

# Real-time cross-cultural lie detection system via multimodal fusion: microexpression enhancement and adversarial defense for forensics

*Zhi Li*

School of Computer Science, Guangdong University of Foreign Studies South China Business College, Guangzhou, China

13927110317@163.com

---

**Abstract.** Traditional methods, such as polygraphs, suffer from limitations including single-modality vulnerability, cultural bias, and adversarial attacks. This paper presents a novel Multimodal Physiological-Behavioral Fusion System (MPBFS) that integrates seven modalities: microexpressions, speech, text, eye movements, galvanic skin response (GSR), cultural features, and adversarial defense mechanisms. The system employs a Cultural Dimension-Physiological Signal Cross-Validation Module to dynamically weight modalities, reducing cultural bias. The Cascade Lightweight for Edge Computing achieves 84.6% accuracy ( $F1=0.87$ ) on cross-cultural datasets with 23 FPS. A Two-Channel Adversarial Defense: Detects DeepFake audio via phase discontinuity analysis and validates microexpressions using Lagrangian biomechanical modeling confirming enhanced cross-cultural robustness. Experiments demonstrate a 31% improvement in Cultural Stability Index and a 42% reduction in adversarial attack success rates, validating the system's robustness. Designed for forensic interrogation and security screening, MPBFS integrates dynamic anonymization and cross-modal validation to ensure ethical AI deployment, addressing key limitations of traditional deception detection methods.

**Keywords:** lie detection, microexpressions, cross-cultural AI, adversarial defense, multimodal fusion

---

## 1. Introduction

Deception detection plays a crucial role in forensic investigations and border security, yet traditional methods remain limited by several fundamental constraints. While polygraph-based approaches have evolved into AI-driven multimodal systems, three persistent challenges demand urgent attention. First, there is the vulnerability of single-modality systems where trained individuals can suppress microexpressions [1]. Second, cultural biases in detection accuracy due to varying behavioral norms between collectivist (e.g., Japan, where eye contact may be avoided) and individualist societies (e.g., the USA, where direct gaze is valued), which can inflate the false-positive rate by 35-40% [2]. Additionally, there is the growing threat of sophisticated adversarial attacks, such as DeepFake audio that circumvents voice stress analyzers by mimicking genuine physiological patterns [3]. These limitations not only reduce reliability in critical scenarios but also hinder cross-cultural applicability.

In response, this study develops an integrated solution that addresses these gaps through an innovative 7-modality fusion (incorporating microexpressions, speech, text, eye movements, GSR, cultural features, and adversarial defense) with real-time cultural adaption. The framework employs a two-tier adversarial defense against synthetic media, and safeguards including dynamic anonymization and dual-blind protocols. The approach combines advanced behavioral analysis with physiological signal processing while incorporating ethical safeguards for deployment. The research contributes to both academic and practical domains by establishing a more accurate, culturally-sensitive, and secure framework for deception detection in high-stakes environments. By bridging these technological and methodological gaps, the research aims to set new standards for trustworthy lie detection systems that meet the demands of modern security and forensic applications.

## 2. Literature review

### 2.1. A review of multimodal deception detection

The field of deception detection has undergone a significant paradigm shift from unimodal to multimodal approaches. Early techniques primarily relied on single indicators such as polygraph measurements of physiological responses or microexpression analysis, both of which exhibited critical limitations. Polygraph testing, while widely adopted in forensic settings, demonstrates

accuracy rates between 65 and 75% in controlled environments and is particularly vulnerable to countermeasures by trained individuals [1]. Microexpression analysis, though theoretically sound based on Ekman's work, faces practical challenges as subjects can consciously suppress these fleeting facial movements [1]. These shortcomings have driven the development of multimodal systems that integrate complementary data streams including visual cues, vocal characteristics, physiological signals, and linguistic patterns. Recent advancements demonstrate that hybrid architectures combining techniques like 3D-ResNet for visual analysis and ECAPA-TDNN for vocal feature extraction achieve superior performance, with documented accuracy improvements of 12.7% over unimodal baselines [4]. However, current multimodal systems still exhibit notable gaps, particularly in accounting for cultural variations in deceptive behavior and defending against increasingly sophisticated adversarial attacks using synthetic media. Most existing frameworks employ static fusion weights that fail to adapt to cultural contexts, resulting in up to 38% higher false-positive rates when applied cross-culturally [2]. Additionally, while some progress has been made in detecting AI-generated content through methods like phase discontinuity analysis, comprehensive solutions for real-time, cross-modal validation of authenticity remain underdeveloped [5].

## 2.2. Cross-cultural lie detection

The challenge of cross-cultural deception detection is fundamentally rooted in the substantial variations of communication norms and stress responses across different societies. Hofstede's cultural dimensions theory provides a robust framework for quantifying these differences, particularly through metrics like individualism-collectivism and uncertainty avoidance, which directly influence deception behaviors [2]. Research consistently shows that individuals from collectivist cultures tend to suppress facial expressions while exhibiting more pronounced physiological indicators when being deceptive, which is precisely the opposite pattern observed in individualist cultures [6]. This cultural divergence creates significant challenges for lie detection systems, as models trained primarily on Western datasets demonstrate error rates that are 35-40% higher when applied to Asian or African populations. Current approaches typically employ static thresholds and uniform weighting schemes that cannot account for these cultural variations, leading to systematic biases in multicultural applications. The problem is further compounded by the lack of dynamically adaptive mechanisms in existing systems - while some studies have incorporated cultural factors as binary classifiers (e.g., East vs. West), none have implemented continuous cultural dimension scaling [7]. This limitation becomes particularly acute in global security applications where subjects may originate from any cultural background. Recent work has begun exploring cultural adaptation in affective computing, but these advances have not been adequately integrated into deception detection pipelines. The absence of standardized protocols for evaluating cross-cultural performance represents another critical gap, with most studies reporting accuracy metrics only for their primary test population while neglecting validation on culturally diverse samples. These shortcomings highlight the urgent need for deception detection systems that can automatically adjust their analytical parameters based on quantifiable cultural dimensions while maintaining robustness against both genuine cultural differences and potential adversarial exploitation of these variations.

## 3. Methodology

### 3.1. Multimodal data acquisition

Data is collected from 6 cultural zones (Hofstede-stratified) using:

Face/eye tracking: Tobii Pro Glasses 3 (120 Hz).

Voice: Anti-noise RawNet3 embeddings.

Text: Linguistic Inquiry and Word Count (LIWC) for deception markers [8].

GSR: Empatica E4 wristband.

### 3.2. Cultural adapter

The cultural adaptation mechanism is implemented through a dedicated CultureAdapter class that dynamically adjusts modality weights based on Hofstede's individualism-collectivism dimension [2]. The class normalizes the Hofstede score (typically ranging from 0-100) into an individualism index bounded between 0.5 (strong collectivism) and 1.5 (strong individualism) using linear scaling. This normalized index then modulates the relative contribution of visual versus textual features in the fusion pipeline. Visual features are weighted proportionally to the individualism score, while textual features receive inverse weighting [9]. Specifically, collectivist cultures (scores <50) emphasize textual analysis up to twice as much as visual cues, reflecting their reliance on verbal context in deception. Conversely, individualist cultures (scores >50) trigger enhanced visual feature weighting, aligning with their more expressive nonverbal behavior patterns. This dynamic adjustment occurs in real-time during inference, requiring less than 2ms additional processing per sample while reducing cross-cultural false positives by 31% in validation tests. The fully differentiable implementation enables end-to-end training with other network components.

A Python-class CultureAdapter adjusts weights based on Hofstede scores:

```

class CultureAdapter:
    def __init__(self, hofstede_score):

self.                                individualism = 1 + (hofstede_score - 50)/100                                (1)

    # Normalize to [0.5,1.5]
    def adjust_weights(self, visual, text):
    return visual*self.individualism, text/self.individualism
    # Collectivist cultures weight text higher

```

### 3.3. Real-time cascade system

The proposed real-time cascade system features a two-tier architecture for efficient deception detection. The lightweight Tier 1 (MobileNetV3 + Mel-Fbanks) enables rapid 5ms screening at 200 FPS, while Tier 2's advanced analysis (3D-ResNet18 + ECAPA-TDNN) activates when the confidence of any modality exceeds 0.7. This modular design achieves 23 FPS processing on embedded platforms, balancing speed and accuracy for real-world interrogation scenarios.

Tier 1 (Screening): MobileNetV3 + Mel-Fbanks (5 ms latency).

Tier 2 (Deep Analysis): Triggered if any modality confidence >0.7:

Microexpressions: Optical flow + 3D-ResNet18.

Voice: ECAPA-TDNN for anti-spoofing.

### 3.4. Adversarial defense

This study develops a dual-channel adversarial defense mechanism to combat synthetic attacks in both audio and visual modalities by exploiting inherent physical inconsistencies in AI-generated content. For audio defense (Channel 1), the system detects DeepFake artifacts through phase discontinuity analysis in STFT spectrograms, where neural vocoders exhibit characteristic breaks in phase coherence at 55 - 75Hz intervals. The visual defense (Channel 2) employs Lagrangian muscle dynamics modeling to authenticate microexpressions, verifying natural facial biomechanics through acceleration patterns of 26 facial landmarks and thermal diffusion characteristics absent in synthetic media [10]. Operating synergistically, this dual approach achieves 92.3% detection accuracy for audio spoofing (tested on ASVspoof 2021) and reduces visual forgery success rates from 68% to 9.2% (tested on FaceForensics++ benchmark), while maintaining real-time performance with only 8ms added latency through optimized parallel processing of both channels [11].

Channel 1: Detects DeepFake audio via phase discontinuities in STFT spectrograms.

Channel 2: Validates microexpressions using Lagrangian muscle dynamics (non-synthetic motion patterns).

## 4. Experiments

### 4.1. Dataset

Cross-Cultural Balance: 16.7% sample proportion per Hofstede dimension (e.g., high vs. low power distance).

Synthetic Data: StyleGAN-v generates diverse microexpressions for underrepresented groups [12].

### 4.2. Results

MPBFS achieved an overall accuracy of 84.6% (F1=0.87), outperforming Polygraph (62.1%) and OpenFace+LSTM (73.8%). As shown in Table 1, the Cross-Cultural Stability Index (CSI) of MPBFS improved from 0.55 (the average of baseline models) to 0.82, indicating a 31% reduction in cultural bias.

Removing the cultural adaptation module reduced CSI by 31% (to 0.56) and reduced accuracy to 75.2%, highlighting its role in mitigating cross-cultural disparities. The two-channel adversarial defense proved critical against synthetic attacks: disabling it increased DeepFake audio detection failure rates by 42% and reduced microexpression verification accuracy from 89% to 61%. This efficiency enables practical deployment in forensic and security scenarios.

**Table 1.** Performance comparison of different lie - detection models

Model	Accuracy	F1	CSI
Polygraph	62.1%	0.59	0.41
OpenFace+LSTM	73.8%	0.72	0.55
MPBFS	84.6%	0.87	0.82

## 5. Discussion

### 5.1. Applications & ethics

The system implements a three-level output protocol designed to match operational urgency with evidentiary reliability. Low-risk results (<30%) are encrypted and logged for review. Intermediate-risk results (30 - 70%) trigger structured reviews with highlighted behavioral markers (e.g., inconsistent microexpression clusters or voice stress patterns) through an augmented interface for forensic specialists, reducing cognitive load by 42% in user studies. High-risk alerts (>70%) activate real-time notifications with configurable escalation paths—in interrogation settings, this manifests as subtle tablet vibrations for investigators, while airport deployments use secure LED indicators at security stations. Crucially, all outputs retain timestamped, multimodal logs to support subsequent judicial review and expert testimony. This tiered approach demonstrated 89% operational appropriateness in field trials across 12 jurisdictions, effectively balancing automation with necessary human oversight in legally sensitive contexts.

### 5.2. Ethical safeguards

To ensure ethical deployment, the system uses dual-blind mode and real-time dynamic blurring. In the dual-blind mode, analysts are intentionally restricted from accessing raw video footage and instead interact solely with processed feature vectors and analytical outputs, effectively preventing potential biases that could arise from visual judgments of subjects' appearances or environments. Complementing this approach, the system implements real-time dynamic blurring that automatically obscures all non-facial regions in the video feed. These combined safeguards not only protect subjects privacy but also ensure compliance with GDPR and other data protection regulations by minimizing the collection and exposure of extraneous personal data, while maintaining the system's analytical capabilities through focused processing of relevant facial features and behavioral cues [13].

## 6. Conclusion

This study proposes MPBFS, an innovative real-time lie detection system that addresses critical challenges in cross-cultural applicability and adversarial robustness through three key advancements. The system first introduces cultural-physiological fusion by incorporating Hofstede's cultural dimensions into feature weighting, reducing cultural bias by 31% (CSI improvement from 0.55 to 0.82) and overcoming limitations of Western-trained models in collectivist cultures. Second, a cascade lightweight architecture combining MobileNetV3 and 3D-ResNet18 achieves an optimal balance between processing speed (23 FPS) and detection accuracy (84.6%), enabling practical deployment on edge devices. Finally, a two-channel adversarial defense mechanism provides robust protection against synthetic media, detecting DeepFake audio with 92.3% AUC through phase discontinuity analysis while verifying microexpressions using Lagrangian kinematics (F1=0.89 against synthetic faces). Together, these innovations establish MPBFS as an advanced solution for reliable deception detection across diverse cultural contexts and security scenarios.

Despite advancements, three challenges remain in this research. Regarding data diversity issues, while StyleGAN-v generates minority groups, some indigenous cultures (e.g., Maori, Inuit) remain underrepresented due to limited training samples. Hardware dependencies on specialized sensors like the Empatica E4 and Tobii Pro for High-precision GSR and eye-tracking limit scalability in resource-constrained settings. Regarding dynamic deception, the system currently focuses on short-interval lies (<5 minutes), while prolonged deception, such as rehearsed alibis may exploit temporal averaging in physiological signals. Future work will extend temporal modeling with Transformer-based architectures to detect deception patterns across hours-long interrogations.

## References

- [1] Ekman, P. (2009). *Telling lies: Clues to deceit in the marketplace, politics, and marriage* (3rd ed.). W.W. Norton.
- [2] Hofstede, G. (1980). *Culture's consequences: International differences in work-related values*. Sage Publications.

- [3] Zhou, Y., Shi, B. E., & Chen, X. (2021). Deception detection in videos using multi-scale spatial-temporal features. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 13348-13357.
- [4] Wang, L., & Patel, R. (2023). Hybrid transformer architectures for real-world multimodal deception detection. *ACM Transactions on Multimedia Computing, Communications, and Applications*, 19(3), 1-24.
- [5] Chen, X., et al. (2023). Audio-visual deepfake defense through phase discontinuity analysis. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (pp. 15268-15278).
- [6] Zhang, Y., Wang, C., & Liu, H. (2023). Cross-cultural microexpression recognition via deep metric learning. *IEEE Transactions on Affective Computing*, 14(2), 1123-1136.
- [7] Hofstede Insights. (2023). Cultural bias in AI-based deception detection: Global benchmark report 2023. Hofstede Press.
- [8] Zhou, L., Burgoon, J. K., & Twitchell, D. P. (2004). Automated linguistics based cues for detecting deception in text-based communication. *Group Decision and Negotiation*, 13(1), 81-106.
- [9] Global Deception Research Team. (2022). Cross-cultural patterns of deceptive behavior. *Journal of Cross-Cultural Psychology*, 53(4), 389-412.
- [10] Matsumoto, D., Hwang, H. S., & Frank, M. G. (2021). Facial expressions as behavioral markers of deception. In *Handbook of emotions* (4th ed., pp. 211-234). Guilford Press.
- [11] ASVspoof Consortium. (2021). ASVspoof 2021 challenge evaluation report. ISCA Archive.
- [12] ECCV Workshop. (2022). Microexpression grand challenge 2022: Synthetic-to-real generalization (LNCS Vol. 13688). Springer.
- [13] Cox, T., Rogers, H., Simmons, O., & Pum, M. Ethical AI Governance Models in Global Financial Institutions: A Cross-Cultural Analysis.