

# Apéndice C: Recolección y caracterización del conjunto de datos CIC-IDS2017

M. A. Pérez Ávila  
*Departamento de Computación*  
*Instituto Tecnológico y de Estudios Superiores de Monterrey*  
Ciudad de México, México  
A01369908@tec.mx

A. M. Sánchez Serratos  
*Departamento de Computación*  
*Instituto Tecnológico y de Estudios Superiores de Monterrey*  
Ciudad de México, México  
A01771843@tec.mx

**Resumen** - Este apéndice presenta la recolección y caracterización del dataset CIC-IDS2017, detallando la infraestructura de red utilizada, la captura de tráfico mediante port mirroring y la generación de flujos bidireccionales con atributos estadísticos y temporales.

## I. INFRAESTRUCTURA DE RED PARA LA CAPTURA DE DATOS

La infraestructura de red diseñada para la recolección de datos del dataset CIC-IDS2017 fue construida con el objetivo de simular un entorno corporativo realista y heterogéneo, combinando tráfico legítimo generado por usuarios humanos con ataques controlados desde múltiples vectores. La red se organizó en dos segmentos principales: la red de atacantes, denominada Outsiders, y la red de víctimas, denominada Insiders. Estos segmentos estaban interconectados a través de un firewall central, que implementaba funciones de filtrado, monitorización y Network Address Translation (NAT), garantizando el registro completo del tráfico entrante y saliente [1].

El segmento de atacantes estaba compuesto por una máquina Kali Linux, con dirección IP 205.174.165.73, utilizada principalmente para ejecutar ataques de fuerza bruta, denegación de servicio (DoS y DDoS) y escaneo de puertos. Complementariamente, se disponía de tres estaciones Windows identificadas con las direcciones 205.174.165.69, 205.174.165.70 y 205.174.165.71, que generaban tráfico malicioso adicional, ejecutaban scripts automatizados y simulaban usuarios internos comprometidos. Cada una de estas máquinas contaba con software especializado para pruebas de intrusión, incluyendo herramientas de ataque web, clientes FTP y generadores de tráfico de red, permitiendo reproducir escenarios de ataque realistas y variados. Todo este conjunto se encontraba interconectado mediante un router y un switch que estructuraban la red ofensiva [2].

Por otra parte, la red de víctimas estaba conformada por tres servidores, un firewall, dos switches y diez PCs interconectadas mediante un controlador de dominio (DC) y Active Directory. Cada equipo cumplía funciones específicas dentro del entorno corporativo, con sistemas operativos Linux, Windows y macOS. Además, uno de los puertos del switch principal fue configurado como puerto espejo, permitiendo la captura completa de todo el tráfico enviado y recibido en la red, garantizando la integridad y

completitud de los datos para su posterior análisis [2]. A continuación, se enlistan los equipos que conforman la red:

1. Web Server 16 Public: Ubuntu Server, IP interna 192.168.10.50, IP pública 205.174.165.68, alojando servicios web accesibles externamente.
2. Ubuntu Server 12 Public: Ubuntu Server, IP interna 192.168.10.51, IP pública 205.174.165.66, ofreciendo servicios FTP y SSH.
3. Ubuntu 14.4 32-bit: Servidor interno, IP 192.168.10.19, configurado para pruebas de tráfico normal y logs de sistema.
4. Ubuntu 14.4 64-bit: Servidor interno, IP 192.168.10.17, usado para servicios internos y almacenamiento de datos de prueba.
5. Ubuntu 16.4 32-bit: Servidor interno, IP 192.168.10.16, ejecutando servicios de SSH y aplicaciones internas.
6. Ubuntu 16.4 64-bit: Servidor interno, IP 192.168.10.12, configurado como servidor de aplicaciones y pruebas de servicios de red.
7. Windows 7 Pro 64-bit: Estación de trabajo interna, IP 192.168.10.9, simulando actividad de usuarios finales, navegación web y correo electrónico.
8. Windows 8.1 64-bit: Estación de trabajo interna, IP 192.168.10.5, generando tráfico de usuario y realizando transferencias de archivos.
9. Windows Vista 64-bit: Estación de trabajo interna, IP 192.168.10.8, utilizada para pruebas de compatibilidad y tráfico legado.

10. Windows 10 Pro 32-bit: Estación de trabajo interna, IP 192.168.10.14, simulando usuarios finales actuales en actividades diarias.
11. Windows 10 64-bit: Estación de trabajo interna, IP 192.168.10.15, adicional para balance de tráfico y pruebas de red.
12. macOS: Estación interna, IP 192.168.10.25, utilizada para almacenamiento de archivos y ejecución de aplicaciones corporativas.

El servidor DNS interno, con dirección IP 192.168.10.3, proporcionaba resolución de nombres dentro de la red, asegurando la conectividad entre dispositivos y servicios internos.

El firewall principal conectaba ambos segmentos de red, operando con su interfaz WAN configurada en 205.174.165.80 y la interfaz LAN en 172.16.0.1. Este dispositivo filtraba el tráfico, aplicaba reglas de seguridad y registraba cada paquete, permitiendo un historial completo de las comunicaciones y garantizando la captura integral de datos para su análisis posterior. La topología diseñada reproducía escenarios corporativos realistas, combinando tráfico interno legítimo, interacción con redes externas y ataques dirigidos desde diversos vectores, lo que aseguraba que el dataset resultante reflejara tanto comportamientos normales como situaciones de riesgo auténticas.

A continuación, se muestra una tabla que resume la infraestructura de red utilizada en el dataset CIC-IDS2017, incluyendo todos los servidores, estaciones de trabajo, firewall y PCs de los atacantes, junto con sus sistemas operativos y direcciones IP privadas y públicas [2].

TABLA I  
INFRAESTRUCTURA DE LA RED EN CIC-IDS2017

Categoría	Máquina / SO	IPs
<b>Servidores de la red víctima</b>	Windows Server 2016 (DC y DNS)	192.168.10.3
	Ubuntu (Servidor Web)	16 192.168.10.50 - 205.174.165.68
	Ubuntu 12	192.168.10.51 - 205.174.165.66
<b>PCs de la red víctima</b>	Ubuntu 14.4 (32 y 64 bits)	192.168.10.17 - 192.168.10.19
	Ubuntu 16.4 (32 y 64 bits)	192.168.10.12 - 192.168.10.16

Categoría	Máquina / SO	IPs
	Windows 7 Pro	192.168.10.9
	Windows 8.1 64 bits	192.168.10.5
	Windows Vista	192.168.10.8
	Windows 10 Pro (32 y 64 bits)	192.168.10.14 - 192.168.10.15
	Mac	192.168.10.25
<b>Firewall</b>	Fortinet	—
<b>Atacantes</b>	Kali	205.174.165.73
	Windows 8.1	205.174.165.69 - 205.174.165.71

Asimismo, la Figura 1 presenta la arquitectura del entorno de prueba utilizado para la recolección de datos en el dataset CIC-IDS2017. Esta imagen ilustra la topología completa de la red, incluyendo los segmentos de atacantes y víctimas, así como la interconexión de routers, switches, firewall, servidores, estaciones de trabajo y el puerto espejo configurado para la captura de tráfico mediante port mirroring [2].

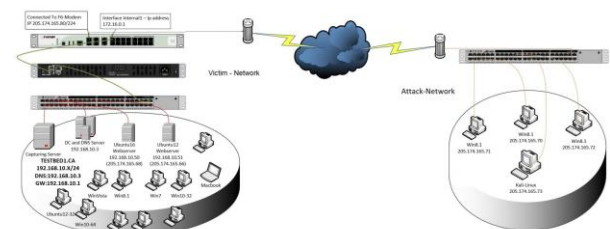


Fig.1 Arquitectura del entorno de prueba (Fuente: Sharafaldin et al., 2018).

## II. CAPTURA DE TRÁFICO Y PORT MIRRORING

La captura del tráfico de red en la infraestructura diseñada para CIC-IDS2017 se realizó utilizando la técnica de port mirroring, también conocida como Switched Port Analyzer (SPAN). Esta metodología permite replicar todo el tráfico que pasa por un puerto de un switch hacia otro puerto conectado a un dispositivo de captura, de manera que se puedan registrar todos los paquetes transmitidos sin afectar el funcionamiento de la red. Port mirroring garantiza la integridad de los datos capturados, ya que permite obtener tanto los headers como los payloads de cada paquete, asegurando que la información contenida en la comunicación sea completa y apta para análisis detallados.

### III. DINÁMICA DE RECOLECCIÓN Y CRONOGRAMA

La recolección de datos del dataset CIC-IDS2017 se llevó a cabo durante cinco días consecutivos, del lunes 3 al viernes 7 de julio de 2017, combinando tráfico legítimo generado por usuarios humanos con ataques planificados según un cronograma previamente definido. Este enfoque permitió registrar tanto la actividad normal de la red como los distintos tipos de ataques contemporáneos, generando un conjunto de flujos bidireccionales representativos de un entorno corporativo realista.

El primer día, lunes 3 de julio, se capturó únicamente tráfico benigno. Durante este periodo, los usuarios realizaron actividades comunes, incluyendo navegación web, transferencias FTP, sesiones SSH y uso de correo electrónico. Esta sesión inicial permitió establecer una línea base del tráfico normal, obteniendo un total de 529,918 flujos y un volumen de datos de 11 GB. A partir del martes 4 de julio se introdujeron ataques específicos según un cronograma estructurado. El martes se ejecutaron ataques de fuerza bruta sobre servicios FTP y SSH mediante herramientas automatizadas, mientras que el miércoles 5 de julio se realizaron ataques de Denegación de Servicio, incluyendo DoS Hulk, DoS GoldenEye, Slowhttptest y la explotación de la vulnerabilidad Heartbleed. Durante el jueves 6 de julio se llevaron a cabo ataques web y de infiltración, divididos en sesiones de mañana y tarde, que incluyeron fuerza bruta web, SQL Injection, XSS y descarga de exploits. Finalmente, el viernes 7 de julio se registraron actividades de botnet y escaneo de puertos, combinadas con ataques DDoS.

La captura total a lo largo de estos cinco días produjo aproximadamente 2,830,743 flujos y un volumen de datos cercano a los 50 GB en formato PCAP. La siguiente tabla resume la actividad por día, incluyendo el tipo de tráfico registrado, la descripción de las actividades y el tamaño y cantidad de flujos obtenidos [3]:

TABLA II  
DISTRIBUCIÓN DE FLUJOS POR ARCHIVO Y TIPO DE ATAQUE EN CIC-IDS2017

Nombre del archivo	Ataque	Recuento ataques	Proporción	Total Incidentes
Monday-WorkingHours.pcap_ISCX.csv	Benign	529,918	100%	529,918
Tuesday-WorkingHours.pcap_ISCX.csv	Benign	432,074	96.90%	445,909
	SSH-Patator	5,897	1.32%	
	FTP-Patator	7,938	1.78%	
	Benign	440,031	64.06%	

Wednesday-WorkingHours.pcap_ISCX.csv	DoS Hulk	231,073	33.64%	170,366
	DoS GoldenEye	10,293	1.50%	
	DoS Slowhttptest	5,499	0.80%	
	Heartbleed	11	0.00%	
Thursday-WorkingHours-Morning-WebAttack.pcap_ISCX.csv	Benign	168,186	98.72%	170,366
	Web Attack - Brute Force	1,507	0.88%	
	Web Attack - Sql Injection	21	0.01%	
	Web Attack - XSS	652	0.38%	
Thursday-WorkingHours-Afternoon-Infiltration.pcap_ISCX.csv	Benign	288,566	99.99%	288,602
	Infiltration	36	0.01%	
Friday-WorkingHours-Morning.pcap_ISCX.csv	Benign	189,067	98.97%	191,033
	Bot	1,966	1.03%	
Friday-WorkingHours-Afternoon-PortScan.pcap_ISCX.csv	Benign	127,537	44.52%	286,467
	PortScan	158,930	55.48%	
Friday-WorkingHours-Afternoon-DDos.pcap_ISCX.csv	Benign	97,718	43.29%	225,745
	DDoS	128,027	56.71%	

#### IV. ESTRUCTURA DE LAS MUESTRAS

Cada registro del dataset CIC-IDS2017 corresponde a un flujo bidireccional de red, también denominado FlowID, que representa la comunicación completa entre dos nodos de la red. Un flujo se define mediante la combinación de cinco elementos fundamentales: la dirección IP de origen, el puerto de origen, la dirección IP de destino, el puerto de destino y el protocolo de comunicación. Esta definición permite identificar de manera única cada “conversación” dentro de la red, distinguiendo el tráfico entre distintos pares de hosts y servicios.

Cada flujo captura información en ambas direcciones: forward, del origen al destino, y backward, del destino al origen. Esta agregación bidireccional permite registrar patrones completos de la comunicación, incluyendo la secuencia de paquetes enviados, los tiempos de envío y recepción, los tamaños de los paquetes, y los flags TCP (como SYN, ACK, FIN, RST, PSH y URG) que reflejan el estado y control de la conexión. Al consolidar estos datos en un solo flujo, se puede analizar la actividad completa de la conexión, identificando comportamientos anómalos que no serían evidentes a nivel de paquete individual.

#### V. EXTRACCIÓN DE CARACTERÍSTICAS

La extracción de características del dataset CIC-IDS2017 se realizó utilizando el software CICFlowMeter, una herramienta especializada en análisis de tráfico de red que convierte archivos de captura de paquetes (PCAP) en flujos bidireccionales y genera atributos estadísticos y temporales de cada flujo. CICFlowMeter funciona siguiendo un enfoque de procesamiento por flujo, donde primero agrupa los paquetes según el FlowID (combinación de IP de origen, puerto de origen, IP de destino, puerto de destino y protocolo), de modo que cada flujo representa una “conversación” completa entre dos nodos de la red.

El software analiza tanto la dirección forward (del origen al destino) como backward (del destino al origen), permitiendo capturar la comunicación de manera integral. Durante este procesamiento, CICFlowMeter calcula atributos de tres tipos principales: temporal, estadístico y de control de conexión, los cuales reflejan distintos aspectos del tráfico de red:

1. Características temporales: Incluyen la duración total del flujo y los tiempos inter-arribo de paquetes, que indican la diferencia temporal entre paquetes consecutivos. Se registran valores mínimos, máximos, promedio y desviación estándar, permitiendo identificar patrones de envío irregular o ráfagas de actividad que pueden ser indicativos de ataques como DoS o escaneos de puertos.
2. Características estadísticas de tamaño y volumen: CICFlowMeter calcula el tamaño de los paquetes en ambas direcciones, incluyendo el mínimo, máximo, promedio, desviación estándar y varianza. También determina la velocidad de transmisión en bytes por segundo y paquetes por segundo. Estos atributos permiten evaluar la carga y el comportamiento del flujo, diferenciando tráfico normal (como navegación web o correo electrónico) de tráfico malicioso con patrones inusuales o volumétricos.
3. Flags TCP y control de conexión: Se registran los estados de los paquetes TCP, incluyendo SYN, ACK, FIN, RST, PSH y URG, que reflejan la forma en que se establece, mantiene y finaliza la comunicación. Estos flags son fundamentales para identificar ataques de fuerza bruta, intentos de escaneo de puertos y conexiones anómalas.
4. Actividad e inactividad: CICFlowMeter mide periodos de actividad y pausas dentro de cada flujo, proporcionando información sobre ráfagas de tráfico o retrasos inusuales que podrían indicar comportamiento anómalo.
5. Tamaños de ventana TCP iniciales: Se registran para ambas direcciones, forward y backward, lo que permite evaluar la capacidad de transmisión y la configuración de los sistemas involucrados, siendo útil para detectar ciertas formas de manipulación de la conexión.

En total, cada flujo puede generar hasta 84 atributos, que combinan información temporal, estadística y de control de la conexión. Esta riqueza de atributos convierte a CICFlowMeter en una herramienta poderosa, ya que permite transformar el tráfico de red bruto en un conjunto de datos estructurado y listo para análisis, proporcionando la base necesaria para entrenar modelos de aprendizaje automático y evaluar sistemas de detección de intrusiones.

#### REFERENCES

- [1] A. Rosay, E. Cheval, F. Carlier y P. Leroux, "Network Intrusion Detection: A Comprehensive Analysis of CIC-IDS2017", *Proceedings of the 8th International Conference on Information Systems Security and Privacy (ICISSP 2022)*, Le Mans, Francia, 2022, pp. 25-36. [Online]. Available: <https://www.scitepress.org/Papers/2022/107740/107740.pdf>.
- [2] I. Sharafaldin, A. H. Lashkari y A. A. Ghorbani, "Toward Generating a New Intrusion Detection Dataset and Intrusion Traffic Characterization", *Proceedings of the 4th International Conference on Information Systems Security and Privacy (ICISSP 2018)*, Lisboa, Portugal, 2018, pp. 108-116. [Online]. Available: <https://www.scitepress.org/papers/2018/66398/66398.pdf>.
- [3] Z. I. Khan, M. M. Afzal y K. N. Shamsi, "A Comprehensive Study on CIC-IDS2017 Dataset for Intrusion Detection Systems", *International Research Journal on Advanced Engineering Hub*, vol. 2, no. 2, pp. 254-260, Feb. 2024. [Online]. Available: [https://www.researchgate.net/publication/378709289\\_A\\_Comprehensive\\_Study\\_on\\_CIC-IDS2017\\_Dataset\\_for\\_Intrusion\\_Detection\\_Systems](https://www.researchgate.net/publication/378709289_A_Comprehensive_Study_on_CIC-IDS2017_Dataset_for_Intrusion_Detection_Systems).