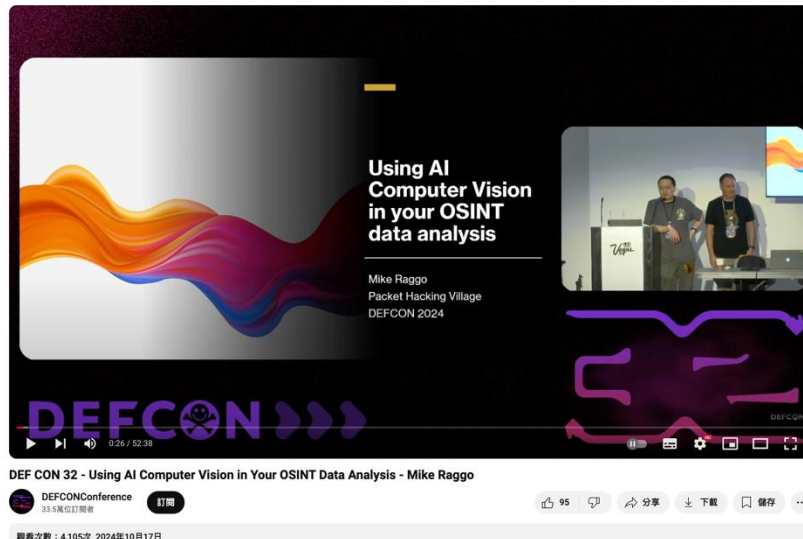# Using AI Computer Vision in Your OSINT Data Analysis

Team 11

110590453　黃台茗

111599004　陳勝舢

112C53010　林幸慧

112c72011　許貽昇

113598023　廖哲霈

113c53019　李宇揚

# Before the beginning

- Our presentation builds on *Def Con 32's Using AI Computer Vision in Your OSINT Data Analysis*.

- We applied the mentioned techniques to real-world cybersecurity incidents, including the <span style="color:red">I-Soon data leak incidents.</span>
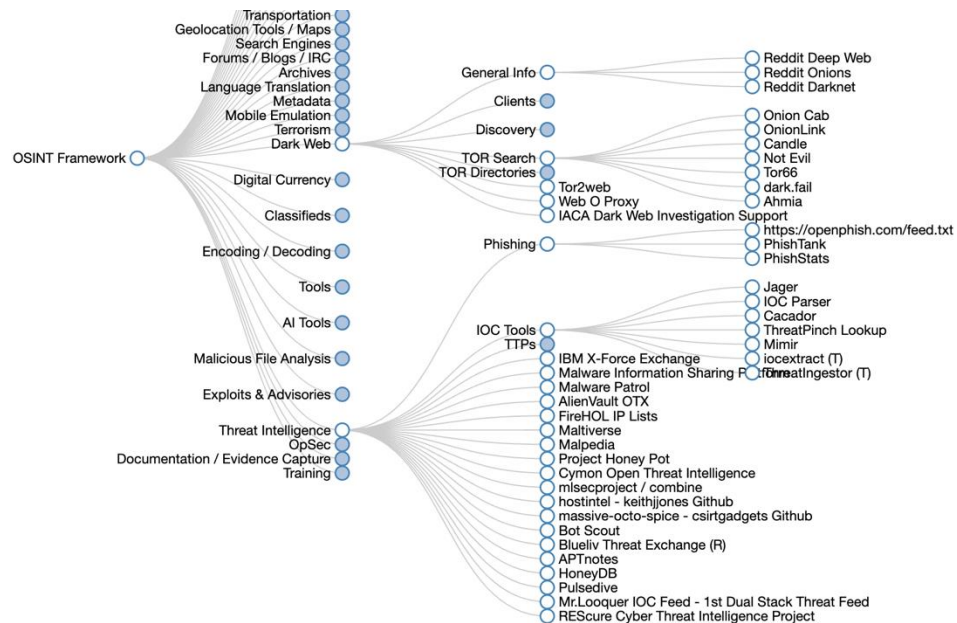
# Outlines

- Introduction
  - OSINT and its challenge
  - AI Computer Vision for OSINT Analysis
- Methodology
  - Optical Character Recognition (OCR)
  - Object Detection
  - Speech  Recognition
  - Dashboard Summary
- Limitation and Vulnerability
- Conclusion
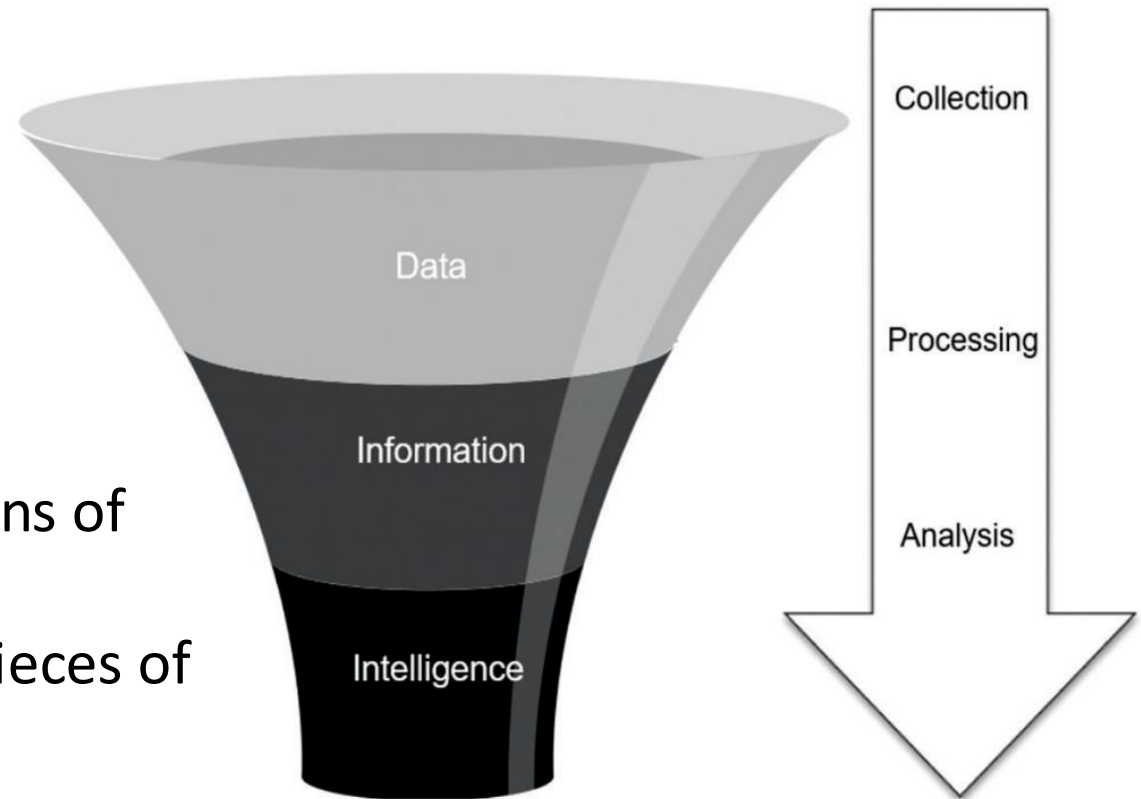- Team Contribution

# Introduction

# What is OSINT?

Open-source intelligence is generated from publicly available information, regularly collected, extracted, and supplied to meet specific intelligence requirements.
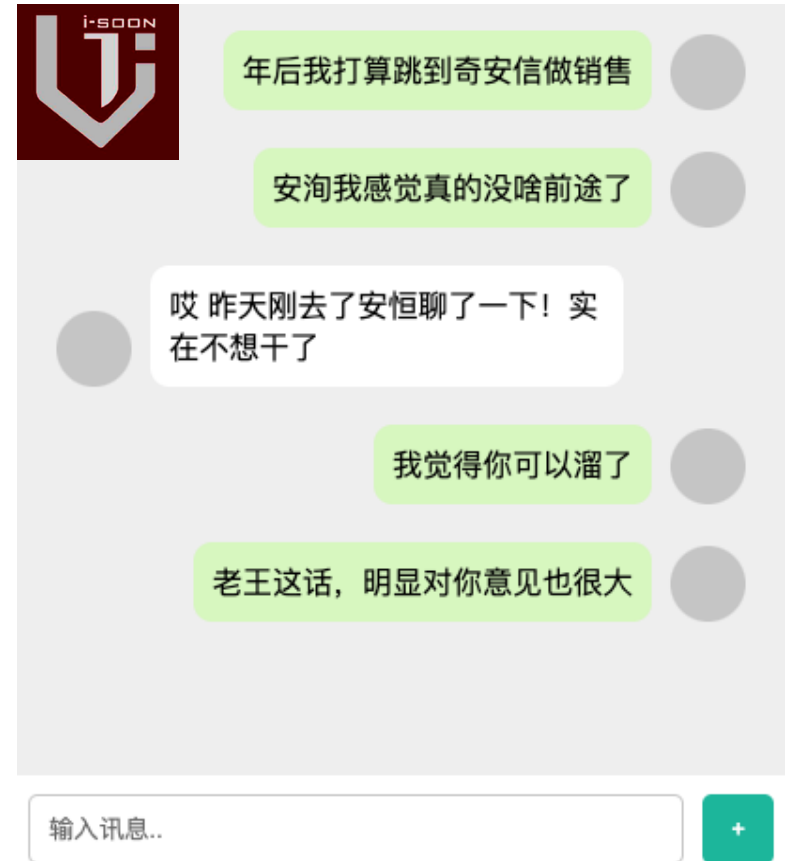
# OSINT Insights

- Unprocessed information
  - e.g., IP addresses, logs, etc.
- Information
  - Structured and filtered data
- Intelligence
  - Mitigation strategies or descriptions of cybersecurity incidents
  - Events correlated from multiple pieces of information

# Analysis Target - I-soon Leak

- Background
  - On February 16, 2024, internal documents from the Chinese company "I-Soon Information Technology" were leaked on GitHub.

- Cause Reason
  - Suspected employee retaliation due to dissatisfaction with compensation and the company.

- Leaked Items
  - 2020-2022 conversation records, victim information, product white papers, salary, and employee lists.



年后我打算跳到奇安信做销售

安洵我感觉真的没啥前途了

哎 昨天刚去了安恒聊了一下！实在不想干了

我觉得你可以溜了

老王这话，明显对你意见也很大

输入讯息..

# Challenges in Analyzing Data Leak (I-Soon)

- The I-Soon Data Leak file contains over 15,743 rows, making the analysis time-consuming...

# OSINT Data Challenge

- Too much data, spend 90% of your sifting through a wasteland of data

- Arduous manual review processes

- Lack of actionable reporting

This begs the need for AI Computer Vision

# AI Computer Vision for OSINT Analysis

- Use Optical Character Recognition (OCR) to identify license plates, signs, and content in documents (passports, licenses, etc.)

- Use Object Detection to sort images by objects identified and narrow your scope in an automated

- Use Video Computer Vision to identify objects, people, signs, and even Audio analysis for conversations.

- Automate flow into reporting, dashboard, map plotting, etc.

# Optical Character Recognition

# Optical Character Recognition (OCR)

- Tesseract OCR is an optical character reading engine developed by HP laboratories in 1985 and open sourced in 2005.

- Tesseract has Unicode (UTF-8) support and can recognize more than 100 languages "out of the box" and thus can be used for building different language scanning software .

# OSINT image OCR using Azure - 1

- Seamless Integration with Azure Cognitive Services

- Automated Intelligence Gathering

# OSINT image OCR using Azure - 2

- Suspected use of tools developed by APT41 (Chengdu 404)
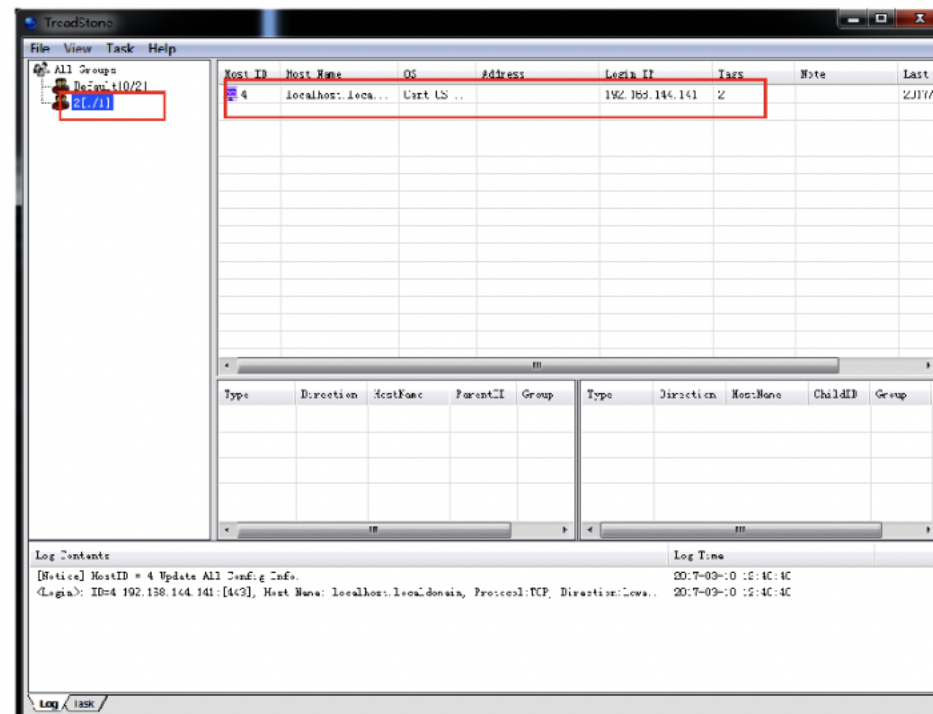


"treadstone." The "treadstone" malware controller software was designed to work with Winnti malware which, at the time, was used only by a small group of hackers – hackers such as QIAN and JIANG, and others they associated with.

（Linux 远程控制系统界面图）

# OSINT image OCR using Azure - 3

- In late September 2022, malicious attackers launched a supply chain attack using the chat program Comm100.

- Windows product name
- Events with event ID 6005 (the event log service was started), events with Event ID 6006 (the event log service was stopped)
- Any TCP endpoints listening on ports 8090, 8091, 8092, 8093, 8094, 8095, 8096
- Any TCP endpoint established to 8.218.67.52:18024
- Antivirus product name
- If Skype or Telegram is installed
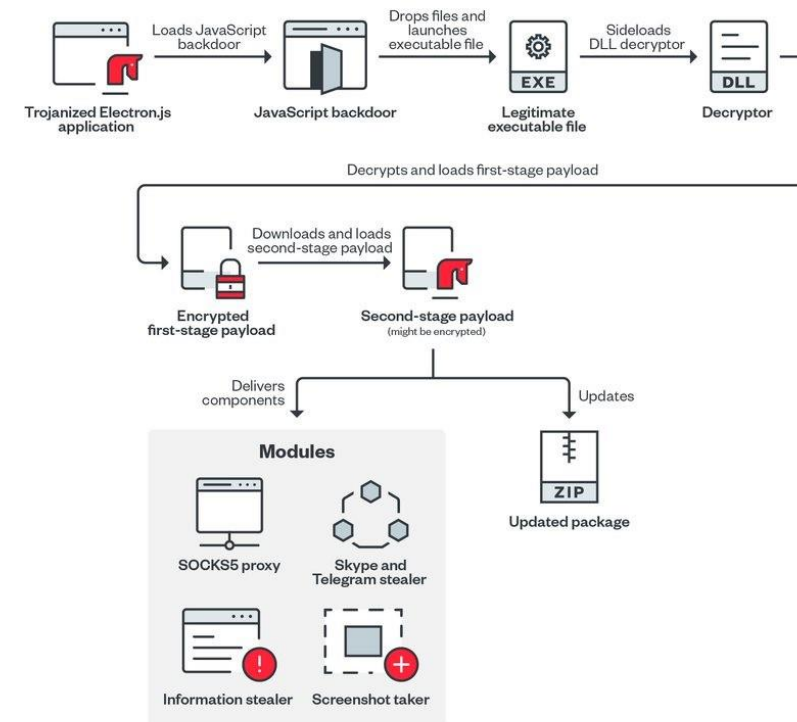- Number of connected monitors
- Product ID value from Windows registry

扬州那边问个人pc的通道

[捂脸]

现在能给不

【彩宝贝】【代理】8.218.67.52:27011 【TCP隧道】8.218.67.52:17011 【账号】admin 【密码】88888888

嗯嗯



©2022 TREND MICRO

# OSINT Take OCR to Word Cloud

The data leak incident lies in its ability to visualize the high-frequency keywords from the leaked conversation, revealing key discussion points and topics.

We can see...
- Words like **"company," "project," "client"** may involve sensitive information related to business activities.
- Words like **"problem," "estimate," and "technology"** might reflect technical details or potential issues discussed internally.
- High-frequency words such as **"we" and "now"** indicate internal collaboration contexts and concerns.

# Object Detection

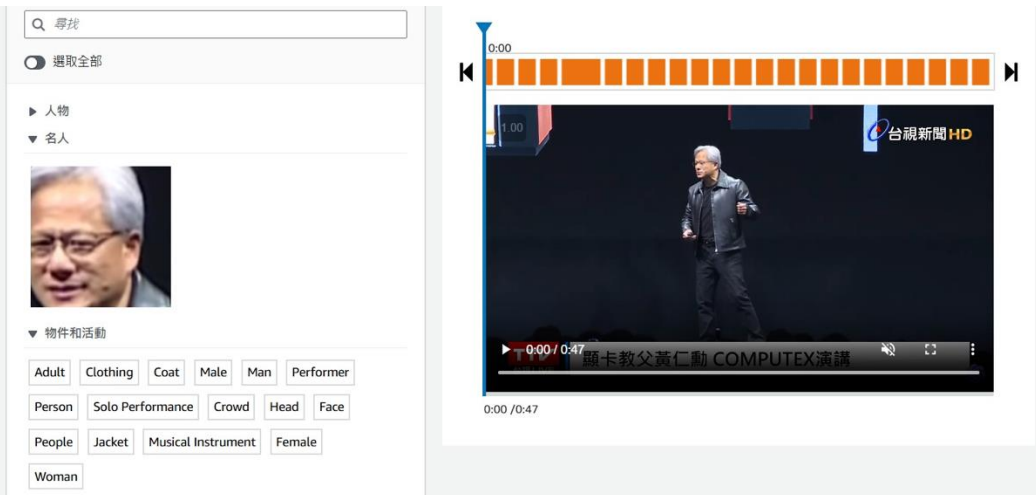# OSINT Object Detection

OSINT Object Detection involves identifying and locating objects of interest within images or videos collected from publicly available sources.



Example using AWS

# OSINT Object Detection with Azure

- Detect and extract bounding boxes based on thousands of recognizable objects and living beings.

- We can confirm identity by comparing similar photos from different documents.

# OSINT Object Detection Train Model

- Azure Object Detection uses optimized Convolutional Neural Networks (CNNs) for feature extraction and object recognition.

# CNN Model

- CNN (Convolutional Neural Network) is a deep learning model designed explicitly for processing structured grid data, such as images.

# OSINT Object Detection Testing Model

- The model may have misclassified the car as a tank due to an insufficient or imbalanced dataset.

- The CNN model likely over-relied on low-level features such as shape or color, failing to distinguish the high-level semantics.

# Why we need OCR and OD?

- Analyze the location of the victim and add geographical coordinates, which is digital identification.



locate

# OSINT video analysis using AWS

Video analysis using a modified version of AWS Rekognition can be used to identify

- The objective throughout the entire video

- Frequency of that objective

- Location in the timeline

- Missed the human eye due to subtleties and backgrounds

# OSINT video analysis using Azure

- Video search and summarization uses a combination of natural language processing and computer vision techniques to analyze the content of a video.

# Speech Recognition

# OSINT Speech Recognition

- Automatic Speech Recognition (ASR), also known as Speech to Text (STT), is the task of transcribing a given audio to text. It has many applications, such as voice user interfaces.

# Speech Recognition Flow

- In the I-Soon Leak Incident, we don't collect related speech files.
- We only introduce the system flow.

# Limitation and Vulnerability

# OCR Problem

There are several common reasons why OCR might fail to recognize text

- **Image quality**
  - Low-resolution scans, poor image clarity, or text distortion can hamper OCR accuracy.
- **Complex fonts and characters**
  - OCR may struggle with non-standard or decorative fonts, unusual characters, or handwriting.
- **Background images**
  - Pages with distracting background images may interfere with text recognition.
- **Improper OCR settings**
  - Incorrect language selection, low confidence thresholds, or other misconfigurations can impact results.

# Speech Recognition Problem

There are still numerous challenges that researchers and developers face when it comes to perfecting speech recognition systems.

- **Variability in Speech Patterns**
  - People speak at different speeds, with varying accents, dialects, and pronunciations.
- **Background Noise and Environmental Factors**
  - Real-world environments are often filled with ambient noise, such as traffic sounds, machinery, or people talking in the background.
- **Vocabulary and Contextual Understanding**
  - The technology must also consider the context in which words are spoken to ensure accurate transcription.
- **Speaker Independence**
  - variations in pitch, tone, and pronunciation between different speakers can impact the algorithm's performance.

# Conclusion

# I-Soon Leak Summary



Sankey Diagram of Communication

- shutd0wn＝吴海波=Wu HaiBo=吴總
- lengmo＝陈诚=Chen Cheng=C總

# I-Soon Leak Summary – Primary Attack Actors

- Who are the primary attack actors?

The primary attack actors identified in the leaked documents are:

1. **I-Soon**:

   - **Role**: I-Soon appears to be a central player in the coordination and execution of cyber attacks, as evidenced by multiple sources mentioning its involvement in various operations.

   - **Collaboration**: This group is noted for its potential cooperation with Chinese government agencies, indicating a likely state-sponsored nature of their activities.

2. **Chengdu 404**:

   - **Role**: This group is another significant threat actor, known for its involvement in various cyber espionage activities and attacks on foreign entities.

   - **Collaboration**: Chengdu 404 is also suspected of having connections with Chinese governmental bodies, particularly those involved in national security.

我们和404没啥合作关键吧?

关联

没有

陪标算不算

都是喝酒理事会成员单位

14号在马来还逮捕了2个同伙

输入讯息..

+

# I-Soon Leak Summary – Business Mode

- Counter-Terrorism Data Support
- Technical Services
- Initiating Talent Cultivation Programs
- Comprehensive Laboratory Development

菲那个，按月付的这种方式

合同怎么签?

直接签一年

看你咋操作

5个w多少个箱子?

你那边儿方便    一个月更新两次

可以选

8-10

5是你的成本价，还是有你自己的利润

成本价

# I-Soon Leak Summary – Product



2.1.5 产品图片

（WiFi 抵近攻击系统（基础版）产品实物图）

（WiFi 抵近攻击系统（mini 版）产品实物图）

# I-Soon Leak Summary – Victim

Taiwan Units Targeted by I-SOON Attacks:
- NTU Hospital: Patient ratio data
- Yuan-Ze Education Foundation: Intranet control
- NTU Institute of Applied Mechanics: Hash passwords, web shell
- Tamkang University: Email privileges

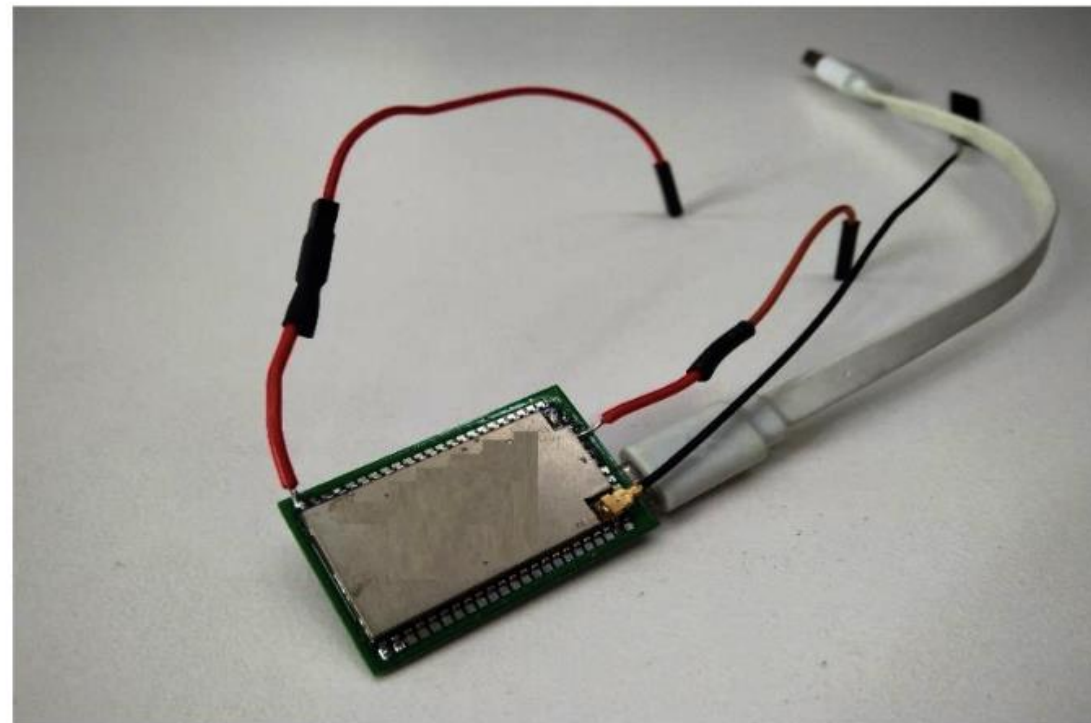| 国家区域 | National area | 目标类型 | Target type | 目标名称 | Target name | Domain Name |
|---|---|---|---|---|---|---|
| 巴基斯坦 | Pakistan | 运营商 | Operator | Zong | Zong | N/A |
| 哈萨克斯坦 | Kazakhstan | 运营商 | Operator | Kcell通讯公司 | Kcell Communication Company | kcell.kz |
| 吉尔吉斯斯坦 | Kyrgyzstan | 运营商 | Operator | megacom | megacom | N/A |
| 马来西亚 | Malaysia | 政府 | government | 工程部 | Engineering department | kkr.gov.my |
| 马来西亚 | Malaysia | 政府 | government | 内政部 | Ministry of the Interior | moha.gov.my |
| 马来西亚 | Malaysia | 政府 | government | 外交部 | Ministry of Foreign Affairs | kln.gov.my |
| 蒙古 | Mongolia | 政府 | government | 警察局 | Police station | N/A |
| 蒙古 | Mongolia | 政府 | government | 外交部 | Ministry of Foreign Affairs | mfa.gov.mn |
| 尼泊尔 | Nepal | 政府 | government | N/A | N/A | N/A |
| 台湾 | Taiwan | 医疗 | Medical care | 台大医院 | National Taiwan University Hospital | ntuh.gov.tw |
| 泰国 | Thailand | 运营商 | Operator | CAT | Cat | N/A |
| 土耳其 | Türkiye | 科技 | science and technology | 科学技术研究理事会 | Council of Science and Technology Research | tubitak.gov.tr |
| 印度 | India | 医疗 | Medical care | 阿波罗医院 | Apollo Hospital | apollohospitals.cc |
| 印度 | India | 政府 | government | 印度出入境 | Indian immigration | UCF |

# Conclusion

- Computer Vision vastly improves accuracy, efficiency, and time-to-evidence.

- Allows the OSINT team to focus on the relevant data for validating narratives and nefarious activities.

- Automates reporting to leverage the human-in-the-loop for focusing on the last mile of analysis and results.

# Team Contribution

# Team Contribution

- ## 110590453 黃台茗
  - developing an OCR tool, video editing, reporting
- ## 111599004 陳勝舢
  - presentations, system implementation, organizing work, building Azure services
- ## 112C53010 林幸慧
  - researching object detection using AWS services
- ## 112c72011 許貽昇
  - developing AI models and APIs
- ## 113598023 廖哲霈
  - web scraping related to I-Soon Leak incidents on Chinese news
- ## 113c53019 李宇揚
  - web scraping related to I-Soon Leak incidents on English websites

# Thanks!