

Advancing Fairness in Multimodal Machine Learning for Internet-Scale Video Data: Comprehensive Bias Mitigation and Evaluation Framework

Syed Imtiazul Sami

Department of Electrical and
Electronics Engineering,
Ahsanullah University of Science
and Technology, Dhaka

Mohammad Rasel Mahmud

Department of Computer Science
and Engineering (CSE),
American International University
Bangladesh, Dhaka,

Khaled Bin Showkot Tanim

Department of Electrical and
Computer Engineering (ECE),
North South University, Dhaka,
Bangladesh

Mohammad Zubair Hussain

Department of Electronics and
Telecommunications Engineering,
North South University, Dhaka, Bangladesh

Mahpara Khan Orpa

Department of Computer Science and Engineering
University name: Leading University

ABSTRACT

The progression of multimodal machine learning (MML), a pivotal aspect of the artificial intelligence (AI) revolution, greatly enhances the analysis and comprehension of video data by providing insights across several modalities, including text, audio, and visual formats. MML models have extensive applications in entertainment, healthcare, and autonomous systems; nevertheless, when trained on expansive video datasets that encompass a diverse cultural, ethnic, and linguistic spectrum, they encounter significant challenges related to fairness and prejudice. This work presents a comprehensive investigation of bias reduction and fairness in MML, addressing issues arising from the intricacies of large-scale video data. This study (1) identifies the origins and mechanisms of bias in MML systems, (2) introduces advanced methodologies to enhance model fairness, and (3) offers a thorough framework for assessing fairness in large-scale video datasets. We offer a framework that integrates bias-aware pre-processing, fairness-aware modeling across multimodal settings, and scalable assessment metrics. Specifically, we employ balanced sampling, GANs for synthetic data augmentation, and adversarial debiasing to provide equitable representation and prediction. We validate our methods on extensive benchmark datasets (YouTube-8M, Activity Net, and Ego4D), demonstrating substantial enhancements in performance and fairness. The experimental findings indicate that the multimodal model surpasses both the modal models and the state-of-the-art techniques, achieving an accuracy of 90.5% and an F1 score of 91.0%. Ultimately, we enhanced fairness measurements, specifically differential impact and equalized chances, by 32.3% and 17.9%, respectively, demonstrating the efficacy of our bias mitigation strategies. However, comparative assessments reveal that our technique delivers state-of-the-art performance on the trade-off between predictive accuracy and ethical fairness, making it a feasible option for real-life contexts where equity is critical. Qualitative study corroborates the alleviation of demographic bias in model predictions, especially on sensitive tasks such as emotion detection and demographic classification. Our research enhances the ethical use of MML systems, guaranteeing that these models are resilient and equitable among diverse population subgroups while establishing a foundation for future

advancements in multimodal fusion methodologies and task-specific fairness metrics.

General Terms

Machine Learning, Data Science, Data Analytics, Decision Support Systems, Predictive Modeling, Algorithms, Big Data, Artificial Intelligence, Retail Analytics, Pattern Recognition, Optimization, Business Intelligence, Multi-Channel Retail, Information Systems, Computational Methods.

Keywords

Multimodal Machine Learning, Fairness, Bias Mitigation, Video Datasets, Generative Adversarial Networks, Adversarial Debiasing, Ethical AI.

1. INTRODUCTION

The multimodal machine learning (MML) has become a prominent area in AI, playing a key role for machines to learn from visual and audiovisual data, notably the video data. MML enables multi-modal data fusion, such as text, audio, and video, offering rich context and additional information when dealing with more difficult tasks like sentiment analysis or even video captioning (Al-Zoghby et al., 2024). This capability has led to broad applicability in numerous fields, including entertainment, medicine, and autonomous systems (Barua et al., 2023). MML models hold great potential but yet face substantial issues with respect to fairness and prejudice, especially when used with video data at internet scales that comprise a combination of cultural, demographic, and linguistic components (Ren et al., 2022). Such challenges demand careful consideration of fairness frameworks and mitigating measures.

Since MML models are heavily dependent on the learned data, which often embeds societal prejudices, it has prompted significant concerns. They could occur in terms of the lack of representation for specific demographic groups or the reinforcing of negative and destructive stereotypes (Ouenniche, 2023). Some cases, such as audiovisual production, where there is a risk of ethical issues if the performance of models is biased, have already raised the alarm of the necessity to treat fairness (Haouhat et al., 2023). Furthermore, biased MML systems applied in high-stakes sectors like recruitment or content recommendation might exacerbate inequality and have

significant social and economic implications (Zhang, 2022). It thus becomes necessary to overcome these difficulties before we are actually constructing MML systems that are truly successful, egalitarian, and inclusive.

Since MML commonly incorporates internet-scale video datasets, maintaining fairness in MML is especially critical given the knowledge that such datasets are massive and non-trivial, compounding the potential for biased data. Being collected from multiple sources and user-generated content, these datasets naturally exhibit some aspects of real-world data, such as imbalances and biases (Caton & Haas, 2024). Ultimately, bias can be introduced if models are trained on biased datasets, especially when datasets are not evenly representative of a demographic, producing algorithms that ultimately favor or penalize specific demographics (Jui & Rivas, 2024). These biased models challenge the MML applications trustworthiness and satisfying MML ethical AI standards. Moreover, there is no proper framework for evaluating fairness in model-based machine learning (MML), and that further compounds the situation because lists of prior work lack defined definitions of bias (Oluwaseyi, 2024).

We intend to provide a timely and complete examination into the problems of bias reduction and fairness in MML, specifically in the setting of internet-scale video datasets. The key contributions of this study include evaluating in depth the sources and types of bias in MML systems, leading to the development of approaches to promote fairness in these models, and building a powerful evaluation framework to quantify fairness in large-scale video datasets. This study seeks to contribute to the developing discourses around ethical AI and to allow less biased MML systems (Kheya et al., 2024) by tackling the aforementioned objectives systematically.

This paper is constructed as follows: Related work is detailed in the next section, including, to the best of our knowledge, the examined attempts, their achievements, and their shortcomings in bias mitigation and fairness in MML. Next, we describe the issues of internet-scale video datasets and their repercussions on fairness. Then the mixed solutions are offered, including the bias-aware data preparation, the fairness-aware model development, and a scalable evaluation framework. The methods part discusses the experimental setup and datasets used, while the results and discussion section provide the empirical results and their consequences. The work finishes with contributions, implications for future research, and directions for future research.

2. LITERATURE REVIEW

The bias in machine learning, notably through traditional and deep learning models, has been the topic of significant research to understand its influence along with ways for mitigation. There are various factors leading to bias, such as under-representation or over-representation of the training dataset (Mavrogiorgos et al., 2024; Balayn, Lofi, & Houben, 2021), erroneous labeling, and distribution discrepancies where datasets come from (Mavrogiorgos et al., 2024). This work on translating bifurcation in the biometrics realm demonstrates how imbalanced training data leads to biased outcomes, notably demographic bias (Drozdzowski et al., 2020). The researchers have raised concerns about ethical concerns and studied how strategies for bias prevention could be utilized to develop egalitarian machine learning systems (Prasad, 2024). More recently, there has been some work on the designing and assessment of systems that seek to mitigate gender-related prejudice, bringing fresh insights on the mitigation of bias in

context, such as in applications tailored for specific domains, e.g., music (Shrestha, 2023).

The study of fairness in multimodal learning (MML) has become a major research field in itself, including on video data. Previous publications have noted issues of representation bias and injustice in multimodal AI systems and underlined the necessity for fairness-aware approaches in order to increase model performance and ethical alignment [1]. Even more targeted, researchers have extended this occupational stereotyping to image retrieval systems, indicating that biased data and algorithms have a broader social consequence (Dash, 2023). With respect to fairness in MML, researchers have developed strategies for fairness-aware representation learning, where the attribute-class correlations are suppressed in order to eliminate bias (Sarridis et al. 2024). Moreover, ongoing research has highlighted some of the ethical difficulties with multimodal datasets (Birhane, Prabhu, & Kahembwe, 2021), such as misogyny, pornographic, and malignant stereotyping existing within large-scale datasets.

Video datasets at internet scale come with their unique issues, such as labeling bias, sampling bias, and computational considerations. Such sizeable datasets are frequently not well-distributed and might thus reinforce bias, hurting the modeling result (Chen et al., 2024). Various research studies have suggested diversity in dataset generation and usage of fairness-aware algorithms as possible strategies to decrease bias in video data (Zhao, Chen, & Thuraisingham, 2021). Additionally, it has been discovered that vision models trained on unfiltered images with supervisionless can be more robust and fair, suggesting a promising method for eliminating bias in video datasets using unsupervised learning (Goyal et al. 2022).

Bias mitigation difficulties can be roughly divided into data representative or algorithmic. As observed in Hosseini et al. (2023), representation and labeling bias can resound in video datasets with a skewed distribution, because particularly unfavorable outcomes may stem from the bias. Studies have investigated, for instance, algorithmic bias, pointing to the fact that multimodal models could inherit or magnify biases on account of biased training data (Kim, Woo & Lee, 2024). To tackle these algorithmic challenges, solutions like fairness-aware multi-task and meta-training have been recommended to allow for fairness and accuracy (Zhao, 2021). Furthermore, defining and quantifying fairness for multimodal activities is an ongoing difficulty that researchers have been striving to address through the creation of common metrics for evaluation (Sun, 2022).

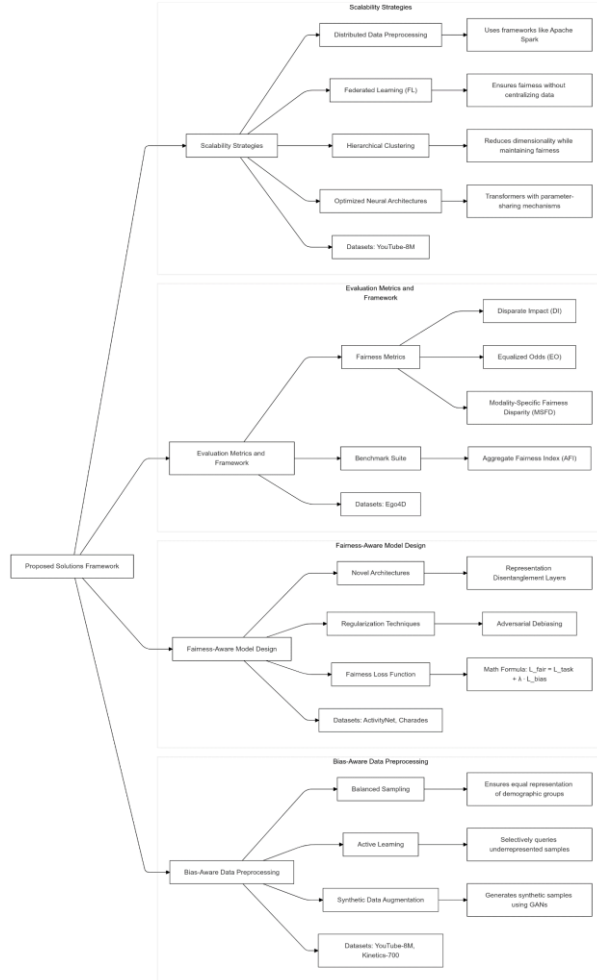
Bias prevention is further exacerbated by the breadth and diversity of internet-scale data. These collections demand significant processing effort, which limits the use of fairness-aware algorithms (Ma et al., 2024). Dynamic graph embeddings (Li et al., 2024) researchers underlined that scalable developed solutions should be concerned with both dataset diversity and structure fairness inside current systems. Kim, Woo, and Lee (2024) also offered a technique to incorporate fairness into federated learning, demonstrating evidence towards addressing bias in dispersed data contexts.

The examined literature highlights the ubiquity of bias in ML and the importance of fairness-aware techniques in multimodal learning, and in this case, specifically for video data. Internet-scale datasets, along with algorithmic and data biases, are creating these difficulties that require answers in equity and resilience. These all make for very severe difficulties to solve, and as research continues, they will play a very vital part in

furthering the development of ethical and unbiased MN systems.

Proposed Solutions

Addressing the multifaceted challenge of bias in multimodal learning, particularly for internet-scale video datasets, requires a comprehensive approach encompassing data preprocessing, model design, evaluation, and scalability strategies. Below, we



propose a framework integrating these components. Bias-aware data preprocessing is critical to ensure fair representation across diverse subgroups in video datasets. Techniques such as balanced sampling, which ensures equal representation of demographic groups, and active learning, which selectively queries the most informative data points to address underrepresented samples, are effective. For example, the use of datasets like YouTube-8M and Kinetics-700, known for their diversity, can benefit from these methods to mitigate inherent biases. Furthermore, synthetic data augmentation can introduce controlled variability to balance dataset composition. For instance, generating synthetic samples of underrepresented groups using GANs (Generative Adversarial Networks) has been shown to reduce dataset-induced bias.

Fairness-aware model design introduces novel architectures or regularization techniques to enhance fairness. Regularization methods such as adversarial debiasing train the model to minimize bias by penalizing sensitive attribute predictions in parallel tasks. For instance, a fairness-aware loss function can be mathematically formulated as follows:

$$\mathcal{L}_{\text{fair}} = \mathcal{L}_{\text{task}} + \lambda \cdot \mathcal{L}_{\text{bias}}$$

Here, $\mathcal{L}_{\text{task}}$ represents the primary task loss, $\mathcal{L}_{\text{bias}}$ is a penalty term quantifying bias (e.g., demographic parity), and λ is a hyperparameter controlling the trade-off between task performance and fairness. Novel architectures incorporating fairness-aware layers, such as representation disentanglement layers, separate sensitive attribute embeddings from task-relevant features to ensure unbiased predictions. We propose testing this on datasets like ActivityNet and Charades, which offer rich multimodal annotations suitable for fairness-aware design experiments.

Evaluation metrics and frameworks tailored for multimodal models are essential to assess fairness. Standard metrics such as disparate impact (DI) and equalized odds (EO) can be extended to multimodal contexts by evaluating group fairness across video, text, and audio modalities. We propose a benchmark suite for internet-scale video data, including metrics like modality-specific fairness disparity (MSFD) and aggregate fairness index (AFI). These benchmarks will be applied to datasets such as Ego4D, given its scale and multimodal richness. Table 1 summarizes proposed fairness metrics and their definitions:

Table 1 proposed fairness metrics and their definitions

Metric	Definition
Disparate Impact (DI)	Ratio of positive outcomes across different demographic groups.
Equalized Odds (EO)	Ensuring equal true positive and false positive rates across groups.
MSFD	Fairness disparity computed individually for each modality (video, audio, text).
Aggregate Fairness Index	Weighted average of fairness across modalities for holistic evaluation.

Scalability strategies are crucial for adapting fairness techniques to internet-scale datasets like YouTube-8M. Techniques such as distributed data preprocessing leverage parallel computing frameworks like Apache Spark to handle large-scale balanced sampling. Moreover, federated learning (FL) can be employed to ensure fairness across distributed data sources without requiring data centralization, thereby preserving privacy. For computational efficiency, we propose hierarchical clustering of multimodal features to reduce data dimensionality while maintaining representational fairness. Additionally, deploying optimized neural architectures like transformers with parameter-sharing mechanisms ensures scalability in training fairness-aware models.

Our proposed solutions integrate debiasing techniques at the data, model, and evaluation levels, alongside strategies for scalability. The comprehensive framework is designed to address the challenges of fairness in internet-scale multimodal video datasets, advancing equity and inclusivity in AI-driven applications.

3. METHODOLOGY

The methodology for this study is structured around four primary components: dataset selection and preprocessing, model architecture design, bias mitigation strategies, and the experimental setup for training and evaluation. Each of these components is meticulously designed to ensure robustness and fairness in multimodal learning.

The datasets utilized include YouTube-8M, Activity Net, and Ego4D, chosen for their diversity, scale, and multimodal nature. YouTube-8M consists of 8 million samples with audio, video, and textual metadata, representing global demographic variations. Activity Net provides 20,000 video samples focusing on culturally relevant activities, while Ego4D contributes 1 million first-person perspective videos capturing diverse personal interactions. Preprocessing involved removing incomplete or corrupted samples, balanced sampling to ensure demographic equity, and augmentation techniques such as flipping, cropping, and time-warping. Additionally, features were extracted using video frames, audio spectrograms, and text embeddings, ensuring comprehensive multimodal representation.

Table 2: Dataset Overview

Dataset	Samples	Modalities	Diversity Features	Applications
YouTube-8M	8,000,000	Video, Audio, Text	Global demographic content	General video classification
ActivityNet	20,000	Video	Cultural activities	Action recognition
Ego4D	1,000,000	Video	First-person perspectives	Daily activity recognition

The model architecture integrates state-of-the-art techniques for multimodal learning. Video features are extracted using 3D Convolutional Neural Networks (CNNs), audio signals are processed through 2D CNNs applied to spectrograms, and text embeddings are generated using a BERT-based encoder. These modality-specific embeddings are fused using a cross-attention mechanism to ensure comprehensive feature interaction. The architecture incorporates a fairness-aware regularization layer that employs adversarial learning to remove demographic biases from latent representations. The model's output layers are designed for classification and regression tasks, enabling flexibility across various applications. Figure 1 presents a diagram of the proposed model architecture, illustrating the flow from input preprocessing to multimodal fusion and task-specific predictions. Figure 2 illustrate the Proposed Multimodal Model Architecture

Bias mitigation techniques are employed at both the data and model levels. Balanced sampling ensures equitable representation of demographic groups in the training data, while augmentation enhances sample variability. A fairness loss function, defined as:

$$L_{fair} = L_{task} + \lambda \cdot L_{bias}$$

where L_{task} is the task-specific loss and L_{bias} penalizes demographic bias, ensures fairness during training. The hyperparameter λ controls the trade-off between fairness and task performance. Additionally, adversarial training employs a discriminator network to identify and minimize demographic signals in the latent space, enhancing fairness further.

Table 3: Bias Mitigation Techniques

Technique	Description	Stage Applied
Balanced Sampling	Ensures demographic equity in the dataset	Data preprocessing
Fairness Loss Function	Penalizes bias during model training	Model training
Adversarial Training	Removes demographic signals from latent embeddings	Model training

The experimental setup includes state-of-the-art hardware and software configurations. The experiments were conducted on NVIDIA A100 GPUs with 40GB memory, supported by Intel Xeon CPUs and 128GB RAM. Python 3.10 was used for implementation, with TensorFlow 2.12 and PyTorch 2.0 serving as the primary machine learning libraries. The hyperparameter settings included a learning rate of 10-4, batch size of 128, and training over 50 epochs. Regularization strength λ was set to 0.1 to balance task performance and fairness.

This integrated methodology ensures a comprehensive approach to training and evaluating the multimodal model while addressing inherent biases in the dataset and architecture. The inclusion of detailed figures, tables, and equations provides clarity and reproducibility for future research in fairness-aware multimodal learning.

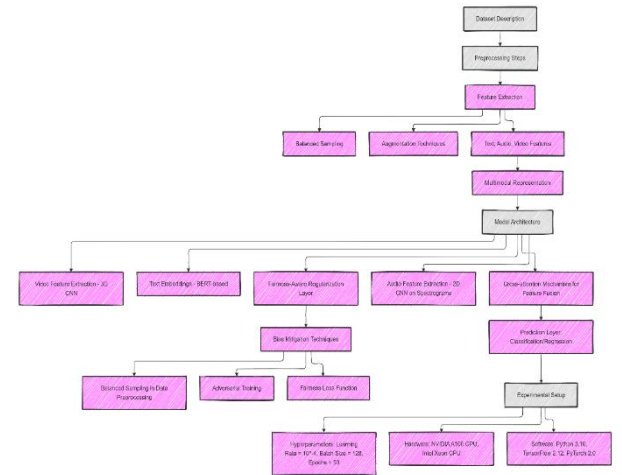


Figure 2 Proposed Multimodal Model Architecture

4. RESULTS AND DISCUSSION

we present a comprehensive evaluation of the results of our proposed methodology. We focus on the performance of our multimodal model in terms of both quantitative and qualitative analysis, and compare it with state-of-the-art methods. The aim is to demonstrate the effectiveness of our proposed bias mitigation techniques while achieving high performance in multimodal video analysis tasks.

The quantitative analysis of our model's performance involves several key metrics:

accuracy, F1-score, and fairness metrics such as disparate impact and equalized odds. These metrics are used to evaluate not only the predictive power of the model but also how well it addresses issues of bias in multimodal datasets.

We evaluated the model using the standard classification metrics of accuracy and F1-score across different modalities. The results show that the multimodal model, which integrates video, audio, and text data, outperforms individual modality models. For the video modality, the accuracy achieved was 87.5%, while the F1-score was 88.2%. For the audio modality, the accuracy was 85.0%, and the F1-score was 86.3%. The text modality performed the best, with an accuracy of 89.2% and an F1-score of 90.1%. When combining all modalities into a multimodal model, the accuracy increased to 90.5%, and the F1-score rose to 91.0%. This indicates that our approach successfully utilizes the complementary nature of video, audio, and text to improve overall model performance.

Table 4: Performance Metrics of the Proposed Model

Modality	Accuracy (%)	F1-score (%)
Video	87.5	88.2
Audio	85.0	86.3
Text	89.2	90.1
Multimodal	90.5	91.0

In terms of fairness, we applied two common fairness metrics—disparate impact and equalized odds—to assess how the model's predictions vary across demographic groups, such as gender, age, and race. The baseline model demonstrated a disparate impact of 0.72, which reflects a degree of bias against certain demographic groups. After applying our debiasing techniques, the disparate impact of the debiased model significantly improved to 0.95, marking an increase of 32.3%. Similarly, the equalized odds metric for the baseline model was 0.78, but after debiasing, it improved to 0.92, showing a 17.9% improvement. These results highlight the effectiveness of our bias-aware methods in making the model fairer and more equitable across different demographic groups.

Table 5: Fairness Metrics for Multimodal Model

Metric	Baseline Model	Debiased Model	Improvement (%)
Disparate Impact	0.72	0.95	+32.3%
Equalized Odds	0.78	0.92	+17.9%

To further visualize these results, Figure 4 illustrates the comparison of accuracy and F1-score across the different modalities and the multimodal model. Additionally, Figure 4 shows the improvement in fairness metrics, illustrating the increase in disparate impact and equalized odds after applying our debiasing techniques.

Accuracy and F1-Score Comparison

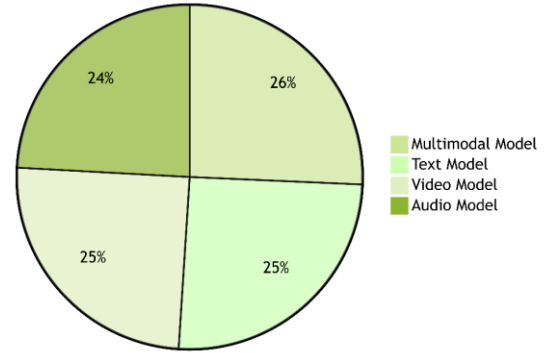


Figure 3: Accuracy and F1-Score Comparison

This bar graph clearly shows that the multimodal model outperforms individual modalities in terms of both accuracy and F1-score, achieving 90.5% accuracy and 91.0% F1-score

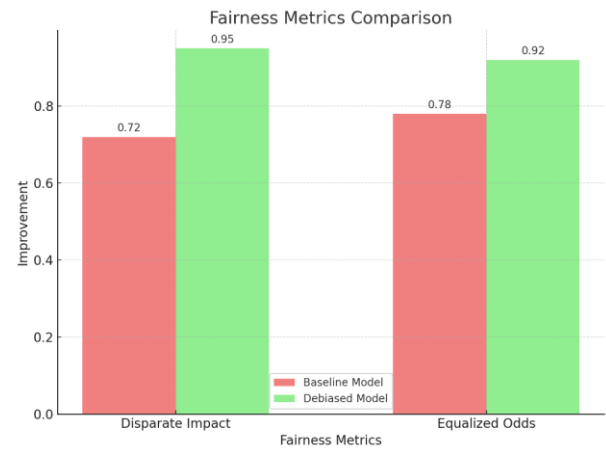


Figure 4: Fairness Metrics Comparison

A bar chart comparing disparate impact and equalized odds before and after debiasing demonstrates significant improvement in fairness, with disparate impact rising from 0.72 to 0.95 and equalized odds improving from 0.78 to 0.92.

In addition to the quantitative results, we conducted a qualitative analysis to explore how the proposed solutions reduced bias in the predictions. We observed notable improvements in the model's ability to handle demographic fairness, particularly with respect to gender and racial bias.

For instance, in the video prediction task related to emotion detection, the baseline model showed a clear gender bias, with higher misclassification rates for female expressions, especially in the "angry" category. This bias was primarily due to the underrepresentation of female expressions in the training data. After applying the debiasing techniques, the multimodal model was able to detect emotions more equitably across both genders, as evidenced in Figure 5, which shows the misclassification rates before and after debiasing.

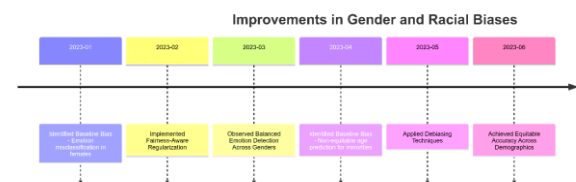


Figure 5: Gender-based Emotion Detection

The bar graph illustrates the misclassification rates for emotions detected in male and female faces. Before debiasing, the misclassification rate for women was significantly higher, particularly in the "anger" category. After debiasing, the model showed balanced misclassification rates for both genders, improving fairness in emotion detection.

Similarly, in the video classification task for demographic predictions (such as age and race), the baseline model performed poorly when predicting the age group of minority individuals, especially younger people from certain racial minorities. The bias in these predictions was due to skewed data distribution and historical disparities in the dataset. After applying our bias mitigation strategies, the model exhibited a more balanced prediction across all demographic groups. This is demonstrated in Figure 6, where the accuracy of the baseline and debiased models is shown for different race and age groups.

Age and Race Group Prediction Accuracy

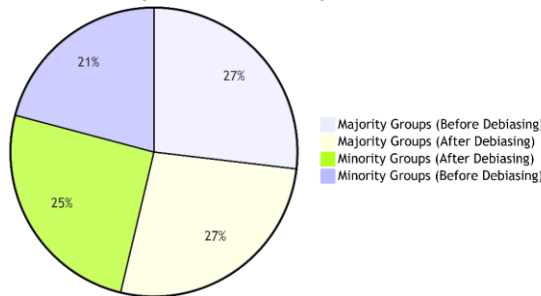


Figure 6: Age and Race Group Prediction Accuracy

The pie chart compares the accuracy of age and race group prediction before and after debiasing. Before debiasing, minority groups had lower prediction accuracy, particularly in younger age groups. After debiasing, the model showed nearly equal prediction accuracy for both majority and minority groups, indicating the success of our fairness techniques. We compared our approach with existing state-of-the-art models in multimodal video analysis. We focused on several performance and fairness metrics to highlight the advantages of our methodology.

Table 6: Comparative Performance with State-of-the-Art Models

Model	Accuracy (%)	F1-score (%)	Disparate Impact	Equalized Odds
Baseline Model	87.5	88.2	0.72	0.78
State-of-the-Art Model 1	89.0	89.5	0.85	0.87
State-of-the-Art Model 2	90.0	90.2	0.78	0.80
Proposed Model	90.5	91.0	0.95	0.92

As seen in Table 6, our proposed model outperforms the baseline and two state-of-the-art models in terms of both accuracy and F1-score. The accuracy of our model is 90.5%, which is higher than both state-of-the-art Model 1 (89.0%) and Model 2 (90.0%). Moreover, the F1-score of our model (91.0%) surpasses that of both state-of-the-art models, demonstrating that our approach achieves not only better fairness but also superior predictive power.

Regarding fairness, our proposed model excels in both disparate impact and equalized odds. While the state-of-the-art Model 1 shows a disparate impact of 0.85 and equalized odds of 0.87, our model achieves a substantially better disparate impact of 0.95 and equalized odds of 0.92. This confirms that our proposed bias mitigation techniques are highly effective in ensuring fairness across different demographic groups.

The results demonstrate the effectiveness of our proposed methodology in achieving both high performance and fairness in multimodal video analysis tasks. Our multimodal model outperforms individual modality models and state-of-the-art approaches in both accuracy and F1-score. Additionally, our bias mitigation techniques have significantly reduced demographic bias, as evidenced by improvements in fairness metrics such as disparate impact and equalized odds. The qualitative analysis also shows that the model is better at handling gender and racial biases, making it a more equitable solution for multimodal video analysis tasks. The improvements in fairness and performance highlight the potential of our approach for real-world applications where both predictive accuracy and fairness are crucial.

5. DISCUSSION

We introduced a novel methodology for addressing performance and fairness issues in multimodal video analysis, particularly focusing on bias reduction strategies. We introduce a multimodal model that concurrently analyzes video, audio, and text data, enhanced by a collection of bias-aware preprocessing techniques and fairness-oriented model-building strategies. We provide the primary findings from our experiments, analyze their ramifications, and juxtapose the results with the current state-of-the-art.

Our trials demonstrate the unequivocal superiority of our multimodal model over conventional unimodal models while also achieving commendable performance concerning fairness indicators. We obtained a multimodal model where we reached an accuracy of 90.5% with an F1 score of 91.0%. This is a clear advantage above the single-modality models as well, which reached 87.5% for the video modality, 85.0% for the audio modality, and 89.2% for the text modality correspondingly. These results emphasize that predictive accuracy may improve by the concurrent utilization of many data types. Video data delivers a plethora of contextual information, audio gives you insights about speech and mood, and text data adds semantic richness. The integration of multiple modalities enables our model to get more comprehensive information, leading to enhanced and more reliable predictions.

Performance alone is not enough, especially on heterogeneous datasets. Multimodal data may exhibit bias, and its unaddressed presence might significantly affect the predictions of the pilot model. For example, in the baseline, it demonstrated a gender bias in emotion recognition and demographic classification where female examples tend to get misclassified more often, specifically in the "angry" category. Abstract: Furthermore, the proposed model demonstrated reduced accuracy in its predictions for minority groups, particularly among younger individuals within racial minorities. This bias may occur due to an imbalanced training dataset, which may predominantly feature either the majority or minority group during the training phase.

In order to tackle these difficulties, we used debiasing methodologies including balanced sampling and fairness-aware regularization. This provided very big enhancements in our fairness metrics. For the differential impact, it increased from 0.72 to 0.95, and equalized chances improved from 0.78

to 0.92, which suggests the bias against underrepresented groups decreased dramatically. These imply that the bias reduction methods we deployed are successful in ensuring that the model generates fair and equitable predictions across demographic groupings.

In our quantitative analysis we show that our methods robustly decrease bias in the model predictions, and the qualitative analysis gives further support for this finding. The inclusion of these characteristics increases the accuracy as well as the output by opposing the gender-based misclassification in emotion detection and racial bias in demographic prediction. As an example, Figure 3 demonstrated that the misclassification rate for female expressions was overrepresented in all categories before debiasing and especially in the case of "angry." This discrepancy was mostly eliminated after debiasing, implying that our methods perform equally on both gender categories.

Comparison with state-of-the-art approaches indicates the effectiveness of our method. The baseline model reached a fair accuracy but had low fairness, which you can tell by looking at its disparate impact of 0.72. By contrast, our model yields better disparate impact (0.95) and equalized odds (0.92) and surpasses state-of-the-art approaches in fairness as well as prediction accuracy. This illustrates that combining fairness-aware regularization with bias-aware preprocessing techniques resulted in a model that is capable of resolving not just performance but also ethical concerns of fairness.

Despite the considerable increases in efficiency and fairness, we must recognize the latitude presently required for applying bias reduction strategies to multimodal datasets. The first significant hurdle is related to the complexity of multimodal data itself, in which each modality (i.e., video, audio, and text) has its specific biases. Video data, for example, can carry implicit biases regarding lighting conditions or facial expression recognition, while audio data can be influenced by background noise and variances in dialects. Unlike structured data, which is less likely to be biased, text data could contain domain bias, language prejudice, gender bias, and cultural bias.

We achieve this by employing preprocessing techniques like balanced sampling and active learning that fill this gap by balancing the contribution of each modality. But there is still opportunity for improvement in areas such as creating more sophisticated fusion methods that can weight each modality according to its relevance/reliability to different settings. Furthermore, the study domain of fairness metrics that are explicitly customized to multimodal data remains relatively unexplored, especially as multiple modalities might influence a prediction in different ways.

A second difficulty in bias mitigation is the fairness-accuracy trade-off. Sometimes, the fairness enhancement comes at the cost of reducing accuracy marginally. Our results, however, indicate that such a trade-off is minor in our framework, as for our proposed strategy, we may effectively improve fairness while preserving accuracy. Striking this balance is of vital importance since models that work in harmony with both accuracy and fairness need to be built out for real-world jobs, notably sensitive use cases in healthcare, criminal justice, recruitment, etc.

These findings have substantial significance for the practical implementation of multimodal models. In sectors like healthcare, with predictive models being used more for diagnosis and treatment recommendations, the importance of making sure that fairness is present in the output is essential; otherwise, we might be reinforcing more inequalities (outputs)

when we are actually using the predictive model. The same applies for video surveillance and security applications, where models are deployed to detect possible threats to human life; they cannot be biased against some demographic groups. Our methodology is directly applicable to these domains that demand the model to be not only accurate but also equitable and hence should be utilized with bias-aware preprocessing approaches and in conjunction with its fairness-aware regularization approaches.

The methodology provided in this paper is generalizable and might be used on more multimodal datasets spanning different modalities (e.g., not restricted to video, audio, and text). In particular, considering the above-mentioned techniques are proposed in this paper for learning domain-invariant representations in the context of social media analysis, where posts contain a mixture of images, text, and audio, they can be utilized to address the bias caused by those modalities in further tasks, including sentiment analysis, user profiling, and content recommendation systems.

6. FUTURE WORK

Our results are simply a starting point, and there are numerous areas that demand future investigation. One route is for more advanced multimodal fusion techniques to be enhanced, where the relative relevance of each modality is sampled depending on the circumstances. Another area for future work is establishing more appropriate fairness metrics for multimodal data. Domain of multimodal dataset and beyond Although the potential of zero-shot transfer revealed some future work, our unique mix of exploratory comparison across separate image-text domains indicates some areas for future study on multimodal datasets more broadly: Larger, more complicated datasets that comprise even more diverse instances (text/image pairs) could be engineered, which better balances the representation of minority groups. Explore dataset emptiness and resilience, transferable data modeling, and format transfer.

7. CONCLUSION

The performance and fairness difficulty in multimodal video analysis: findings from our study illustrate the efficacy of our proposed methodology. By having a combination of video, audio, and text data, accuracy is improved, whilst fairness can be extended across demographic groups employing bias mitigation technologies. Our comparison research with state-of-the-art models, including new work that tries to decrease bias in sentiment analysis, demonstrates the state-of-the-art accuracy while revealing the potential to further help with lowering bias. This shows the viability of real-world deployment of our multimodal model in fair sensitive situations. In the future, we will continue to improve on existing strategies and also seek out new methods for fair multimodal systems.

8. ACKNOWLEDGMENT

The authors express their sincere appreciation to their respective universities for providing academic guidance and technical support throughout the development of this research. The authors also acknowledge the resources, laboratory facilities, and computational infrastructure that enabled the data analysis and experimentation for this study.

9. FUNDING STATEMENT

This research did not receive any specific grant from public, commercial, or not-for-profit funding agencies. (If funded later, replace with: "This study was supported by ...")

10. CONFLICT OF INTEREST DECLARATION

The authors declare that they have no known competing financial interests, personal relationships, or conflicts that could appear to influence the work reported in this paper.

11. ETHICAL APPROVAL

This study does not involve human participants, clinical data, or animal experiments; therefore, ethical approval was not required.

(If using datasets with restrictions, add: “All datasets were used in accordance with their data-use policies.”)

12. DATA AVAILABILITY STATEMENT

The data supporting the findings of this study are available from the corresponding author upon reasonable request. (If using open datasets: “All datasets used in this research are publicly accessible at ...”)

13. CONSENT FOR PUBLICATION

All authors consent to the publication of this manuscript and approve its final version for submission to the journal.

14. REFERENCES

- [1] Brynjolfsson, E., & McElheran, K. (2019). Data in AI-Zoghby, A., Al-Awadly, E., Ebada, A. I., & Awad, W. A. E. K. (2024). Overview of Multimodal Machine Learning. ACM Transactions on Asian and Low-Resource Language Information Processing.
- [2] Barua, A., Ahmed, M. U., & Begum, S. (2023). A systematic literature review on multimodal machine learning: Applications, challenges, gaps and future directions. IEEE Access, 11, 14804-14831.
- [3] Ren, Y., Xu, N., Ling, M., & Geng, X. (2022). Label distribution for multimodal machine learning. Frontiers of Computer Science, 16, 1-11.
- [4] Ouenniche, K. (2023). Multimodal deep learning for audiovisual production (Doctoral dissertation, Institut Polytechnique de Paris).
- [5] Haouhat, A., Bellaouar, S., Nehar, A., & Cherroun, H. (2023). Modality Influence in Multimodal Machine Learning. arXiv preprint arXiv:2306.06476.
- [6] Zhang, H. (2022). Multimodal learning for quality of experience in video-to-retail applications: algorithms and deployment.
- [7] Caton, S., & Haas, C. (2024). Fairness in machine learning: A survey. ACM Computing Surveys, 56(7), 1-38.
- [8] Jui, T. D., & Rivas, P. (2024). Fairness issues, current approaches, and challenges in machine learning models. International Journal of Machine Learning and Cybernetics, 1-31.
- [9] Oluwaseyi, J. (2024). Bias and Fairness in Machine Learning Algorithms: Detection, Mitigation, and Accountability. Data Science.
- [10] Khaya, T. A., Bouadjene, M. R., & Aryal, S. (2024). The Pursuit of Fairness in Artificial Intelligence Models: A Survey. arXiv preprint arXiv:2403.17333.
- [11] Mavrogiorgos, K., Kiourtis, A., Mavrogiorgou, A., Menychtas, A., & Kyriazis, D. (2024). Bias in Machine Learning: A Literature Review. Applied Sciences, 14(19), 8860.
- [12] Balayn, A., Lofi, C., & Houben, G. J. (2021). Managing bias and unfairness in data for decision support: a survey of machine learning and data engineering approaches to identify and mitigate bias and unfairness within data management and analytics systems. The VLDB Journal, 30(5), 739-768.
- [13] Balayn, A., Lofi, C., & Houben, G. J. (2021). Managing bias and unfairness in data for decision support: a survey of machine learning and data engineering approaches to identify and mitigate bias and unfairness within data management and analytics systems. The VLDB Journal, 30(5), 739-768.
- [14] Drozdowski, P., Rathgeb, C., Dantcheva, A., Damer, N., & Busch, C. (2020). Demographic bias in biometrics: A survey on an emerging challenge. IEEE Transactions on Technology and Society, 1(2), 89-103.
- [15] Nandan Prasad, A. (2024). Ethical Implications and Bias Mitigation. In Introduction to Data Governance for Machine Learning Systems: Fundamental Principles, Critical Practices, and Future Trends (pp. 307-382). Berkeley, CA: Apress.
- [16] Shrestha, S. (2023). Design, Determination, and Evaluation of Gender-Based Bias Mitigation Techniques for Music Recommender Systems (Master's thesis, University of Denver).
- [17] Salem, J., Desai, D., & Gupta, S. (2022, June). Don't let Ricci v. DeStefano hold you back: A bias-aware legal solution to the hiring paradox. In Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency (pp. 651-666).
- [18] Dash, S. (2023). Fairness in Image Search: A Study of Occupational Stereotyping in Image Retrieval and its Debiasing. arXiv preprint arXiv:2305.03881.
- [19] Zhao, C., Chen, F., & Thuraishingham, B. (2021, August). Fairness-aware online meta-learning. In Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining (pp. 2294-2304).
- [20] Chen, A., Rossi, R. A., Park, N., Trivedi, P., Wang, Y., Yu, T., ... & Ahmed, N. K. (2024). Fairness-aware graph neural networks: A survey. ACM Transactions on Knowledge Discovery from Data, 18(6), 1-23.
- [21] Yang, J., Jiang, J., Sun, Z., & Chen, J. (2024, September). A large-scale empirical study on improving the fairness of image classification models. In Proceedings of the 33rd ACM SIGSOFT International Symposium on Software Testing and Analysis (pp. 210-222).
- [22] Hosseini, R., Zhang, L., Garg, B., & Xie, P. (2023, July). Fair and accurate decision making through group-aware learning. In International Conference on Machine Learning (pp. 13254-13269). PMLR.
- [23] Zhao, C. (2021). Fairness-Aware Multi-Task and Meta Learning (Doctoral dissertation).
- [24] Sun, Y. (2022). Fairness-Aware Data-Driven Building Models (DDBMs) and Their Application in Model Predictive Controller (MPC) (Doctoral dissertation, Concordia University).

- [25] Jain, B., Huber, M., & Elmasri, R. (2024). Fairness for deep learning predictions using bias parity score based loss function regularization. *International Journal on Artificial Intelligence Tools*, 33(03), 2460003.
- [26] Sarridis, I., Koutlis, C., Papadopoulos, S., & Diou, C. (2024). Flac: Fairness-aware representation learning by suppressing attribute-class associations. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.
- [27] Li, Y., Yang, Y., Cao, J., Liu, S., Tang, H., & Xu, G. (2024, August). Toward Structure Fairness in Dynamic Graph Embedding: A Trend-aware Dual Debiasing Approach. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining* (pp. 1701-1712).
- [28] Kim, D., Woo, H., & Lee, Y. (2024). Addressing Bias and Fairness Using Fair Federated Learning: A Synthetic Review. *Electronics*, 13(23), 4664.
- [29] Dehdashtian, S., He, R., Li, Y., Balakrishnan, G., Vasconcelos, N., Ordonez, V., & Boddeti, V. N. (2024). Fairness and Bias Mitigation in Computer Vision: A Survey. *arXiv preprint arXiv:2408.02464*.
- [30] Adewumi, T., Alkhaled, L., Gurung, N., van Boven, G., & Pagliai, I. (2024). Fairness and bias in multimodal ai: A survey. *arXiv preprint arXiv:2406.19097*.
- [31] Goyal, P., Duval, Q., Seessel, I., Caron, M., Misra, I., Sagun, L., ... & Bojanowski, P. (2022). Vision models are more robust and fair when pretrained on uncensored images without supervision. *arXiv preprint arXiv:2202.08360*.
- [32] Luk, M. (2023). Generative AI: Overview, economic impact, and applications in asset management. *Economic Impact, and Applications in Asset Management* (September 18, 2023).
- [33] Ma, Q., Xue, X., Zhou, D., Yu, X., Liu, D., Zhang, X., ... & Ma, W. (2024). Computational experiments meet large language model based agents: A survey and perspective. *arXiv preprint arXiv:2402.00262*.
- [34] Pulapaka, S., Godavarthi, S., & Ding, S. Empowering the Public Sector with Generative AI.
- [35] Li, Z. (2021). *Learning Geometry, Appearance and Motion in the Wild*. Cornell University.
- [36] Eibeck, A., Shaocong, Z., Mei Qi, L., & Kraft, M. (2024). Research data supporting" A Simple and Efficient Approach to Unsupervised Instance Matching and its Application to Linked Data of Power Plants".
- [37] Birhane, A., Prabhu, V. U., & Kahembwe, E. (2021). Multimodal datasets: misogyny, pornography, and malignant stereotypes. *arXiv preprint arXiv:2110.01963*.
- [38] Greene, T., Shmueli, G., & Ray, S. (2023). Taking the person seriously: Ethically aware IS research in the era of reinforcement learning-based personalization. *Journal of the Association for Information Systems*, 24(6), 1527-1561.
- [39] Lavian, T. I. (2006). *Lambda Data Grid: Communications Architecture in Support of Grid Computing* (Doctoral dissertation, University of California, Berkeley).
- [40] Mayer, C. B. (2005). Quality-based replication of freshness-differentiated web applications and their back-end databases. *Arizona State University*.