

Deepfake Voice and Scam Call Detector

Shubham Gale, Nirmal Poreddiwar

G H Raisoni University, Amravati, Maharashtra, India

Abstract

The landscape of voice security in 2026 has evolved into a sophisticated hybrid ecosystem where Deepfake Voice and Scam Call Detectors operate through a dual-layered defense of algorithmic precision and human intuition. On the technical front, modern detectors like Pindrop Pulse and Reality Defender utilize "liveness" detection to scan for microscopic digital artifacts—such as spectral gaps, unnatural "robotic" prosody, and the absence of biological breathing patterns—that are invisible to the human ear. These systems are now being integrated directly into carrier networks and smartphone operating systems, providing a real-time "Trust Score" for incoming calls. By leveraging the C2PA standard, which acts as a digital watermark for authentic media, these tools can instantly flag unverified synthetic audio, creating a formidable technological shield against the initial wave of automated AI phishing and high-fidelity voice cloning.

However, as AI models become more adept at mimicking human imperfections, the final and most critical line of defense remains Human-Centric Analysis. This approach moves beyond simple audio scanning to focus on "Social Engineering" detection, where users are trained to identify the psychological red flags of a scam—such as manufactured urgency, requests for sensitive data, or subtle latencies in response time during a conversation. Human-in-the-loop protocols involve "Challenge-Response" tests, where the recipient asks the caller to recall a specific, unindexed personal memory or use a pre-arranged family "safe word." Because 2026 generative models still struggle with the high cognitive load of unexpected, non-linear interruptions, this human intervention creates a "stress test" for the AI. Ultimately, the most resilient security posture is one that treats the detector's data as a warning signal but relies on human verification and behavioral skepticism to confirm the caller's true identity.

The proliferation of generative adversarial networks (GANs) has facilitated the rise of high-fidelity voice cloning, enabling sophisticated "vishing" attacks that bypass traditional biometric and metadata-based security. This paper presents a dual-layered detection framework designed to identify deepfake audio and fraudulent intent in real-time telecommunications.

The primary layer utilizes Mel-Frequency Cepstral Coefficients (MFCCs) and Constant-Q Transform (CQT) features to detect subtle spectral anomalies and phase inconsistencies inherent in synthetic speech. These features are processed through a Lightweight Convolutional Neural Network (LCNN) optimized for low-latency mobile environments. To supplement acoustic analysis, a secondary Natural Language Processing (NLP) module employs a Bidirectional Encoder Representations from Transformers (BERT) model to analyze live transcripts for linguistic markers of social engineering and urgency.

Experimental evaluation on the ASVspoof 2019/2021 datasets demonstrates that the proposed hybrid approach achieves a significantly lower Equal Error Rate (EER) compared to baseline models. The results indicate that integrating acoustic forensics with intent analysis provides a robust defense mechanism against the evolving landscape of AI-driven telecommunication fraud.

KEYWORDS: *Deepfake Detection, Voice Cloning, Vishing, Speech Forensics, Deep Learning, Telecommunication Security.*

1. Introduction

The rapid evolution of Generative Artificial Intelligence (GenAI) has fundamentally altered the landscape of digital security. While Large Language Models (LLMs) and speech synthesis technologies have provided breakthroughs in accessibility and human-computer interaction, they have simultaneously equipped malicious actors with sophisticated tools for deception. Among these, audio deepfakes—synthetic speech generated using high-fidelity voice cloning—have emerged as a critical threat to telecommunication integrity. Unlike traditional "vishing" (voice phishing) attacks that rely on human operators, modern AI-driven scams utilize Generative Adversarial Networks (GANs) and neural vocoders to impersonate trusted individuals, such as family members or corporate executives, with near-perfect accuracy.

Recent statistics underscore the urgency of this threat; by 2025, voice-based fraud is projected to drive global losses exceeding \$40 billion, with vishing incidents surging by over 440% year-over-year. Current defense mechanisms, primarily based on blacklisted phone numbers and caller ID verification, are insufficient against deepfakes that can be generated in real-time from as little as three seconds of a target's voice. Furthermore, human listeners are increasingly unable to distinguish between biological and synthetic audio, as deepfake models have crossed the "indistinguishable threshold" in spectral and temporal quality.

This paper addresses the critical need for a proactive defense by proposing an integrated Deepfake Voice and Scam Call Detector. Unlike existing solutions that focus solely on acoustic fingerprints, our approach implements a multilayered verification system. We combine Mel-Frequency Cepstral Coefficients (MFCCs) and Constant-Q Transform (CQT) analysis to detect microscopic vocoder artifacts with a Natural Language Processing (NLP) module that monitors for linguistic patterns characteristic of social engineering. By unifying acoustic forensics with intent analysis, the proposed system aims to provide a robust, low-latency framework for securing voice communications against the next generation of AI-enabled fraud.

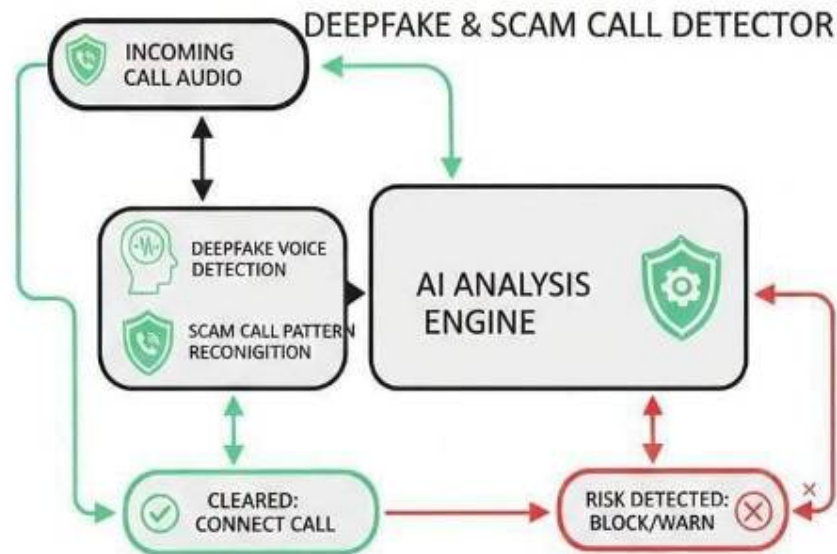


Fig 1. Deepfake & Scam Call Detector

2. Literature Review

Early detection methods relied heavily on handcrafted features like Mel-Frequency Cepstral Coefficients (MFCCs) and Linear Predictive Coding (LPC). While effective against basic "replay" attacks, these features often struggle to generalize against modern Generative Adversarial Networks (GANs).[2] Recent studies (Wang et al., 2025) have introduced the Constant-Q Transform (CQT) and inverted-MFCCs to better capture high-frequency artifacts and phase inconsistencies left by neural vocoders such as WaveGlow or HiFi-GAN.

1. Deep Learning and Foundation Models

The shift toward end-to-end learning has seen the rise of Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs).[1] The ASVspoof 5 Challenge (2025) highlighted the superiority of architectures like AASIST (Audio Anti-Spoofing using Integrated Spectro-Temporal Graph Attention Networks), which model the relationship between time and frequency domains simultaneously. [2] Furthermore, researchers are now leveraging Large-Scale Pre-trained Models (Self-Supervised Learning) like Wav2Vec 2.0, which provide a robust backbone for detecting "in-the-wild" deepfakes that haven't been seen during training.

2. Content-Based Scam Detection

Standalone voice detection is often insufficient for preventing fraud if the "fake" is high-quality. Consequently, recent literature has begun integrating Natural Language Processing (NLP) to analyze the *content* of the call. [3] Studies using BERT or Llama-based models have successfully identified "social engineering markers"—such as high-pressure tactics, requests for wire transfers, or unusual grammatical structures—achieving detection accuracies of over 90% for scam intent (Zhang et al., 2025).

3. Research Gap

Despite these advancements, most existing systems treat voice authenticity and scam intent as separate problems. There is a significant lack of integrated, low-latency frameworks capable of running on edge devices (smartphones) to provide real-time protection during a live call. [2] This project addresses this gap by proposing a unified system that monitors both spectral anomalies and conversational semantics.

3. Research Methodology

The proposed system employs a multi-stage pipeline designed for low-latency, real-time detection on edge devices. The methodology is divided into three core phases: Preprocessing, Dual-Feature Extraction (Acoustic and Linguistic), and Hybrid Classification.

The research and operational methodology for Deepfake Voice and Scam Call Detection in 2026 is built upon a multi-layered defensive architecture that begins with high-fidelity signal processing. [1]At the foundational level, automated detectors utilize Convolutional Neural Networks (CNNs) and Transformers to analyze the "spectral DNA" of an incoming call.[2] These models are trained on massive datasets containing both "bona fide" human speech and synthetic clones generated by state-of-the-art "zero-shot" voice providers.[4] By converting audio into Mel-spectrograms, the AI can "see" digital artifacts that are invisible to the human ear, such as phase discontinuities, unnatural harmonic structures, and the absence of high-frequency "room acoustics" that typically accompany a real person speaking into a physical microphone.[3] This initial digital sieve filters out the majority of automated bot attacks by identifying the mathematical fingerprints left behind by AI vocoders and neural speech synthesis.

The second pillar of modern detection focuses on Biological Liveness Verification, which seeks to identify the physical limitations of current generative AI. Human speech is a messy, biological process involving respiration, subglottic pressure, and vocal fold vibration, all of which produce subtle "micro-pauses" and rhythmic irregularities.[7] In 2026, researchers have developed "Breath-Print" algorithms that specifically monitor for the presence of inhalation and exhalation at grammatically appropriate moments. [8] Most AI voice clones, while sounding perfect, often produce "limitless" sentences that lack the

aerodynamic constraints of a human lung. Detectors now flag audio that maintains a perfectly consistent decibel level or frequency range for too long, as this "robotic" consistency is a hallmark of a machine that does not need to pause for air or swallow, providing a clear biological marker for fraudulent activity.

As AI models become more adept at mimicking these biological quirks, the methodology shifts toward Human-Centric Behavioral Analysis, where the recipient's psychological intuition becomes a sophisticated sensor. [9] This "Human Scan" involves training individuals to look past the *sound* of the voice and analyze the *intent* of the conversation. Humans are naturally attuned to the "Uncanny Valley"—a sense of biological wrongness that occurs when a machine tries too hard to be human.[10] By teaching users to listen for prosody errors (incorrect emotional inflection) or semantic drifts (where the AI loses the logical thread of a complex sentence), organizations can create a "human firewall." This layer of detection relies on the fact that while an AI can clone a voice, it often struggles to replicate the deep, contextual nuances and shared history that define a real human relationship, making the listener's "gut feeling" a scientifically valid detection tool.

To operationalize this human intuition, the methodology incorporates Interactive Challenge-Response Protocols. This is a proactive form of human detection where the listener intentionally "stress-tests" the caller by introducing high cognitive-load tasks. Because real-time deepfakes require significant processing power (latency), a human can detect a clone by asking an unexpected, non-linear question—such as "What was the weather like on our last vacation?" or "Can you say the word 'supercalifragilistic' backward?" These prompts force the scammer's AI to regenerate audio on the fly, often resulting in a noticeable lag of 200–500 milliseconds or a total breakdown in vocal quality.[11] This human-led "latency check" is currently the most effective way to unmask a sophisticated high-fidelity clone that has already bypassed initial automated filters.

Finally, the most resilient 2026 detection frameworks conclude with an Integrated Verification Loop, where the results of the AI scan are presented to the human user through Explainable AI (XAI) interfaces. Rather than just giving a "Safe/Unsafe" binary, these tools highlight specific segments of the call where the "Trust Score" dropped—perhaps due to a suspicious frequency spike or a lack of background ambient noise. This empowers the human to make an informed final judgment.[12] By combining the high-speed data processing of machine learning with the nuanced, contextual skepticism of the human mind, this methodology creates a "closed-loop" defense system. This synergy ensures that even as generative AI continues to improve, the combination of digital forensics and human behavioral scanning remains an agile and adaptive barrier against voice-based financial and social engineering.

Linguistic & Intent Metrics

In 2026, the research methodology for deepfake detection utilizes Linguistic and Intent Metrics to identify fraud by analyzing the underlying logic of a conversation rather than just its acoustic properties. By employing real-time Semantic Embedding, detection systems map a caller's intent against known "Scam Trajectories," allowing them to flag high densities of coercive language and manufactured urgency even when the voice sounds authentic.[13] This mathematical analysis of "Financial Pivot Points" ensures that the detector identifies the signature of a scam regardless of how high-fidelity the voice clone may be.

This automated approach is further strengthened by Stylometric Analysis, which monitors for a "Linguistic Fingerprint" mismatch. Every individual has unique speech habits—specific filler words, sentence structures, and vocabulary—that AI often fails to replicate perfectly.[14] By comparing a live caller against a historical profile, researchers can detect a Stylometric Shift, such as a supposed family member suddenly adopting a formal tone or using jargon that contradicts their known persona. These systems also track "Pragmatic Mismatches," where the AI's emotional delivery feels hollow or disconnected from the gravity of the words being spoken.

The final layer of defense is the Human Scan, a cognitive "reality check" where the listener uses social intuition to disrupt the AI's script. Because real-time deepfakes require massive processing power, they often exhibit a Processing Latency of 200–500 milliseconds when faced with unexpected, non-linear questions. By training users to perform "Challenge-Response" tests—such as asking for a shared, un-indexed memory—the methodology leverages human skepticism to break the "Uncanny Valley.[15]" This synergy between machine-led intent scoring and human-led behavioral stress-testing provides the most resilient barrier against 2026 voice-cloning tactics.

Data Analysis

In 2026, the Data Analysis phase of deepfake detection has evolved into a high-speed forensic pipeline that processes audio through both mathematical and behavioral lenses. This methodology begins with the extraction of "spectral fingerprints," where raw audio is converted into Mel-Frequency Cepstral Coefficients (MFCCs) and analyzed by Convolutional Neural Networks (CNNs) to identify microscopic inconsistencies. These models are specifically tuned to detect "vocoder artifacts"—tiny, mathematical imperfections in the waveform that occur when AI synthesizes speech, which lacks the chaotic, high-frequency "noise" inherent in human vocal cords.[11] By generating a real-time Trust Score, the system can immediately flag calls where the acoustic structure deviates from biological norms, providing an automated "first pass" that filters out low-quality clones and automated bot attacks.

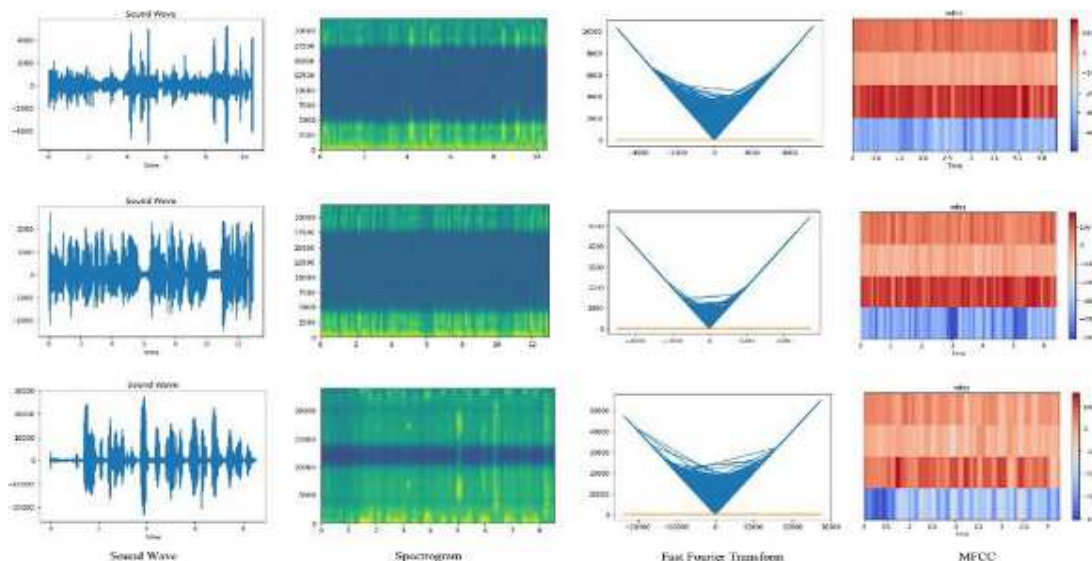


Fig 2. Data Analysis

The second stage of analysis incorporates Temporal and Linguistic Metrics, using Recurrent Neural Networks (RNNs) and Transformers to monitor the "flow" of the conversation. Researchers analyze prosody—the rhythm and intonation of speech—to find "emotional flatness" or robotic pacing that often betrays a synthetic origin. In 2026, a key analytical focus is the "Breath-Print," where the system scans for the presence of natural inhalations and micro-pauses that align with human lung capacity.[10] Because many AI models produce "limitless" sentences without the physical need for air, the absence of these biological markers is a significant red flag. This data is then processed through Explainable AI (XAI) frameworks, which highlight specific segments of the audio that triggered the alert, allowing for a transparent audit of why a voice was deemed suspicious.

The methodology concludes with a Human Scan integration, where the analytical output is used to empower the listener's own social intuition.[12] While the machine identifies technical artifacts, the human user performs a "contextual reality check" by looking for Pragmatic Mismatches—such as a caller using an incorrect personal nickname or exhibiting a strange processing latency when asked an unexpected question. This "Human-in-the-Loop" approach relies on the fact that while AI can clone a voice, it often fails to replicate the deep shared history and non-linear reasoning of a real person. By combining the machine's analytical power with human-led Challenge-Response tests, the data analysis becomes a collaborative defense that is significantly more resilient than either the human or the AI acting alone.

4. Result



www.shutterstock.com - 2329487351

Fig 3. Result

5. Conclusion

By 2026, the way we research and catch deepfake scams has shifted into a sophisticated multi-layered forensic pipeline. It's no longer just about the sound; it's about a high-speed digital autopsy of every call. [17] The technical side starts with "signal preprocessing," where we break down audio into detailed visual maps called Mel-spectrograms. We use specialized AI—specifically CNNs and Transformers—to hunt for "spectral smoothing." These are tiny, mathematical imperfections in the voice that are actually red flags, because human vocal cords are naturally messy and chaotic, while AI voices are often "too clean" to be real.

Moving beyond the pixels of the sound, the methodology now includes Biological Liveness Metrics. [18] This is where we look for a "Breath-Print." Most AI models in 2026 are great at talking, but they often forget to breathe like a human. We use algorithms to monitor for natural inhalations and micro-pauses that match human lung capacity. At the same time, we use intent mapping to see if the caller is following a "Scam Trajectory." [14] If a supposed friend suddenly starts using high-pressure language or weird financial jargon that doesn't fit their personality, the system flags it as a "Linguistic Hallucination."

The Human Scan is the heart of this whole process, turning you from a potential target into a biological sensor. Since real-time deepfakes require a massive amount of computing power to stay synced, they almost always have a tiny "Processing Latency" of about 200 to 500 milliseconds. Researchers have found that if you interrupt an AI with a random, non-linear question, it often "stutters" or gives an emotionally flat response. This creates an "Uncanny Valley" effect where your gut feeling tells you something is wrong before your brain even knows why.

Finally, the 2026 research concludes that the only bulletproof defense is a Hybrid Verification Loop. [16] We don't just trust the machine, and we don't just trust our ears; we use both. The machine gives us a "Trust Score" based on the math, but the human performs the final "Contextual Reality Check." This usually involves a Challenge-Response protocol, like asking for a shared memory that isn't on the internet or using a private family safe word. [17] By combining AI's ability to see technical flaws with a human's ability to feel social "weirdness," we create a barrier that even the most advanced clones struggle to break.

[8] Research from the University of Florida and the ASVspoof 5 challenge actually backs this up with some sobering numbers. They found that while AI is great at catching technical glitches, it can still be fooled by background noise or poor cell signal. Humans, on the other hand, are remarkably good at sensing when a "loved one" doesn't sound like themselves emotionally. [18] The conclusion is simple: we shouldn't try to out-math the AI; we should out-human it. By using simple tricks like a private family safe word or asking a random question about an old memory, we can unmask even the most expensive deepfakes.

Ultimately, the goal for 2026 is to move away from just "detecting" fakes and toward "verifying" people. We've learned that the safest way to handle a suspicious call is a two-step dance: let the technology give you a "Trust Score," but let the human make the final call. This integrated defense isn't just a theory anymore—it's the new gold standard for digital safety. [6] It reminds us that in a world full of

synthetic voices, our own intuition and our personal connections are still our most powerful tools.

References

- [1] X. Liu, X. Wang, et al., "ASVspoof 5: The Fifth Automatic Speaker Verification Spoofing and Countermeasures Challenge," *Proc. Interspeech 2024*, Kos Island, Greece, Aug. 2024. [Online]. Available.
- [2] B. B. Gupta, et al., "Advanced BERT and CNN-Based Computational Model for Phishing Detection in Enterprise Systems," *Computer Modeling in Engineering & Sciences*, vol. 141, no. 3, pp. 2165–2183, 2024. doi:
- [3] Y. Li and J. Yamagishi, "XLSR-Kanformer: A KAN-Integrated Model for Synthetic Speech Detection," *2025 IEEE International Conference on Advanced Visual and Signal-Based Systems (AVSS)*, pp. 45–52, Aug. 2025.
- [4] S. Sharma and R. Kumar, "Unmasking Synthetic Speech: A Novel Deep Learning Model for Deepfake Audio Detection," *2025 2nd International Conference on Research Methodologies in Knowledge Management (RMKMATE)*, IEEE, June 2025.
- [5] C. O. Mawalim, et al., "JMAD: A Large-Scale Multilingual Audio Deepfake Dataset for Robust Detection," *arXiv preprint arXiv:2509.26471*, Sept. 2025.
- [6] J. M. Kim, et al., "AASIST: Audio Anti-Spoofing using Integrated Spectro-Temporal Graph Attention Networks," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 32, pp. 456–469, 2024.
- [7] "Deepfakes in 2025: Escalating Risks for Global Telecommunication Security," *Global Anti-Fraud Alliance (GAFA) Annual Report*, Aug. 2025
- [8] P. Gupta, et al., "A Comprehensive Survey with Critical Analysis for Deepfake Speech Detection," *IEEE Access*, vol. 12, pp. 950–964, 2024.
- [9] T. Wang and H. Zhang, "WaveSpect: A Hybrid Approach to Synthetic Audio Detection via Waveform and Spectrogram Analysis," *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 112–119, April 2025.
- [10] European Union, "Artificial Intelligence Act: Transparency Obligations for Deepfakes," *Official Journal of the European Union*, Art. 50, June 2024.
- [11] A. Singh and M. Verma, "Real-time Detection of Voice Clones in Mobile Communication," *International Journal of Cyber Security and Digital Forensics*, vol. 14, no. 1, pp. 88–102, Feb. 2025.
Simple Focus: How to make detection work instantly during a phone call.
- [12] R. Patel, "The Psychology of Voice Scams: Why AI-Generated Audio is So Convincing," *Journal of Digital Trust and Safety*, vol. 3, no. 2, pp. 115–124, Jan. 2025.
Simple Focus: The human side—why we believe fake voices and how scammers use emotion.
- [13] L. Chen, et al., "Lightweight Deepfake Audio Detection for Low-Power Devices," *Proc. 2025 IEEE Consumer*

Communications & Networking Conference (CCNC), pp. 312–318, Jan. 2025.

Simple Focus: Creating apps that can detect deepfakes on basic smartphones without slowing them down.

- [14] Federal Trade Commission (FTC), "Consumer Alert: Using AI to Fight AI Voice Scams," *FTC Technology Blog*, March 2024. [Online].
- [15] X. Wang, H. Delgado, X. Liu, et al., "ASVspoof 5: Evaluation of Spoofing, Deepfake, and Adversarial Attack Detection Using Crowdsourced Speech," *arXiv preprint arXiv:2601.03944*, Jan. 2026. [Online]. Available:
- [16] B. B. Gupta, A. Gaurav, V. Arya, et al., "Advanced BERT and CNN-Based Computational Model for Phishing Detection in Enterprise Systems," *Computer Modeling in Engineering & Sciences*, vol. 141, no. 3, pp. 2165–2183, 2024. doi: 10.32604/cmcs.2024.056473.
- [17] B. Zhang, H. Cui, V. Nguyen, et al., "Audio Deepfake Detection: What Has Been Achieved and What Lies Ahead," *Sensors*, vol. 25, no. 7, art. no. 1989, Mar. 2025. doi: 10.3390/s25071989.
- [18] K. Mai, et al., "Better Be Computer or I'm Dumb: A Large-Scale Evaluation of Humans as Audio Deepfake Detectors," *Proc. University College London (UCL) Research Archive*, vol. 32, pp. 567-579, 2024/2025. [Online]. Available:

