

Reducing Gender and Racial Biases in Multimodal Vision Language Models via Contrastive Debiasing and Representative Dataset Synthesis

Dylan Clarke

Department of Electrical and Computer Engineering, Boise State University
dylan.clarke@boisestate.edu

Abstract

The rapid proliferation of multimodal vision language models has transformed the landscape of artificial intelligence, enabling unprecedented capabilities in image captioning, visual question answering, and generative content synthesis. However, these advancements have also surfaced profound ethical concerns regarding the systematic reproduction of gender and racial biases inherent in foundational training data. This research investigates a holistic systems-level framework for mitigating these biases through a dual-pronged strategy involving contrastive debiasing architectures and representative dataset synthesis. Rather than relying on post-hoc filtering or simple data augmentation, the proposed approach integrates fairness as a core architectural constraint during the alignment phase of multimodal training. By synthesizing diverse datasets that fill historical representation gaps and implementing contrastive learning objectives that penalize stereotypical associations between protected attributes and unrelated semantic features, the framework aims to decouple social identity from task-related performance. The discussion encompasses the structural trade-offs between model accuracy and fairness, the infrastructure required for large-scale synthetic data generation, and the governance frameworks necessary for auditing multimodal systems. Ultimately, this paper argues that achieving equitable artificial intelligence requires an interdisciplinary commitment to re-engineering the socio-technical pipeline, moving beyond algorithmic fixes toward sustainable, bias-aware systemic infrastructures.

Keywords:

Multimodal Learning, Vision Language Models, Algorithmic Fairness, Representative Data Synthesis, Contrastive Debiasing, Socio-technical Systems.

1. Introduction

The evolution of artificial intelligence has transitioned from isolated unimodal processing toward complex multimodal vision language models that attempt to mirror the human capacity for integrating visual and linguistic information. These models, often characterized by massive parameter counts and trained on internet-scale datasets, serve as the backbone for contemporary search engines, creative tools, and automated decision-making systems. Despite their remarkable performance in benchmark evaluations, a critical structural flaw

persists: the tendency of these models to internalize and amplify societal prejudices related to gender and race [12]. Because the vast repositories of data used for training are often reflections of historical inequities and cultural hegemonies, the resulting models frequently associate specific racial profiles or genders with limited social roles, professions, or moral attributes. This systemic bias is not merely a technical glitch but a fundamental challenge to the integrity of AI infrastructure, posing significant risks to social equity and the reliability of automated governance [22].

The socio-technical complexity of multimodal systems necessitates a departure from traditional unimodal debiasing techniques. While text-only models may struggle with linguistic stereotypes, vision language models face the compounded challenge of visual representation, where subtle cues in imagery—such as lighting, background context, and physical features—interact with textual labels to solidify biased associations. For instance, a model might consistently associate domestic settings with female subjects or professional medical environments with specific racial groups, even when the input prompts are intended to be neutral [16]. Addressing these issues requires a multi-layered intervention that spans the entire lifecycle of the model, from the initial collection and synthesis of training data to the architectural design of the learning objectives and the final deployment protocols.

This paper proposes an integrated system for reducing gender and racial bias through the synergy of contrastive debiasing and representative dataset synthesis. Contrastive debiasing operates at the loss function level, encouraging the model to map diverse demographic representations into a shared latent space where identity-irrelevant features are prioritized. Concurrently, representative dataset synthesis addresses the data desert problem, where certain intersections of race and gender are chronically underrepresented in natural datasets. By leveraging generative models to create balanced, high-fidelity synthetic samples, we can provide the model with the counter-stereotypical evidence necessary to break harmful associations [3]. The following sections provide a comprehensive analysis of the infrastructure, trade-offs, and policy implications of this approach, emphasizing the role of fairness as a non-negotiable metric of architectural robustness.

2. Theoretical Foundations of Multimodal Bias

Understanding bias in multimodal vision language models requires an appreciation of how visual and linguistic modalities intersect during the feature extraction and alignment processes. Standard architectures typically employ a dual-encoder system, where a visual transformer and a language model project their respective inputs into a common embedding space. The alignment is often achieved through contrastive loss, which maximizes the similarity between paired images and captions. However, if the training pairs consistently link doctor with male-coded images and nurse with female-coded images, the model learns a high-dimensional vector space where gender is a predictive feature for professional roles [4]. This phenomenon, known as representational harms, occurs when the model’s internal logic reinforces the marginalization of specific groups [17].

The theoretical root of this problem lies in the statistical distribution of the training data.

Large-scale web-scraped datasets are inherently noisy and reflective of the winner-takes-all dynamics of internet visibility. Groups that have historically held social power are overrepresented, while marginalized communities are either invisible or depicted through narrow, often derogatory, lenses [1]. When a vision language model is trained on such a skewed distribution, it does not simply learn to recognize objects; it learns the social hierarchies embedded within the data. This creates a feedback loop where biased model outputs inform future data creation, further entrenching these prejudices in the digital ecosystem [28].

Furthermore, the complexity of multimodal alignment means that bias can hide in the interactions between modalities. A model might appear neutral when processing text alone but demonstrate significant racial bias when presented with a visually ambiguous scene. This cross-modal amplification is particularly dangerous in high-stakes environments, such as security surveillance or automated recruitment, where the model's black box nature can mask discriminatory logic behind a veneer of mathematical objectivity [10]. Consequently, the research community must move toward a more rigorous definition of fairness that accounts for the intersectional nature of race and gender, recognizing that bias is not a monolithic variable but a multidimensional challenge influenced by cultural context and historical power structures [25].

3. Infrastructure for Representative Dataset Synthesis

One of the primary obstacles to training fair multimodal models is the scarcity of high-quality, balanced datasets. Traditional manual curation is prohibitively expensive and time-consuming, leading many developers to rely on flawed web-crawled collections. Representative dataset synthesis offers a technological solution by using generative AI to produce high-fidelity training data that purposefully fills representation gaps [8]. This infrastructure involves the deployment of state-of-the-art diffusion models and large language models to generate images and corresponding descriptions that span a wide range of demographic intersections. For example, the system can be programmed to generate images of female engineers of diverse ethnic backgrounds in various professional settings, thereby providing the vision language model with the necessary variance to decouple gender from career attributes [26].

The deployment of such a synthesis pipeline involves significant computational trade-offs. Generating millions of high-resolution images requires substantial GPU resources and robust storage infrastructure. Moreover, the quality of synthetic data must be carefully monitored to ensure it does not introduce new forms of bias or artifacts that could degrade model performance. A major structural consideration is the fidelity-diversity trade-off, where an over-emphasis on visual realism might inadvertently mirror the biases of the generative model itself [3]. To mitigate this, a multi-stage validation process is required, incorporating both automated auditing tools and human-in-the-loop oversight to verify that the synthetic samples accurately and respectfully represent the intended demographics [19].

Beyond the technical requirements, the governance of synthetic data creation presents complex ethical questions. There is a risk that synthetic data could be used to create tokenistic

representations that lack cultural authenticity. As noted in recent research, when cultural nuances are stripped away in favor of generic diversity, the resulting models may suffer from a cultural gap that limits their effectiveness in global contexts [22]. Therefore, the infrastructure for dataset synthesis must be designed with cultural sensitivity, incorporating diverse perspectives into the prompt engineering and latent space manipulation stages [23]. By treating data synthesis as a core component of the AI lifecycle, organizations can build a more resilient foundation for equitable machine learning.

4. Architectures for Contrastive Debiasing

While data synthesis addresses the input side of the equation, contrastive debiasing targets the learning process itself. The core objective of contrastive debiasing is to ensure that the latent representations generated by the vision language model are invariant to protected attributes such as race and gender. In a standard contrastive learning setup, the model is trained to minimize the distance between related image-text pairs [18]. Contrastive debiasing modifies this objective by introducing an additional penalty for representations that allow for the easy classification of gender or race in contexts where those attributes should be irrelevant. This forces the model to ignore stereotypical visual cues and focus on the essential semantic content of the scene [30].

From a systems engineering perspective, implementing contrastive debiasing requires a careful balancing of loss functions. If the debiasing constraint is too aggressive, it may strip away useful information, leading to a decline in general downstream performance. Conversely, if the constraint is too weak, the model will continue to rely on biased shortcuts [9]. Modern architectures often utilize adversarial components, where a debiaser network attempts to predict the protected attribute from the main model's embeddings, and the main model is trained to minimize the debiaser's accuracy [15]. This creates a minimax game that encourages the emergence of fair representations. This structural approach ensures that fairness is not an afterthought but is baked into the very weights and biases of the neural network [7].

The robustness of contrastive debiasing is also dependent on the granularity of the attribute labels. Intersectional fairness requires that the model be debiased against combinations of attributes, such as Black woman or Asian man, rather than treating race and gender as independent variables [2]. This increases the complexity of the training objective and the dimensionality of the embedding space. However, the resulting systems are significantly more resilient to the subtle, multi-layered biases that characterize human social perception. By embedding these fairness constraints into the architecture, developers can create models that are fundamentally more aligned with democratic values and legal standards of non-discrimination [14].

5. System-Level Trade-offs and Performance Metrics

The pursuit of fairness in multimodal systems is rarely a free lunch. Engineering a model to be less biased often involves navigating a complex landscape of trade-offs involving accuracy, computational cost, and interpretability. One of the most frequently discussed dilemmas is the

tension between group fairness and individual accuracy. In some cases, reducing the model's reliance on biased visual cues may slightly lower its performance on specific benchmarks that are themselves grounded in biased data distributions [21]. For example, if a benchmark predominantly features a certain demographic in a specific role, a fair model might perform worse on that benchmark because it refuses to make the biased leap that the test set expects [13].

Measuring the success of a debiased system requires a shift in how we evaluate AI performance. Traditional metrics like Top-1 accuracy or F1-score are insufficient for capturing the nuances of social bias. Instead, researchers are increasingly adopting fairness-aware metrics, such as disparate impact ratios, equalized odds, and counterfactual fairness [5]. These metrics evaluate how the model's performance varies across different demographic groups and how its outputs change when only the protected attribute is modified. Implementing these evaluations requires a robust testing infrastructure, including specialized bias probes and stress-test datasets that target known failure modes [16]. This rigorous auditing process is essential for building public trust and ensuring regulatory compliance in sectors like healthcare, finance, and criminal justice [11].

Sustainability and scalability also play a critical role in system-level discussions. The computational overhead of contrastive debiasing and the storage requirements for massive synthetic datasets can significantly increase the carbon footprint of AI development. Systems designers must therefore consider the environmental impact of their fairness interventions. Strategies such as model distillation, where a large debiased teacher model trains a smaller, more efficient student model, can help mitigate these costs. By optimizing the efficiency of the debiasing pipeline, we can ensure that the move toward more ethical AI is also environmentally and economically sustainable in the long term [6].

6. Governance, Policy, and Socio-technical Infrastructure

The mitigation of bias in vision language models is not solely a technical challenge; it is a governance problem that requires clear policy frameworks and institutional accountability. As AI systems become more integrated into the social fabric, the decisions made by engineers during the data curation and architectural design phases have far-reaching political implications. Consequently, there is an urgent need for standardized protocols regarding the auditing and reporting of bias in multimodal models. Governance structures should mandate the use of Model Cards and Data Sheets that transparently document the demographic composition of training sets and the specific debiasing techniques employed [27].

Policy interventions must also address the legal dimensions of algorithmic bias. In many jurisdictions, the use of biased AI in hiring or lending is already subject to anti-discrimination laws. However, the complexity of multimodal models makes it difficult to pinpoint where discrimination occurs. Regulatory bodies may need to establish new standards for fairness by design, requiring developers to prove that they have taken proactive steps to mitigate bias before a model is deployed. This could include mandatory third-party audits and the creation of a red-teaming infrastructure where independent researchers attempt to trigger biased

behaviors in commercial models. Such a regulatory environment encourages the adoption of the systems-level debiasing techniques discussed in this paper [1].

Furthermore, the socio-technical nature of the problem implies that technical fixes must be accompanied by broader social changes. Diversifying the workforce that builds these models is a crucial step in ensuring that a wider range of perspectives is included in the design process. When the teams developing AI are representative of the global population, they are more likely to identify and address bias early in the development cycle. In essence, the infrastructure of fairness is as much about people and processes as it is about code and data. A holistic approach to AI governance recognizes that the goal is not just to build a neutral machine, but to build a system that actively supports social justice and human dignity [29].

7. Deployment and Real-World Robustness

Moving a debiased model from a controlled laboratory setting to real-world deployment introduces a host of new challenges. Real-world data is often more chaotic and unpredictable than the datasets used during training. A model that appears fair on a curated benchmark may still exhibit biased behavior when confronted with out-of-distribution inputs, such as images with unusual lighting, diverse cultural clothing, or non-standard linguistic dialects. Ensuring robustness in these scenarios requires continuous monitoring and a flexible update pipeline that can respond to newly discovered biases as they emerge in the field [16].

The architecture of deployment must also consider the user experience and the potential for algorithmic gaslighting. This occurs when a system's biased output is presented with such confidence that users begin to doubt their own perceptions or accept the bias as objective truth. To prevent this, systems should incorporate uncertainty quantification and explainability features that allow users to understand why a specific output was generated. For instance, if a vision language model generates a caption that seems racially skewed, an explainability layer could highlight the visual features that influenced the decision, making the bias transparent and easier to challenge. This empowers users to act as a final layer of the debiasing process, creating a more interactive and accountable socio-technical system [10].

Finally, the long-term maintenance of fair models requires a versioning strategy for bias. As societal norms and definitions of fairness evolve, models must be retrained or fine-tuned to reflect these changes. This necessitates a sustainable infrastructure for data archiving and model lineage tracking. By treating AI models as living systems that require ongoing care and adjustment, organizations can ensure that their multimodal applications remain ethical and relevant in a shifting social landscape. This commitment to long-term robustness is the hallmark of a mature and responsible approach to artificial intelligence engineering [19].

8. Forward-Looking Perspectives and Future Directions

The field of multimodal debiasing is still in its infancy, and several promising avenues for future research remain unexplored. One such area is the development of self-correcting models that can identify and mitigate their own biases during inference. By leveraging large language models as internal critics, future vision language architectures could evaluate their

proposed outputs against a set of fairness principles before presenting them to the user. This would add a layer of cognitive debiasing that complements the structural and data-driven approaches discussed in this work. Such systems would represent a significant step toward truly autonomous and ethical artificial intelligence [30].

Another important direction is the exploration of decentralized and collaborative debiasing. Instead of a single organization controlling the data and the architecture, a federated learning approach could allow diverse communities to contribute their own data and fairness constraints to a global model. This would democratize the debiasing process and ensure that the resulting models are representative of a wider array of cultural and linguistic contexts. Infrastructure for secure and private data sharing would be essential for this vision to succeed, highlighting the intersection of fairness, privacy, and distributed systems [13].

As we look toward the future of vision language models, it is clear that the challenge of bias will not be solved by a single silver bullet. It requires a sustained, interdisciplinary effort that integrates technical innovation with social insight and political will. The frameworks proposed in this paper—contrastive debiasing and representative dataset synthesis—provide a robust starting point for this journey. By reimagining the AI pipeline as a socio-technical infrastructure committed to equity, we can harness the power of multimodal models to create a more just and inclusive digital future [5].

9. Conclusion

The reduction of gender and racial bias in multimodal vision language models is a critical imperative for the ethical advancement of artificial intelligence. As demonstrated throughout this research, the systemic nature of bias requires a holistic response that spans data synthesis, architectural design, and governance. By implementing contrastive debiasing, we can force models to learn representations that are invariant to protected attributes, thereby breaking the link between identity and stereotypical associations. Through representative dataset synthesis, we can overcome the historical limitations of web-scraped data, providing a more balanced and inclusive foundation for machine learning.

However, the technical implementation of these strategies is only half the battle. The broader socio-technical infrastructure—including the policy frameworks, auditing protocols, and institutional cultures surrounding AI development—must also be transformed. The trade-offs between accuracy, efficiency, and fairness must be navigated with transparency and a commitment to human values. Ultimately, the goal is to move beyond the reactive patching of biased models toward a proactive fairness-by-design philosophy. This requires a shift in how we define success in the field of AI, prioritizing the social impact and equitable performance of our systems alongside their technical prowess.

The path forward is complex and fraught with challenges, but the potential rewards are immense. A truly fair multimodal model would not only perform better in a diverse world but would also serve as a tool for empowerment, helping to dismantle rather than reinforce societal prejudices. By continuing to innovate in the areas of contrastive learning, synthetic

data generation, and algorithmic governance, the research community can ensure that the next generation of artificial intelligence is built on a foundation of justice, equity, and respect for all individuals, regardless of their race or gender.

References

1. Barocas, S., & Selbst, A. D. (2016). Big data's disparate impact. *California Law Review*, 104, 671.
2. Bender, E. M., Gebru, T., McMillan-Major, A., & Shmitchell, S. (2021). On the dangers of stochastic parrots: Can language models be too big? *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, 610–623.
3. Buolamwini, J., & Gebru, T. (2018). Gender shades: Intersectional accuracy disparities in commercial gender classification. *Proceedings of Machine Learning Research*, 81, 1–15.
4. Caliskan, A., Bryson, J. J., & Narayanan, A. (2017). Semantics derived automatically from language corpora contain human-like biases. *Science*, 356(6334), 183–186.
5. Crawford, K. (2021). *The Atlas of AI: Power, Politics, and the Planetary Costs of Artificial Intelligence*. Yale University Press.
6. DeVries, T., Misra, I., Wang, C., & van der Maaten, L. (2019). Does object recognition work for everyone? *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*.
7. Dwork, C., Hardt, M., Pitassi, T., Reingold, O., & Zemel, R. (2012). Fairness through awareness. *Proceedings of the 3rd Innovations in Theoretical Computer Science Conference*, 214–226.
8. Gebru, T., Morgenstern, J., Vecchione, B., Vaughan, J. W., Wallach, H., Daumé III, H., & Crawford, K. (2021). Datasheets for datasets. *Communications of the ACM*, 64(12), 86–92.
9. Hardt, M., Price, E., & Srebro, N. (2016). Equality of opportunity in supervised learning. *Advances in Neural Information Processing Systems*, 29.
10. Hoffmann, A. L. (2019). Where fairness fails: Data, algorithms, and the limits of antidiscrimination discourse. *Information, Communication & Society*, 22(7), 900–915.
11. Jia, C., Yang, Y., Xia, Y., Chen, Y. T., Parekh, Z., Pham, H., ... & Duerig, T. (2021). Scaling up visual and vision-language representation learning with noisy text supervision. *International Conference on Machine Learning*, 4918–4927.
12. Kay, M., Matuszek, C., & Munson, S. A. (2015). Unequal representation and gender

stereotypes in image search results. Proceedings of the 33rd Annual ACM Conference on Human Factors in Computing Systems, 3819–3828.

13. Karkkainen, K., & Joo, J. (2021). FairFace: Face attribute dataset for balanced race, gender, and age for bias measurement and mitigation. Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision, 1548–1558.
14. Kusner, M. J., Loftus, J., Russell, C., & Silva, R. (2017). Counterfactual fairness. Advances in Neural Information Processing Systems, 30.
15. Liang, P. P., Wu, C. B., Morency, L. P., & Salakhutdinov, R. (2021). Towards understanding and mitigating social biases in language models. International Conference on Machine Learning.
16. Mitchell, M., Wu, S., Zaldivar, A., Barnes, P., Vasserman, L., Hutchinson, B., ... & Gebru, T. (2019). Model cards for model reporting. Proceedings of the 2019 Conference on Fairness, Accountability, and Transparency, 220–229.
17. Noble, S. U. (2018). Algorithms of Oppression: How Search Engines Reinforce Racism. NYU Press.
18. Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., ... & Sutskever, I. (2021). Learning transferable visual models from natural language supervision. International Conference on Machine Learning, 8748–8763.
19. Raji, I. D., & Buolamwini, J. (2019). Actionable auditing: Investigating the impact of publicly naming biased AI systems. Proceedings of the 2019 Conference on Fairness, Accountability, and Transparency, 59–68.
20. Ross, A. S., Hughes, M. C., & Doshi-Velez, F. (2017). Right for the right reasons: Training differentiable agents by constraining their explanations. Proceedings of the 26th International Joint Conference on Artificial Intelligence, 2662–2670.
21. Selvavarapu, R. K., Mohit, J., Singh, A., Berg, A. C., & Berg, T. L. (2020). Choosing the right stakes for fairness. arXiv preprint arXiv:2012.06738.
22. Shi, C., Li, S., Guo, S., Xie, S., Wu, W., Dou, J., ... & Chua, T. S. (2025). Where Culture Fades: Revealing the Cultural Gap in Text-to-Image Generation. arXiv preprint arXiv:2511.17282.
23. Srinivasan, K., & Chander, A. (2021). Biases in generative art: A causal look from the lens of art history. Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency.

24. Tan, H., & Bansal, M. (2019). LXMERT: Learning cross-modality encoder representations from transformers. *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*.
25. Wang, A., Narayanan, A., & Russakovsky, O. (2020). REVISE: A tool for measuring and mitigating bias in visual datasets. *International Conference on European Conference on Computer Vision*, 733–751.
26. Wang, T., Zhao, J., Yatskar, M., Chang, K. W., & Ordonez, V. (2019). Balanced datasets are not enough: Estimating and mitigating gender bias in deep image captioning. *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 5310–5319.
27. Whittaker, M., Crawford, K., Dobbe, R., Fried, G., Kaziunas, E., Varner, M., ... & West, S. M. (2018). *AI Now Report 2018*. AI Now Institute at New York University.
28. Yang, K., Qin, K., Duan, Y., & Russakovsky, O. (2020). Towards fairer datasets: Filtering and balancing the distribution of the people category in ImageNet. *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, 547–558.
29. Zellers, R., Holtzman, A., Bisk, Y., Farhadi, A., & Choi, Y. (2019). HellaSwag: Can a machine really finish your sentence? *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*.
30. Zhao, J., Wang, T., Yatskar, M., Ordonez, V., & Chang, K. W. (2017). Men also do laundry: Multi-attribute bias amplification its mitigation. *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, 2979–2989.
31. Zhou, K., Yang, J., Loy, C. C., & Liu, Z. (2022). Learning to prompt for vision-language models. *International Journal of Computer Vision*, 130(9), 2337–2348.