

# Enhancing Deepfake Detection Through Diverse Data and Augmentation

## 1. Introduction

Deepfakes are images, videos, or audio created or manipulated using AI tools to depict real or fictional people. This synthetic media poses various risks, including spreading misinformation, fraud, and manipulating public opinion. These threats can impact individuals, businesses, and governments, potentially resulting in identity theft, defamation, and financial losses.

Despite extensive research on preventing and identifying deepfakes, challenges persist. A key issue highlighted by several papers is the **overreliance on specific datasets** that lack diversity in languages, cultures, and scenarios. Most current research focuses on datasets like FaceForensics++, primarily including English, Japanese, and Chinese content. (*Salman, Shamsi, & Qureshi, 2023*) This limitation hinders the generalizability of detection methods across other languages and cultural contexts. Furthermore, these datasets often concentrate on simpler deepfakes, potentially failing to represent more complex forgeries adequately.

We aim to evaluate and comprehend the effect of dataset limitations on deepfake research. We'll use a blend of technical methods and qualitative approaches to achieve this while also investigating potential solutions.

## 2. Significance of Deepfake Studies

Deepfakes are important because they introduce both opportunities and dangers, requiring society to adapt through new policies, technologies, and ethical standards to mitigate their negative effects while exploring their potential benefits.

Deepfake technology has become integral to the digital world, cutting film production costs, enabling virtual education, and offering new entertainment experiences. Driven by AI advancements and market demand, its applications span entertainment, marketing, politics, and social media, creating stunning visuals while spreading misinformation. While it offers benefits, deepfakes pose risks like fraud, misinformation, and illegal content, potentially eroding media trust and harming social stability.

Deepfakes are relevant to everyone as they affect how we interact with media, trust information, and engage online. While they offer certain benefits, they also threaten privacy, identity, and democracy. As technology advances, understanding and addressing the socio-technical issues surrounding deepfakes becomes crucial for maintaining an informed society.

## 3. Literature Overview

Previous research on deepfakes explores both the technical advancements in their generation, the challenges of detection, and the ethical issues they raise, alongside assessments of current detection methods. On the technical side, *Salman et al. (2023)*<sup>[1]</sup> offer a comprehensive overview of how deepfakes are generated using techniques like Generative Adversarial Networks (GANs), which allow for the creation of highly realistic and complex deepfake content. This paper

highlights the need for multimodal approaches to address audio and visual deepfakes while acknowledging the challenges posed by limited datasets.

Detection methods have made significant progress. Works by *Wang et al. (2020)*<sup>[5]</sup> and *Wang & Huang (2024)*<sup>[6]</sup> enhance detection by analyzing neuron activation patterns and speech production characteristics.

Assessments such as those by *Passos et al. (2023)*<sup>[9]</sup> point out the limitations in existing detection techniques, including the overreliance on specific datasets and the challenges of cross-dataset validation. *Rana et al. (2022)*<sup>[8]</sup> also emphasize that deep learning-based techniques like CNNs have shown strong detection capabilities. They struggle to consistently detect more complex deepfakes generated by advanced techniques like GANs. Lastly, *Mustak et al. (2023)*<sup>[7]</sup> outline the societal dangers of deepfakes, including misinformation, fraud, and manipulation of public opinion, while proposing strategies such as AI-based detection, legal regulations, and corporate interventions. These works emphasize the importance of a multi-faceted approach to deepfake detection and mitigation, combining technical advancements in generation and detection with better datasets, ethical considerations, and regulatory frameworks to manage the risks and explore the potential benefits of deepfake technology.

#### 4. Limitations & Gaps

**Dataset Limitations:** A key limitation highlighted by several papers is the **overreliance on specific datasets** that lack diversity in terms of languages, cultures, and scenarios. Most current research is focused on datasets like FaceForensics++, which primarily include English, Japanese, and Chinese datasets, limiting the generalizability of detection methods across other languages and cultural contexts. Additionally, these datasets often focus on simpler deepfakes and may not adequately represent more complex forgeries.

**Ethical and Legal Gaps:** Although technical solutions are progressing, there are still significant **gaps in ethical and legal frameworks**. While some papers advocate for AI-based detection technologies, they often overlook the broader ethical and regulatory challenges that need to be addressed. Current regulations, such as GDPR(General Data Protection Regulation), do not adequately cover areas like non-consensual deepfake pornography or political manipulation, leaving individuals and organizations vulnerable to harm.

**Detection Model Vulnerabilities:** include models' vulnerability to challenging conditions, such as noise and visual occlusions, with **DeepSonar** struggling in noisy environments and **ART-AVDF** facing difficulties with extreme head poses. Both models rely heavily on specific criteria or data alignment—**DeepSonar** on neuron coverage criteria and **ART-AVDF** on synchronized audio-visual data—which limits their generalizability.

#### 5. Relevant Theoretical Framework: Ethical Framework

Alongside improving dataset diversity, ethical frameworks will play a critical role in guiding the responsible use of these technologies. Our research draws on three ethical frameworks to ensure that deepfake detection tools are developed and deployed with careful consideration of their social impact. The consequentialism approach emphasizes the importance of evaluating positive outcomes, such as reducing misinformation, and the potential risks, such as privacy violations or

reputational harm from false positives. The non-consequentialism framework focuses on the intentions behind the development and use of detection technologies, advocating for transparency and respect for individual rights. Meanwhile, agent-centered approaches place responsibility on the individuals and organizations deploying these tools, ensuring that ethical oversight is maintained, particularly in sensitive contexts like politics or personal media. Incorporating these diverse theoretical perspectives into the design and implementation of deepfake detection technologies is essential for creating systems that are not only technically effective but also ethically sound.

## **6. Proposed Solution**

### **Improving Deepfake Detection with Diverse Data and Augmentation**

Using diverse data and data augmentation is important for improving deepfake detection models to work well with different media types and cultures. For images, changes like rotating, blurring, adjusting colors, and adding slight noise help models recognize deepfakes in various conditions. In audio, changes like shifting pitch, stretching time, and adding background sounds help models handle different voice manipulations. For videos, skipping frames and adjusting playback speed make it easier for models to detect complex deepfakes. Additionally, gathering more data from social media platforms across different countries and regions ensures models are exposed to diverse types of content, accents, and visual styles, allowing for improved detection across global media.

With diverse data, models become better at handling different situations. By learning from manipulated data, they get better at identifying deepfakes in real-world settings and are less likely to overfit to specific datasets. Adding noisy or adversarial data during training strengthens models against attacks or subtle changes. Also, by combining information from images, audio, and text, models can detect deepfakes more accurately across different media types, helping create a well-rounded and reliable detection system that works across languages, cultures, and complex scenarios, including those encountered in various social media platforms globally.

### **Model Training**

We will first explore collaboration opportunities to use the models from Wang et al. (2020) or Wang & Huang (2024) as our pre-trained model. Alternatively, several pre-trained models such as Xception, Vision Transformer (ViT), and ResNet50 present excellent choices. A detailed comparison of these models based on key attributes is provided in the Appendix A to facilitate selection. After selecting a pre-trained model, we will fine-tune it on our dataset and evaluate its performance using various metrics to assess its effectiveness comprehensively.

### **Ethical and Legal Considerations**

While the primary focus is on technical solutions, ethical and legal considerations are essential for ensuring that deepfake detection tools are used responsibly. Future research should align with existing regulations like GDPR to protect user privacy while addressing deepfake misuse, emphasizing transparency and accountability in media platforms. Additionally, detection models

must prioritize individual privacy and minimize false positives to prevent potential harm to personal reputations, particularly in sensitive areas like politics or personal media.

### Type of Study: Mixed Methods

This research will use a **mixed-methods approach**, combining:

- **Quantitative:** Evaluating the performance of detection models before and after incorporating the extended dataset.
- **Qualitative:** Collecting stakeholder insights through surveys on the ethical and technological aspects of addressing deepfake challenges.

### Research Design:

#### What Data Helps You Answer the Question?

The research will utilize both quantitative and qualitative data to evaluate the effectiveness of augmented deepfake detection models. Quantitative data will include augmented deepfake datasets across various languages and scenarios, spanning image, audio, and video formats. On the qualitative side, expert feedback from AI researchers, policymakers, and legal professionals will provide insights into the legal and ethical implications of deploying these enhanced detection systems.

#### How Will You Collect the Data?

For the quantitative aspect, the baseline dataset DFDC will establish a benchmark for evaluating deepfake detection models. Additionally, videos and images will be collected from social media platforms such as TikTok, YouTube, and Instagram using official APIs, including content from multiple regions to enhance generalization. The collected data will be integrated into the baseline datasets to create an enriched, comprehensive dataset representing various scenarios, including underrepresented languages, accents, cultural styles, and manipulation techniques. Augmented data will be generated using transformations and synthetic deepfake creation techniques, and these datasets will be used to train and test deepfake detection models. Qualitative insights will be gathered through a survey of open-ended and Likert-scale questions. This survey will target legal and AI experts to address concerns about deploying these models in sensitive areas like privacy and political usage. Additionally, ethical and diversity experts will be consulted to identify potential ethical issues and consequences, ensuring responsible and inclusive data collection practices. The details of the survey questionnaire are provided in Appendix B.

#### How Will You Analyze the Data?

The qualitative method uses a survey to gather valuable insights from stakeholders. It is a consultation tool to assess the system's effectiveness in real-world scenarios. These insights highlight the ethical aspects and support adjustments to the quantitative approach as needed. We will then conduct pre- and post-augmentation comparisons utilizing standard metrics, including accuracy, recall, and F1 scores, while also incorporating areas under the AUC-ROC for a more complete performance assessment. The models' robustness will be tested under various

real-world conditions, including scenarios with different noise levels, varying lighting conditions, and content across multiple languages and cultural contexts.

## 7. Limitations of the proposed solutions

### Measuring Real-World Effectiveness

The model's performance can only be fully assessed in real-world scenarios, where factors like user behavior and data variability can impact its accuracy. While it may perform well in controlled tests, it could face challenges in new contexts or with unseen manipulations, highlighting the need for continuous real-world testing and updates.

### Dependence on API Tools

Our approach relies on third-party APIs for data collection, which can introduce limitations such as restricted access to specific data, potential quality issues, and platform policy constraints.

### Challenges in Identifying Authentic Data From Social Media Platforms

Collecting data from social media platforms often involves mislabeled or manipulated content, making it difficult to ensure authenticity. Automated validation tools are prone to errors, and manual verification is impractical at scale. However, some platforms have AI detection models that label content as potentially AI-generated. These labels can provide a useful starting point for filtering and validating data, though they still require further verification to ensure accuracy.

## 8. Conclusion

The proposed solution, using diverse datasets and data augmentation techniques, can significantly improve the real-world performance of deepfake detection systems. By incorporating a broader range of languages, cultures, and media types, we can create more robust models capable of detecting deepfakes across global contexts. This approach is crucial for addressing deepfake-related threats in sectors like social media, politics, and entertainment.

This project highlights the importance of integrating technical advancements with ethical and legal frameworks. Only through a multidisciplinary approach can we develop deepfake detection systems that are both effective and responsible deepfake detection systems to protecting society rather than exacerbating the digit.

This project taught us the importance of a step-by-step approach, from data collection to evaluation, and the value of adapting machine learning models for real-world use. We learned to focus on research methods, exploring possibilities rather than just seeking accurate answers, and recognized that dataset diversity is essential for generalizing detection models. However, ensuring ethical data collection and responsible deployment remains a significant challenge. Additionally, the rapid evolution of deepfake technology highlights the need for constant model updates to maintain effectiveness while ensuring technology is used responsibly to protect privacy, prevent harm, and build trust.

## Contribution

Name	Contribution
Chun Wen Liou	Paper review, Significance of Deepfake Studies ( three status quo paper), Relevant Theoretical Framework, Qualitative Method(Survey Questions), Slides.
Jiajun Huo	Paper review on evaluation method on the current algorithm. Made slides and wrote certain parts of the proposal, including quantitative method, introduction, and overview.
Yun Chiao Cheng	I have completed tasks including a paper review, analyzing related work, designing a solution framework, and implementing a quantitative method to support my research. Additionally, I created a comprehensive slide deck to present my findings effectively.
Yu-Chen Su	Contributed to the paper review (analyzing two deepfake detection models), developed the Related Work and Solution Design sections, and designed the Qualitative Method, including survey questions with a focus on Likert scales. Additionally, I converted the survey into a Google Form and contributed to the creation of presentation slides.

## Reference

- [1] Salman, S., Shamsi, J. A., & Qureshi, R. (2023). Deep fake generation and detection: Issues, challenges, and solutions. \*IT Professional, 25\*(1), 52-59. <https://doi.org/10.1109/MITP.2022.3230353>
- [2] Meskys, E., Liaudanskas, A., Kalpokiene, J., & Jurcys, P. (2020). Regulating deep fakes: Legal and ethical considerations. \*Journal of Intellectual Property Law & Practice, 15,\* 24-31. <https://doi.org/10.1093/jiplp/jpz167>
- [3] van der Sloot, B., & Wagenveld, Y. (2022). Deepfakes: Regulatory challenges for the synthetic society. \*Computer Law & Security Review, 46,\* 105716. <https://doi.org/10.1016/j.clsr.2022.105716>
- [4] Li, M., & Wan, Y. (2023). Norms or fun? The influence of ethical concerns and perceived enjoyment on the regulation of deepfake information. \*Internet Research, 33\*(5), 1750-1773
- [5] Wang, R., Juefei-Xu, F., Huang, Y., Guo, Q., Xie, X., Ma, L., & Liu, Y. (2020). DeepSonar: Towards effective and robust detection of AI-synthesized fake voices. In \*Proceedings of the 28th ACM International Conference on Multimedia (MM '20)\* (pp. 1207-1216). Association for Computing Machinery. <https://doi.org/10.1145/3394171.3413716>

- [6] Wang, Y., & Huang, H. (2024). Audio–visual deepfake detection using articulatory representation learning. \*Computer Vision and Image Understanding, 248,\* 104133. <https://doi.org/10.1016/j.cviu.2024.104133>
- [7] Mustak, M., Salminen, J., Mäntymäki, M., Rahman, A., & Dwivedi, Y. K. (2023). Deepfakes: Deceptions, mitigations, and opportunities. \*Journal of Business Research, 154,\* 113368. <https://doi.org/10.1016/j.jbusres.2022.113368>
- [8] Rana, M. S., Nobi, M. N., Murali, B., & Sung, A. H. (2022). Deepfake detection: A systematic literature review. \*IEEE Access, 10,\* 25494-25513. <https://doi.org/10.1109/ACCESS.2022.3154404>
- [9] Passos, L. A., Jodas, D., Costa, K. A. P., Souza Júnior, L. A., Rodrigues, D., Del Ser, J., Camacho, D., & Papa, J. P. (2024). A review of deep

learning-based approaches for deepfake content detection. \*Expert Systems, 41\*(8), e13570. <https://doi.org/10.1111/exsy.13570>

[10] Booz Allen. (n.d.). Deepfakes pose business risks—Here’s what to know. Retrieved from <https://www.boozallen.com/insights/ai-research/deepfakes-pose-businesses-risks-heres-what-to-know.html>

[11] U.S. Government Accountability Office (GAO). (2024, March 11). \*Science & Tech Spotlight: Combating deepfakes\*. <https://www.gao.gov/products/gao-24-107292>

[12] RecordedFuture. (n.d.). \*2024 deepfakes and election disinformation report: Key findings & mitigation strategies\*. Retrieved from <https://www.recordedfuture.com/research/targets-objectives-emerging-tactics-political-deepfakes>

## AI Disclosure

We used ChatGPT (version GPT-4) for this assignment. Access link: <https://chatgpt.com/>. I used the tool several times to proofread and fix the grammar. We critically reviewed and revised all AI-generated content to ensure accuracy and relevance. We can provide the unedited transcript with prompts, interactions, and output upon request. I take full responsibility for the final work, ensuring it is ours.

## Appendix A Comparison of Open Source Pre-trained Model

Model	Best for	Image Detection	Video Detection	Key Strength
Xception	Image, Video	Excellent	Excellent	Captures pixel-level artifacts
Vision Transformer (ViT)	Video(frame-based)	Good	Excellent	Captures Global context
ResNet50	Image	Good	Limited	Residual learning

## Appendix B: Survey Questionnaire <https://forms.gle/BLsRBrw87c4JQdgw7>

### Section 1: Ethical Challenges and Legal Considerations

1. **What ethical challenges do you foresee in deploying deepfake detection technologies in public and private sectors, especially in sensitive contexts such as politics, journalism, or personal media?**  
*[Open text box for responses]*
2. **What measures should be taken to avoid false positives in deepfake detection systems, particularly in high-stakes areas like political content or private individuals' media?**  
*[Open text box for responses]*
3. **What ethical considerations should be prioritized when deploying deepfake detection systems in sensitive areas such as politics, journalism, or entertainment?**  
*[Open text box for responses]*
4. **What legal challenges do you think will arise in regulating deepfake content, particularly with respect to privacy laws like GDPR, and what changes might be necessary to address these challenges?**  
*[Open text box for responses]*
5. **How should liability be determined when a deepfake is detected, especially if the technology is used for harmful purposes like defamation or fraud? Should the developers of the detection system be held accountable?**  
*[Open text box for responses]*
6. **How do you ensure that deepfake detection systems are culturally sensitive and accurate across diverse geographic regions, and what ethical considerations should**

**be made when collecting or curating datasets?**

*[Open text box for responses]*

7. **What do you think are the potential societal benefits and risks of widespread implementation of deepfake detection systems, particularly with regard to public trust in media and democracy?**

*[Open text box for responses]*

---

## **Section 2: Perception of Deepfake Detection Technologies**

8. **How concerned are you about the impact of deepfake detection technologies on individual privacy rights?**

**Scale:**

- (1) Not Concerned
- (2) Slightly Concerned
- (3) Neutral
- (4) Concerned
- (5) Very Concerned

9. **How important do you believe the following factors are for ensuring dataset diversity in deepfake detection systems?**

**Scale:**

- (1) Not Important
- (2) Slightly Important
- (3) Neutral
- (4) Important
- (5) Very Important

10. **How confident are you that deepfake detection technologies will keep pace with the evolution of deepfakes?**

**Scale:**

- (1) Not Confident at All
- (2) Slightly Confident
- (3) Neutral
- (4) Confident
- (5) Very Confident

11. **To what extent do you agree with the following statement about technical robustness in deepfake detection models?**

*"Deepfake detection models should be able to maintain high accuracy even in challenging conditions such as noisy environments, low-resolution videos, and extreme facial poses."*

**Scale:**

- (1) Strongly Disagree
- (2) Disagree
- (3) Neutral
- (4) Agree
- (5) Strongly Agree

**12. How effective do you believe data augmentation is for improving the performance of deepfake detection models?**

*"Data augmentation techniques, such as adding noise, shifting pitch, and modifying visual conditions, significantly improve the generalization of deepfake detection models."*

**Scale:**

- (1) Not Effective
- (2) Slightly Effective
- (3) Neutral
- (4) Effective
- (5) Very Effective