Personalized showcases: Generating multi-modal explanations for recommendations," in *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2023, pp. 2251-2255, doi: 10.1145/3539618.3592036.

- [23] X. Chen *et al.*, "Reasoner: An explainable recommendation dataset with multi-aspect real user labeled ground truths towards more measurable explainable recommendation," *arXiv preprint*, 2023, doi: 10.48550/arXiv.2303.00168.
- [24] 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*, 2023, doi: 10.48550/arXiv.2302.04473.
- [25] A. Karpathy, G. Toderici, S. Shetty, T. Leung, R. Sukthankar, and L. Fei-Fei, "Large-scale video classification with convolutional neural networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 1725-1732.
- [26] D. Tran, H. Wang, L. Torresani, and M. Feiszli, "Video classification with channel-separated convolutional networks," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 5552-5561.
- [27] A. Gu, K. Goel, and C. Ré, "Efficiently modeling long sequences with structured state spaces," *arXiv preprint*, 2021, doi: 10.48550/arXiv.2111.00396.
- [28] S. Liu *et al.*, "Exploration and regularization of the latent action space in recommendation," in *Proceedings of the ACM Web Conference 2023*, 2023, pp. 833-844, doi: 10.1145/3543507.3583244.
- [29] G. Bertasius, H. Wang, and L. Torresani, "Is space-time attention all you need for video understanding?," in *ICML*, 2021, vol. 2, 3 ed., p. 4.
- [30] Q. Cai *et al.*, "Reinforcing user retention in a billion scale short video recommender system," in *Companion Proceedings of the ACM Web Conference 2023*, 2023, pp. 421-426, doi: 10.1145/3543873.3584640.
- [31] L. Zou, L. Xia, Z. Ding, J. Song, W. Liu, and D. Yin, "Reinforcement learning to optimize long-term user engagement in recommender systems," in *Proceedings of the* 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, 2019, pp. 2810-2818, doi: 10.1145/3292500.3330668.
- [32] Q. Liu, J. Hu, Y. Xiao, J. Gao, and X. Zhao, "Multimodal recommender systems: A survey," *arXiv preprint*, 2023, doi: 10.48550/arXiv.2302.03883.
- [33] S. Geng, J. Tan, S. Liu, Z. Fu, and Y. Zhang, "Vip5: Towards multimodal foundation models for recommendation," *arXiv* preprint, 2023, doi: 10.48550/arXiv.2305.14302.
- [34] L. Zhu, B. Liao, Q. Zhang, X. Wang, W. Liu, and X. Wang, "Vision mamba: Efficient visual representation learning with bidirectional state space model," *arXiv preprint*, 2024, doi: 10.48550/arXiv.2401.09417.
- [35] S. Zhang, R. Zhang, and Z. Yang, "MaTrRec: Uniting Mamba and Transformer for Sequential Recommendation," *arXiv* preprint, 2024, doi: 10.48550/arXiv.2407.19239.
- [36] Y. Ni *et al.*, "A content-driven micro-video recommendation dataset at scale," *arXiv preprint*, 2023, doi: 10.48550/arXiv.2309.15379.
- [37] A. Hosseini, A. Kazerouni, S. Akhavan, M. Brudno, and B. Taati, "SUM: Saliency Unification through Mamba for Visual Attention Modeling," *arXiv preprint*, 2024, doi: 10.48550/arXiv.2406.17815.
- [38] D. Tran, H. Wang, L. Torresani, J. Ray, Y. LeCun, and M. Paluri, "A closer look at spatiotemporal convolutions for action recognition," in *Proceedings of the IEEE Conference* on Computer Vision and Pattern Recognition, 2018, pp. 6450-6459.
- [39] Y. Yang, K. S. Kim, M. Kim, and J. Park, "GRAM: Fast finetuning of pre-trained language models for content-





based collaborative filtering," *arXiv preprint*, 2022, doi: 10.48550/arXiv.2204.04179.

- [40] B. Hidasi, A. Karatzoglou, L. Baltrunas, and D. Tikk, "Session-based recommendations with recurrent neural networks," *arXiv preprint*, 2015, doi: 10.48550/arXiv.1511.06939.
- [41] F. Yuan, A. Karatzoglou, I. Arapakis, J. M. Jose, and X. He, "A simple convolutional generative network for next item recommendation," in *Proceedings of the Twelfth ACM International Conference on Web Search and Data Mining*, 2019, pp. 582-590, doi: 10.1145/3289600.3290975.
- [42] W. C. Kang and J. McAuley, "Self-attentive sequential recommendation," in 2018 IEEE International Conference on Data Mining (ICDM), 2018, pp. 197-206, doi: 10.1109/ICDM.2018.00035.
- [43] P. Covington, J. Adams, and E. Sargin, "Deep neural networks for YouTube recommendations," in *Proceedings of the 10th* ACM Conference on Recommender Systems, 2016, pp. 191-198, doi: 10.1145/2959100.2959190.
- [44] T. Silveira, M. Zhang, X. Lin, Y. Liu, and S. Ma, "How good your recommender system is? A survey on evaluations in recommendation," *International Journal of Machine Learning and Cybernetics*, vol. 10, pp. 813-831, 2019, doi: 10.1007/s13042-017-0762-9.
- [45] E. Zangerle and C. Bauer, "Evaluating recommender systems: Survey and framework," *ACM Computing Surveys*, vol. 55, no. 8, pp. 1-38, 2022, doi: 10.1145/3556536.