

AI in Social Engineering and Phishing Campaigns

~ Janhavi Sonavale
~ Durvaas More

INTRODUCTION



AI is changing the way cyberattacks are carried out. In social engineering and phishing, attackers now use AI to create highly realistic emails, messages, and even fake voices or videos to trick people into giving away sensitive information. These AI-powered scams are faster, harder to detect, and more personalized.

At the same time, cybersecurity experts are using AI to defend against these threats by detecting unusual behavior, blocking phishing emails, and training users. As AI continues to evolve, it is becoming both a powerful tool for attackers and an essential defense for organizations.

WHAT IS SOCIAL ENGINEERING ?

Social engineering refers to the use of manipulation, deception, or trickery to influence individuals into revealing confidential information, such as passwords, account numbers, or access to secure systems. Unlike traditional hacking, which involves exploiting software vulnerabilities, social engineering targets the human element—taking advantage of people's trust, emotions, or curiosity.

WHAT IS PHISHING ?

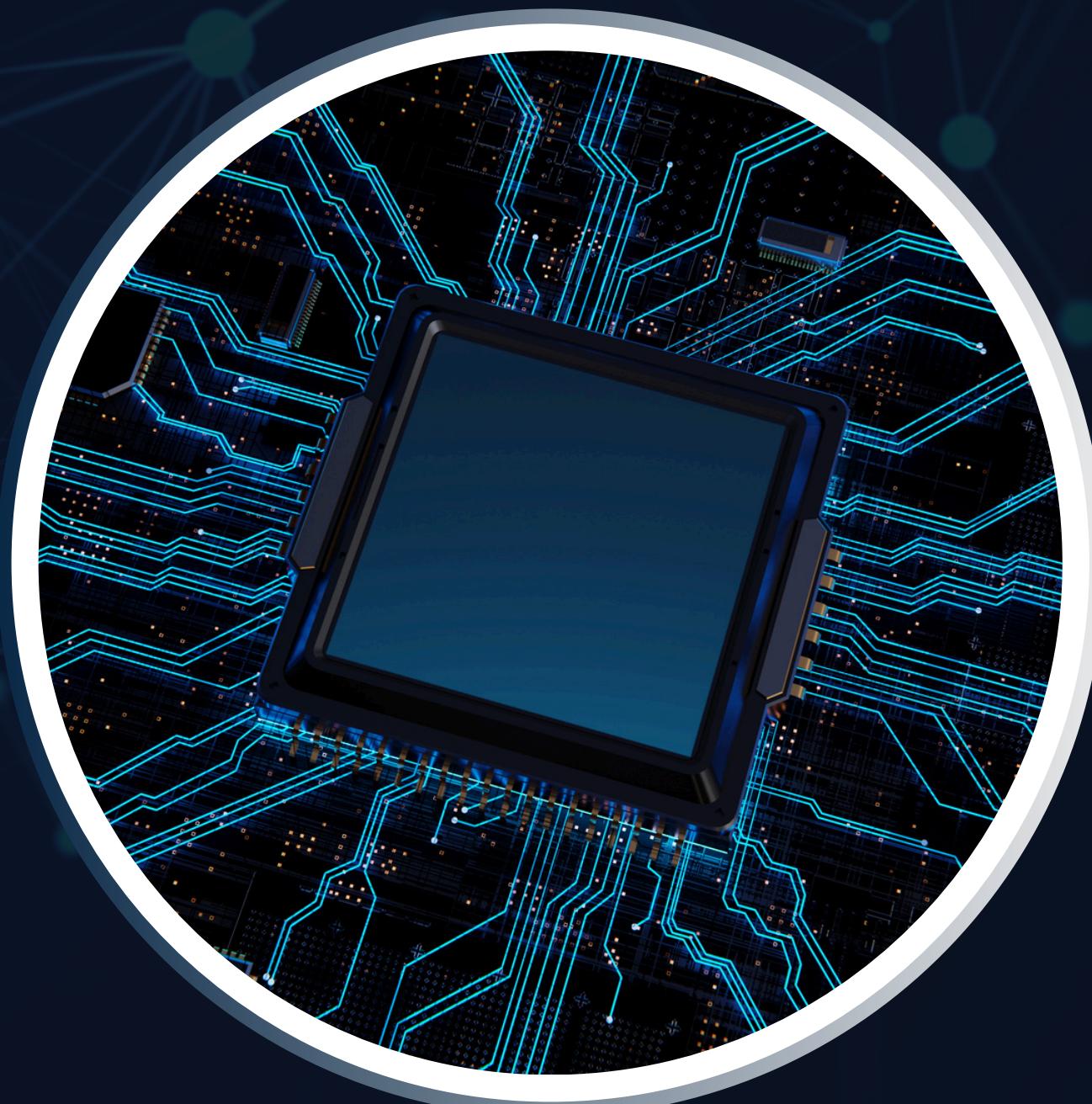
Phishing is a type of cyberattack where an attacker uses fake communications—such as emails, texts, or websites—that appear to come from a trusted source to trick individuals into sharing sensitive information like passwords, credit card numbers, or other personal details.

PROBLEMS

- AI-Enhanced Attacks: Cybercriminals use AI to create personalized phishing emails, deepfake videos, and voice mimicking (vishing) to deceive individuals, making attacks harder to detect.
- Human Vulnerability: Despite advanced technology, people remain the weakest link, with AI exploiting emotions like urgency to trick victims.
- Evolving Attack Methods: Phishing now extends beyond email to SMS (smishing) and phone (vishing) scams, making detection more difficult.

SOLUTIONS

- AI Phishing Detection: AI analyzes emails, websites, and messages to spot suspicious patterns and block phishing attempts.
- Behavioral Analysis: AI monitors user actions to identify anomalies and detect potential breaches early.
- Deepfake Detection: AI identifies deepfake videos and audio, preventing impersonation attacks.
- User Training with AI: AI-driven simulations train users to recognize phishing attempts, strengthening human defenses.



TOOL OVERVIEW

🔍 What It Does:

- Detects whether an email is a phishing attempt or a legitimate message using a trained AI model.
- Generates sample phishing or legitimate emails for testing, training, or demonstrations.

⚙️ Key Features:

- Email content analysis using Natural Language Processing (NLP).
- Random email generation using predefined phishing and legitimate templates.
- CLI (Command-Line Interface) based user interaction.

👤 Use Cases:

- Cybersecurity education & awareness.
- Dataset expansion and AI testing.
- Phishing email simulations.

TOOLS

- 🧠 Language & Libraries:

- Language: Python
- Libraries Used:
 - pandas: Data handling
 - scikit-learn: Machine learning (CountVectorizer + Naive Bayes)
 - joblib: Model saving/loading
 - random & os: General Python utilities

- 📁 File Dependencies:

- phishing_emails.csv – dataset for training
- model.pkl – trained model used for prediction

REAL LIFE CASES

1. Deepfake CEO Voice Scam (UK, 2019)

- What Happened: Criminals used AI-generated audio to mimic the voice of a CEO and called a UK-based energy firm's executive.
- Impact: The executive was convinced he was speaking with his German boss and transferred \$243,000 to a Hungarian bank account.
- Technology Used: AI voice cloning (deepfake audio).
- Significance: Demonstrated how realistic AI-generated voice can bypass normal fraud detection.



2. ChatGPT-Style Phishing Emails

- What Happened: Cybersecurity researchers and threat intelligence firms (like SlashNext and Darktrace) have found that attackers are using AI tools (e.g., ChatGPT) to craft highly convincing phishing emails that are grammatically correct and personalized.
- Impact: Increased click-through rates and trust from targets.
- Technology Used: AI language models to generate human-like phishing content.
- Significance: Removes the telltale signs of phishing (e.g., bad grammar, poor formatting).



FUTURE ENHANCEMENT

1. Multimodal Phishing Detection

- Enhancement: Integrate image, HTML, and attachment analysis alongside text-based email detection.
- Benefit: Detect phishing attempts embedded in logos, spoofed websites, or malicious PDFs/images.

2. Behavioral and Contextual Analysis

- Enhancement: Use behavioral patterns (e.g., time of day, typical communication partners) to spot anomalies.
- Benefit: Adds a contextual layer, reducing false positives and improving detection accuracy.

3. Multilingual Support

- Enhancement: Expand detection capabilities to handle phishing in multiple languages.
- Benefit: Crucial for global organizations facing non-English phishing threats.

4. Integrated Email Client Plugins

- Enhancement: Develop lightweight plugins for Outlook, Gmail, and Thunderbird.
- Benefit: Allows real-time phishing warnings directly within users' email clients.

5. Explainable AI (XAI) Features

- Enhancement: Add transparency by showing why a specific email was flagged.
- Benefit: Increases trust and allows cybersecurity analysts to validate or override decisions.

CONCLUSION

PhishSenseAI is an advanced AI-powered cybersecurity tool designed to detect phishing attacks and generate realistic phishing emails for training and awareness. Leveraging cutting-edge natural language processing (NLP) and machine learning models, it analyzes email content, links, and metadata to identify subtle signs of phishing such as urgency, suspicious tone, and spoofed sender information. It integrates directly with email clients to scan incoming messages in real time, while also monitoring spam folders to detect any overlooked threats. PhishSenseAI also evaluates the reputation and structure of embedded URLs and inspects email headers for irregularities. Beyond detection, the tool enables organizations to generate personalized phishing simulation emails for employee training, helping users learn to recognize and avoid real threats. These simulations can be managed through built-in campaign tools, and each email is assigned a phishing risk score for user clarity. It also supports integration with enterprise-level Security Information and Event Management (SIEM) systems, allowing for real-time alerts and customizable thresholds. With strong privacy safeguards in place, including anonymized data collection and ethical usage policies, PhishSenseAI provides a holistic, responsible, and powerful approach to phishing defense.

**THANK
YOU!**