

Secure Communication Platform With Aiml-Based Threat Detection

Adithya Varshan A¹, Ashvanthgopal Baskaran², Ranjith V³, Dr. V. Ravindra Krishna Chandar⁴

^{1,2,3,4} Dept of Cybers ecurity

⁵HOD, Dept of Cybers ecurity

^{1,2,3,4,5} Dhanalakshmi Srinivasan University Samayapuram, Tiruchirappalli – 621112, Tamil Nadu, India

Abstract- *The increasing prevalence of spam and scam calls poses a significant security threat, particularly in multilingual regions such as India. This paper presents ShieldCall AI, an AI-driven spam call detection system that integrates Machine Learning, Natural Language Processing (NLP), and real-time audio processing to identify fraudulent calls across eight major Indian languages. The system employs a Multinomial Naive Bayes classifier with TF-IDF feature extraction to classify call transcripts as spam or legitimate. Live audio is transcribed via Google Speech Recognition, while language detection and translation modules support Hindi, Tamil, Telugu, Kannada, Malayalam, Bengali, Marathi, and Gujarati by translating regional language transcripts to English prior to classification. The backend is implemented using Python and Flask, and the frontend is a mobile-first web application deployed on Render. A distinctive feature is the Live AI Screener, an autonomous conversational agent that interacts with suspected spam callers using voice synthesis, collects evidence, and dynamically escalates risk scores. Testing across all supported languages achieved high spam detection accuracy with low false-positive rates, demonstrating the practical effectiveness of the system for real-world personal and organizational security applications*

Keywords: spam call detection, machine learning, Naive Bayes, TF-IDF, multilingual NLP, speech-to-text, AI screener, cybersecurity, Indian languages

I. INTRODUCTION

In the modern telecommunications landscape, spam and scam calls have emerged as a pervasive and growing threat. With billions of robocalls placed annually worldwide, individuals—especially in multilingual developing nations like India—face constant harassment, financial fraud, and privacy violations through deceptive phone calls. Conventional call-blocking solutions rely on static blacklists or simple keyword filters that are easily circumvented by attackers who adapt their scripts, spoof caller IDs, or switch languages mid-call.

To address these challenges, this paper presents ShieldCall AI, an AI-Based Spam Call Detection System

designed to intelligently detect, analyse, and block spam calls in real time. The system integrates a machine learning classification pipeline with a multilingual speech processing backend, enabling accurate spam detection across eight major Indian languages: Hindi, Tamil, Telugu, Kannada, Malayalam, Bengali, Marathi, and Gujarati.

At its core, ShieldCall AI employs a Naive Bayes classifier trained on TF-IDF features extracted from a curated multilingual call transcript corpus. Audio from live or simulated calls is captured and transcribed using Google Speech Recognition, then passed through language detection and translation modules before classification. The risk score, intent label, and chain-of-custody log are recorded for every call, providing a transparent and auditable record of all screening events.

The system further incorporates a Live AI Screener—an autonomous agent that proactively engages suspected spam callers using voice synthesis, collecting evidence and dynamically adjusting risk scores based on detected threat signals. This paper describes the design, implementation, testing, and results of the ShieldCall AI platform, demonstrating its effectiveness and scalability for real-world deployment in India and similar multilingual regions.

II. PROBLEM STATEMENT

Despite significant advances in cybersecurity, spam and scam calls remain a largely unsolved problem, particularly in multilingual environments. The scale of the problem in India is significant: hundreds of millions of mobile subscribers receive unsolicited commercial calls and fraudulent vishing attempts daily. Current systems face several critical shortcomings that this work directly addresses:

- **Language Barrier:** Most commercial spam filters are trained on English-language datasets and fail to detect scam calls conducted in regional Indian languages.
- **Static Blacklists:** Traditional call-blocking relies on known spam numbers, which are easily bypassed through caller ID spoofing and number rotation.

- **No Real-Time Audio Analysis:** Existing tools typically analyse calls post-hoc rather than screening them in real time at the moment of the call.
- **Lack of AI-Driven Interception:** No widely available system can autonomously engage a suspected spam caller on behalf of the user, collecting forensic evidence in the process.
- **Absence of Chain-of-Custody Records:** Call screening events are rarely logged in a tamper-evident, auditable format suitable for legal proceedings or regulatory review.
- **Poor Accessibility:** Existing solutions lack user-friendly interfaces optimised for mobile users in rural or low-bandwidth settings across India.

ShieldCall AI is designed specifically to overcome each of these limitations through a unified, deployable platform that combines machine learning, multilingual NLP, autonomous AI screening, and forensic logging.

III. OBJECTIVES

The primary objectives of the ShieldCall AI project are as follows:

- Train and deploy a Naive Bayes + TF-IDF classifier capable of distinguishing spam and legitimate calls with high precision across multiple languages.
- Integrate speech-to-text, language detection, and translation pipelines supporting eight major Indian languages.
- Capture and analyse live call audio in real time, providing immediate risk assessment and intent classification.
- Deploy an autonomous AI screener capable of conversing with suspected spam callers, gathering evidence, and escalating risk scores dynamically.
- Maintain a detailed, timestamped chain-of-custody log of all call screening events for auditing and legal purposes.
- Deliver a mobile-first, intuitive web dashboard for real-time call monitoring and analytics.
- Build a Flask-based REST API architecture capable of supporting multiple concurrent users in a production environment.

IV. LITERATURE SURVEY

A systematic review of existing research was conducted to establish the theoretical and empirical foundations of ShieldCall AI. The survey spans four key areas: spam classification algorithms, multilingual NLP

techniques, real-time processing requirements, and vishing threat trends.

A. Naive Bayes Effectiveness

Multiple studies confirm that Naive Bayes combined with TF-IDF is highly effective for text-based spam classification, offering a strong balance of accuracy and computational efficiency [1][2]. The Multinomial Naive Bayes variant is particularly suited to high-dimensional text feature spaces produced by TF-IDF vectorisation, making it an ideal choice for a system that must operate under tight latency constraints in a live call-screening context.

B. Multilingual NLP Challenges

Research highlights the need for language-agnostic preprocessing pipelines, with translation-based approaches showing strong generalisation across Indian languages [3]. The diversity of scripts, grammatical structures, and code-switching patterns in Indian languages makes direct multilingual classification non-trivial; a pipeline combining language identification and neural machine translation provides the most robust generalisation without requiring per-language labelled training corpora.

C. Real-Time Processing Requirements

Studies emphasise that sub-second classification response time is critical for live call interception, achievable with lightweight ML models [4]. Traditional deep learning approaches, while more expressive, introduce inference latencies incompatible with real-time call screening. The Naive Bayes + TF-IDF pipeline satisfies this constraint with inference times consistently below 80 milliseconds in testing.

D. Voice Phishing Trends

Recent reports document rapid growth in AI-generated vishing attacks, underscoring urgency for automated multilingual detection [5]. The emergence of AI-synthesised voices and script-generating large language models has dramatically lowered the barrier to entry for sophisticated phone scams, making automated real-time detection a national security priority in densely connected mobile markets such as India.

E. Research Gaps

Integration of autonomous AI screeners for live call interception and multilingual support in a single unified platform remains underexplored in existing literature. No

reviewed system combines real-time multilingual speech processing, autonomous conversational interception, risk scoring, and forensic logging within a single deployable application. ShieldCall AI is specifically designed to fill this gap.

V. SYSTEM DESIGN AND METHODOLOGY

A. System Architecture

ShieldCall AI adopts a five-layer client-server architecture designed for scalability, modularity, and deployment on commodity cloud infrastructure. Each layer is independently testable and replaceable, enabling future extension without architectural changes.

- **Frontend Layer:** A mobile-first HTML/CSS/JS dashboard communicating with the backend via AJAX and Fetch API over HTTPS.
- **Backend API Layer:** A Flask REST API handling call simulation, transcript analysis, history management, and AI screener orchestration.
- **ML Inference Layer:** Loads the trained Naive Bayes + TF-IDF model at startup for synchronous, sub-100ms classification.
- **Speech Processing Layer:** Integrates Google STT, langdetect, and googletrans for multilingual audio-to-classified-text conversion.
- **Storage Layer:** Server-side JSON/SQLite persistence for call history and session state.

B. Module Design

The system is decomposed into four principal modules, each encapsulating a distinct processing responsibility:

Module 1 – Speech Processing: Live microphone audio is captured via the Web Speech API, transcribed using Google Speech Recognition, and the resulting text is passed to the language detection module. The langdetect library identifies the language, and if non-English, the googletrans API translates the transcript to English before classification. This design enables a single English-trained classifier to serve all eight supported languages without retraining.

Module 2 – ML Classification: A Multinomial Naive Bayes classifier trained on TF-IDF features from a curated multilingual call transcript dataset performs binary classification (Spam / Legitimate). A numerical risk score (0–100) is computed based on classifier confidence and heuristic keyword signals including urgency language, OTP requests, payment demands, and official-organisation impersonation.

Module 3 – Live AI Screener: An autonomous conversational agent uses voice synthesis (TTS) to engage suspected spam callers. The screener maintains per-call conversation history, collects evidence via multi-turn dialogue, and dynamically escalates the risk score on detection of urgency language, OTP requests, or payment demands.

Module 4 – Web Dashboard: A mobile-first HTML/CSS/JS security dashboard provides real-time call monitoring, call history retrieval, analytics visualisations, and settings management. Role-based UI rendering simulates Investigator, Analyst, and Admin permission tiers.

C. Key Implementation Details

Backend (Flask): RESTful endpoints — POST /analyze for transcript classification, GET /history for call log retrieval, POST /simulate for call simulation, POST /screener for AI screener session management. CORS-enabled for frontend-backend communication. Deployed on Render with Gunicorn WSGI.

Risk Engine: A heuristic risk engine assigns incremental scores for urgency language (+15), OTP/sensitive data requests (+30), payment demands (+25), known scam patterns (+35), official organisation claims on mobile numbers (+45), and caller not in contacts (+10). Positive signals reduce the score. The final score is clamped to 0–100 and mapped to Low / Medium / High risk levels.

Multilingual Pipeline: Google Speech Recognition → langdetect → googletrans. Supported languages: Hindi (hi), Tamil (ta), Telugu (te), Kannada (kn), Malayalam (ml), Bengali (bn), Marathi (mr), Gujarati (gu), and English (en).

VI. TESTING

A. Testing Strategy

A comprehensive multi-tier testing strategy was employed to validate all aspects of the ShieldCall AI system. Unit testing was performed on each module independently using pytest and mock data. Integration testing validated end-to-end call simulation, multilingual transcription, AI screener sessions, and history retrieval. System testing was conducted on the fully deployed Render application with real call simulations across all supported languages. Security testing probed API endpoints for authentication bypasses, injection vulnerabilities, and unauthorised data access. Performance testing measured API latency and classification throughput under concurrent user load.

B. Functional Test Cases

Table 1: Functional Test Cases and Results

Test Case	Expected Result	Status
Spam call simulation (English)	Risk > 70, SPAM, blocked	PASS
Legitimate call (English)	Risk < 30, LEGITIMATE, forwarded	PASS
Hindi spam transcript	Correct translation, classified SPAM	PASS
Tamil legitimate call	Correct translation, LEGITIMATE	PASS
AI Screener session init	Screener greets caller via TTS	PASS
Unauthorised API access	401 error returned	PASS

C. Performance Results

Table 2: Performance Metrics on Deployed System

Metric	Observed Result
ML Inference (text input)	< 80 ms per request
Full Pipeline (audio to result)	< 2.5 s (LAN), < 3.5 s (mobile)
API Response (history query)	< 300 ms
Multilingual Translation	< 1 s per transcript
System Uptime (30 days)	99%+ on Render

VII. RESULTS

The deployed ShieldCall AI system successfully identifies spam calls in real time across all eight supported Indian languages. In testing, the system correctly classified English spam calls with risk scores above 70 and legitimate calls with scores below 30, satisfying the binary classification requirement with a clear decision boundary.

Hindi, Tamil, and Telugu transcripts were correctly translated and classified with no statistically significant degradation in accuracy compared to English inputs, validating the translate-then-classify pipeline architecture. The language detection module correctly identified all nine supported language codes across 200 test transcripts with zero misclassifications.

The Live AI Screener demonstrated its ability to autonomously intercept suspected spam callers using voice synthesis, asking probe questions and dynamically escalating risk scores on detection of urgency language, OTP requests, or payment demands. In one representative simulation, a caller exhibiting 'TRAI disconnection' and 'arrest warrant' patterns received a risk score of 100/100 and was automatically blocked.

The web-based Security Dashboard provided real-time risk display, call history logging, and analytics across all tested scenarios. API response times remained consistently under 300 ms for history queries and under 80 ms for ML inference, satisfying the sub-second performance requirement for live call screening.

The system is publicly accessible at <https://shieldcall-ai-cww7.onrender.com> and has demonstrated stable operation under concurrent user load during the evaluation period, with 99%+ uptime recorded over 30 days of continuous deployment.

VIII. CONCLUSION AND FUTURE ENHANCEMENTS

A. Conclusion

ShieldCall AI represents a significant advancement in AI-driven spam call protection for multilingual populations. The system combines a Naive Bayes + TF-IDF classification engine with real-time multilingual speech processing, an autonomous AI screener, and forensic logging to overcome the limitations of existing spam detection solutions. It performs full call screening within three seconds across eight Indian languages while maintaining high classification accuracy.

The platform's client-server architecture with REST API design enables straightforward integration with third-party telephony systems and institutional security workflows. By combining machine learning, multilingual NLP, autonomous conversational AI, and a mobile-first dashboard into a single deployable system, ShieldCall AI demonstrates the practical feasibility of intelligent, scalable cybersecurity solutions for real-world deployment in India and similarly diverse linguistic environments.

The technical contributions of this work are: (i) a production-ready multilingual spam call detection pipeline achieving sub-3-second end-to-end latency across eight Indian languages; (ii) a novel autonomous AI screener for live caller interception with dynamic risk escalation; and (iii) a forensic chain-of-custody logging mechanism suitable for regulatory and legal contexts.

B. Future Enhancements

- **Telephony API Integration:** Connect ShieldCall AI to Twilio or Exotel for real PSTN call interception without reliance on web audio capture.
- **Blockchain-Based Evidence Logging:** Integrate Ethereum smart contracts for immutable, court-admissible forensic records of all screening events.
- **Expanded Language Support:** Extend multilingual pipeline to Odia, Punjabi, Assamese, and Urdu, covering the remaining major Indian languages.
- **On-Device ML Inference:** Develop a TensorFlow Lite model for local inference without cloud API dependency, enabling offline operation in low-bandwidth settings.
- **Mobile Application:** Build native Android and iOS apps with direct telephony integration for seamless real-time call protection.
- **Federated Learning:** Implement privacy-preserving federated model updates with community-contributed spam samples to continuously improve detection accuracy.

X. APPENCIX

Source Code:

Simulator_api.py

```

"""
Simulator API router — POST /api/simulate and
/api/simulate/auto
"""
import asyncio
from fastapi import APIRouter, BackgroundTasks,
HTTPException
from pydantic import BaseModel
from typing import Optional
from services.simulator import simulate_call, get_all_profiles
from ws.ws_manager import manager

router = APIRouter(prefix="/api", tags=["simulator"])

# Track the active simulation
_active_simulation = {"running": False, "call_id": None}

```

```

class SimulateRequest(BaseModel):
    profile_id: Optional[int] = None
class AutoSimulateRequest(BaseModel)
    enabled: bool

interval_seconds: int = 30

@router.post("/simulate")
async def trigger_simulation(
    request: SimulateRequest,
    background_tasks: BackgroundTasks
):
    """Trigger a single simulated call and stream it to the
    dashboard."""
    if _active_simulation["running"]:
        raise HTTPException(status_code=409, detail="A
simulation is already running. Please wait.")

    _active_simulation["running"] = True
    _active_simulation["call_id"] = None

    async def run():
        try:
            result = await simulate_call(
                profile_id=request.profile_id,
                ws_callback=manager.broadcast
            )

            _active_simulation["call_id"] = result.get("id")
        except Exception as e:
            await manager.broadcast("error", {"message": str(e)})
        finally:
            _active_simulation["running"] = False

    background_tasks.add_task(run)
    return {"status": "started", "message": "Simulation started.
Watch the Live tab."}

@router.get("/simulate/status")
async def simulation_status():
    """Get the status of the current simulation."""
    return {
        "running": _active_simulation["running"],
        "call_id": _active_simulation["call_id"],
    }

@router.get("/simulate/profiles")
async def get_profiles():
    """Get all available caller profiles for the simulator."""
    return {"profiles": get_all_profiles()}
@router.get("/ollama/status")

```

```

async def ollama_status():
    """Check if Ollama is available."""
    from services.ai_agent import check_ollama_available
    available = await check_ollama_available()
    return {
        "available": available,
        "status": "Connected" if available else "Offline (using
fallback mock responses)",
        "model": "llama3" if available else "mock"
    }

```

main.py

```

"""
FastAPI application entry point for ShieldCall AI.
"""
from fastapi import FastAPI
from fastapi.staticfiles import StaticFiles
from fastapi.responses import FileResponse
from database.db import init_db
from routers import websocket, simulator_api, calls_api,
settings_api

app = FastAPI(
    title="ShieldCall AI",
    description="AI-powered call screening and authentication
system",
    version="1.0.0",
)

# Initialize database on startup
@app.on_event("startup")
async def startup_event():
    init_db()
# Mount API routers
app.include_router(websocket.router)
app.include_router(simulator_api.router)
app.include_router(calls_api.router)
app.include_router(settings_api.router)

# Serve static dashboard
app.mount("/static", StaticFiles(directory="static"),
name="static")

@app.get("/")
async def serve_dashboard():
    return FileResponse("static/index.html")

@app.get("/health")
async def health():
    return {"status": "ok", "service": "ShieldCall AI"}

```

REFERENCES

- [1] A. K. Saha, S. Nandi, "Cyber Threat Detection Using Machine Learning Techniques: A Performance Evaluation Perspective," IEEE Xplore.
- [2] scikit-learn, "Naive Bayes," scikit-learn.org, 2024.
- [3] Google Cloud, "Speech-to-Text Documentation," cloud.google.com, 2024.
- [4] M. A. Ferrag et al., "Threat Detection of Internet of Things Based on Machine Learning," IEEE Xplore.
- [5] TRAI, "Unsolicited Commercial Communications Regulation," TRAI, 2023.
- [6] S. Dey, A. Sharma, "Integrated Computer Network Security System: Intrusion Detection and Threat Prediction Using Machine Learning Algorithms," IEEE Xplore.
- [7] K. S. Sahoo, B. K. Mishra, "Machine Learning Based Insider Threat Modelling and Detection," IEEE Xplore.
- [8] R. Vinayakumar et al., "API-Based Ransomware Detection Using Machine Learning-Based Threat Detection Models," IEEE Xplore.