



TRABAJO FIN DE GRADO
INGENIERÍA EN INFORMÁTICA

Implementación en tiempo real de sistemas de identificación de tráfico de red

Subtítulo del proyecto

Autor

Álvaro Maximino Linares Herrera

Directores

Jesús Esteban Díaz Verdejo



ESCUELA TÉCNICA SUPERIOR DE INGENIERÍAS INFORMÁTICA Y DE
TELECOMUNICACIÓN

—
Granada, 3 de septiembre de 2017



Implementación en tiempo real de sistemas de identificación de tráfico de red

Subtitulo del proyecto

Autor

Álvaro Maximino Linares Herrera

Directores

Jesús Esteban Díaz Verdejo

Implementación en tiempo real de sistemas de identificación de tráfico de red

Álvaro Maximino Linares Herrera

Palabras clave: inspección profunda de paquetes, bro, flujos, identificación, clasificación, red, tráfico

Resumen

La identificación de tráfico en red es realmente importante para aplicaciones de ingeniería de tráfico y de seguridad.

En este trabajo se tratará la creación de un programa para un NMS (Network Monitoring System), en este caso se usará BRO, mediante el cual se pueda resolver el emparejamiento de flujos. BRO consiste en un NMS que funciona mediante el terminal en Linux o Mac, una de las peculiaridades de este programa es que para la creación de scripts que nos permitan extender la funcionalidad de la que dispone, tendremos que usar BRO como lenguaje de programación. Es un lenguaje de scripting, el cual está orientado a eventos, que se lanzan cuando ocurre algo relacionado con el control y análisis de redes, es un lenguaje que para los que vienen de C++, Java o Python, no debe de suponer un gran reto, más allá de acostumbrarse a sus sintaxis. Es un lenguaje potente que al estar orientado a redes nos permite obtener mucha información de los flujos que tenemos en la red o en el archivo que vayamos a analizar. En este trabajo mediante implementaciones offline se verificará la eficacia de esta técnica de clasificación de tráfico.

Project Title: Project Subtitle

First name, Family name (student)

Keywords: Keyword1, Keyword2, Keyword3,

Abstract

Write here the abstract in English.

Yo, **Álvaro Maximino Linares Herrera**, alumno de la titulación TITULACIÓN de la **Escuela Técnica Superior de Ingenierías Informática y de Telecomunicación de la Universidad de Granada**, con DNI 76669401M, autorizo la ubicación de la siguiente copia de mi Trabajo Fin de Grado en la biblioteca del centro para que pueda ser consultada por las personas que lo deseen.

Fdo: Álvaro Maximino Linares Herrera

Granada a 3 de septiembre de 2017.

D. **Jesús Esteban Díaz Verdejo**, Profesor del Área de XXXX del Departamento YYYY de la Universidad de Granada.

Informa:

Que el presente trabajo, titulado *Implementación en tiempo real de sistemas de identificación de tráfico de red*, ha sido realizado bajo su supervisión por **Álvaro Maximino Linares Herrera**, y autorizamos la defensa de dicho trabajo ante el tribunal que corresponda.

Y para que conste, expiden y firman el presente informe en Granada a 3 de septiembre de 2017.

Los directores:

Jesús Esteban Díaz Verdejo	Nombre Apellido1 Apellido2 (tu- tor2)
-----------------------------------	--

Agradecimientos

Poner aquí agradecimientos...

Índice general

1. Introducción	15
1.1. Motivación y antecedentes	15
1.2. Objetivos	16
1.3. Metodología	17
1.4. Estructura de la memoria	18
2. Estado del arte	19
2.1. Identificación de tráfico	19
2.1.1. Técnicas de identificación de tráfico	19
2.1.2. Identificación de tráfico basada en flujos	20
2.1.3. DPI	21
2.2. BRO	21
2.2.1. Funcionalidades básicas de BRO	23
2.2.2. Eventos y trazas	24
2.2.3. Incorporación de funcionalidades	24
2.3. Emparejamiento de flujos	25
3. Diseño y arquitectura del sistema	27
3.1. Arquitectura del sistema	27
3.2. Módulo y funciones	28
3.3. Gestión de flujos	29
3.4. Estructuras de datos	30
4. Implementación	33
5. Evaluación y pruebas	37
6. Conclusiones y trabajo futuro	39
Bibliografía	43

Capítulo 1

Introducción

1.1. Motivación y antecedentes

A partir de los años 90 Internet tuvo una gran expansión, pasando de solo compartir información entre investigadores, como se hacía al principio, a permitir realizar compras, hacer videollamadas, pedir cita en el médico, gestionar todo lo relacionado con el banco y demás. Internet llega a miles de millones de personas que consumen recursos de forma masiva. Por lo que es crucial tener cierta calidad de servicio [34] en los recursos que se consumen.

La calidad de servicio es lo más importante tanto a nivel de desarrollo como de consumidor. Como desarrollador, se tendrá que pensar en esta característica como una de las más importantes, pues un mal servicio hará que se pierdan clientes. Como cliente la calidad de servicio es algo deseable, pues si por ejemplo en un servicio de videollamada se experimenta cierta pérdida de calidad, no pasará nada. Por el contrario si en un servicio de visionado de contenido bajo demanda de algún evento deportivo se experimenta la pérdida de calidad en la imagen o se cae el servicio lo más probable es que el usuario que lo experimenta no vuelva a usar esa plataforma. [1]

También hay que tener en cuenta la seguridad, pues como ya se ha mencionado se tratan datos de tipo muy sensible, como los datos bancarios y sanitarios.

Existe una forma de dar prioridad a diferentes tipos de tráfico, mediante la identificación de este y su posterior clasificación. La identificación de tráfico consiste en analizar el tráfico de una red y mediante alguna de las distintas técnicas existentes organizarlo mediante patrones comunes, como podría ser las IP's y los puertos, el contenido de las cabeceras de los paquetes, mediante técnicas heurísticas, etc. Mediante estas técnicas se podrá dar prioridad a un tipo de tráfico sobre otro, mejorando así la calidad de servicio.

Una de las técnicas más usadas consiste en analizar cada uno de los paquetes de la red. Esta técnica es conocida como Inspección Profunda de Paquetes, DPI o *Deep Packet Inspection* [2] por sus siglas en inglés, trata de encontrar dentro de los paquetes cadenas o patrones que se repitan para poder decir que pertenecen al mismo tipo.

Esta técnica, aunque efectiva tiene varios inconvenientes.

- *Escalabilidad.* Ya que se trata de una técnica que tiene que ir paquete a paquete y entrar dentro de ellos requiere de una gran cantidad de recursos. Por lo que si se trata de analizar una red pequeña se obtendrán buenos resultados en cuanto a rendimiento, pero de tratarse de una red amplia, se tendrán malos tiempos de análisis.
- *Privacidad.* Al entrar dentro de los paquetes se produce cierta violación de las políticas de privacidad de la red. Si se trata de una red local no hay problemas. Pero si se trata de una red pública si que se incurre en violaciones de privacidad.

Una solución parcial a este problema se solucionaría con la técnica de emparejamiento de flujos [3], desarrollada por el departamento de TSTC de la Universidad de Granada y probada en entornos de laboratorio [4]. Esta técnica en comparación con otras con las que se realizaron comparaciones es mucho más eficiente y respeta la privacidad, pues no necesita entrar dentro de los paquetes, solamente necesita la información de las IP's de origen y destino y la de los puertos de origen y destino.

Por lo tanto en el presente trabajo se va a implementar esta técnica y se probará más allá de un entorno de laboratorio, llevándola a escenarios reales.

Para llevar esta técnica a un escenario real se precisará de un monitor de redes, NMS o *Network Monitoring System*. A pesar que existen muchos y todos son bastante competentes, existe uno llamado Bro [5], el cual permite incorporar funcionalidades extras a las que ya trae por defecto. Esto se puede realizar mediante la creación de módulos que se lanzarán cuando se inicie el programa.

1.2. Objetivos

El objetivo de este trabajo es el desarrollo de un módulo para un NMS, en este caso Bro. Con el desarrollo del módulo se tratará de demostrar que la técnica de emparejamiento de flujos se puede realizar fuera de un entorno de laboratorio.

Este objetivo se descompone en distintos subobjetivos.

- Conocimiento de la gestión de los flujos en Bro.
- También hará falta gestionar las entradas y salidas. En un principio se hará uso de un archivo *pcap*, de forma que se pueda comprobar que todo funciona correctamente. Una vez terminado será posible realizar una ejecución a tiempo real en una red.
- Implementación de la función de emparejamiento de flujos, así como el control de los distintos eventos para la gestión del tráfico.
- Se realizarán pruebas del funcionamiento correcto del módulo.

Por último, y fuera de los objetivos más técnicos, el desarrollo de esta memoria también se considera un objetivo, por lo tanto se tendrá en cuenta dentro de la temporización.

1.3. Metodología

Para realizar este trabajo se establecen una serie de tareas:

- Estado del Arte.
 - Lectura del artículo del departamento. [4]
 - Búsqueda de información sobre la identificación de tráfico.
 - Análisis de las herramientas.
- Diseñar el módulo.
- Implementar el módulo.
- Evaluación y pruebas del módulo.

En el siguiente diagrama de Gantt se puede ver una temporización de estas tareas.

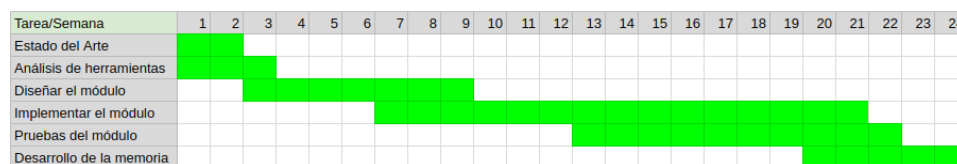


Figura 1.1: Temporización en diagrama de Gantt.

El gasto que requiere este proyecto se desglosa de la siguiente manera.

- *Licencias.* El gasto en licencias será nulo, pues Bro está creado bajo licencia de software libre [5]. El módulo será subido a GitHub [6] con licencia de software libre, por lo que cualquiera podrá usarlo o modificarlo en el futuro.
- *Equipo.* Teniendo en cuenta que la vida útil de un portátil es de unos 4 años y que el desarrollo de este trabajo requiere de unos 6 meses, se tendrá que un portátil de gama media-alta de 800€ generará un gasto de un octavo de la vida útil. Por lo tanto será un coste de unos 100€. Hay que tener en cuenta que Bro solo funciona en Linux o Mac OS X [7], por lo que se está hablando de un portátil y no de un Mac, lo cual supondría un mayor gasto.
- *Programador.* Según se puede comprobar en Internet [8], el precio por hora de un programador está entre los 30€ y los 40€, por lo tanto si en el desarrollo del proyecto se ha necesitado de unas 500 horas y dejando el precio en 30€ la hora, el gasto será de 15000€. Aquí el gasto podría ser mayor, pues no solo se necesita un programador, sino además un director de proyecto, lo cual subiría más este gasto.
- *Varios.* Además de los gastos ya descritos habrá que sumarle una parte de gastos varios como Internet, luz, agua, y demás. Unos 150€ al mes, que al ser 6 meses supone un gasto de 900€.

Teniendo en cuenta todos estos cálculos, se tiene que el coste total del proyecto es de unos 16000€ por 6 meses de trabajo.

1.4. Estructura de la memoria

Esta memoria se organizará de la siguiente forma:

- En el capítulo 2 se hablará de todos los fundamentos teóricos y tecnológicos sobre los que se basa el proyecto.
- En el capítulo 3 se contará cómo se pretende resolver el problema expuesto.
- En el capítulo 4 se encontrará detallado cómo se han implementado los diferentes módulos.
- En el capítulo 5 se realizarán las pruebas, para comprobar que todo funciona como estaba previsto, tanto a nivel funcional como a nivel de aplicación.
- En el último capítulo se hablará de las conclusiones recogidas a lo largo de este proyecto y las posibles opciones que tiene para seguir trabajando sobre él.

Capítulo 2

Estado del arte

Aquí se encuentra la parte teórica y tecnológica del proyecto. Se hablará de la identificación de tráfico y de las distintas técnicas que existen para ello. También se hablará de Bro [5], el analizador de tráfico que se usará para llevar a cabo el desarrollo del módulo. Se introducirá en que consiste, cómo gestiona los eventos y sus funcionalidades básicas, así como la posibilidad de ampliar estas.

Por último se presentará de forma teórica en que consiste el emparejamiento de flujos y de qué forma se podría usar para identificar el tráfico.

2.1. Identificación de tráfico

Hay que tener claro que identificar el tráfico no es clasificarlo. Identificar el tráfico consiste en analizarlo y obtener patrones comunes con los cuales se pueda ordenar el tráfico de forma posterior. Algunos de los patrones que se pueden tener en cuenta a la hora de identificar tráfico son los puertos, las IP's o la clase a la que pertenece.

Por lo tanto la identificación del tráfico de la red es esencial para realizar una clasificación posterior. Dicha clasificación se podrá analizar en busca de amenazas o intrusos, a nivel de seguridad y componentes lentos o defectuosos que afecten al rendimiento del servicio. Para poder identificar el tráfico lo deseable es un método que sea rápido y proteja la privacidad de los distintos usuarios.

2.1.1. Técnicas de identificación de tráfico

Existen varias técnicas de identificación de tráfico, como ya se explicó anteriormente. Estas técnicas son el paso previo a la clasificación del tráfico. Podría considerarse que son el paso intermedio a la clasificación, teniendo

en cuenta que el tráfico es capturado, identificado y clasificado. El tráfico, hasta ahora, se puede identificar de la siguiente forma.

- Por los puertos de la capa de transporte, establecido por *IANA*. [9]
- Por el contenido del paquete o *DPI*. [10]
- Por la aplicación de técnicas de aprendizaje automático sobre estadísticas de tráfico, *machine learning*. [11]

En los siguientes apartados se verán estas técnicas de una forma más detallada.

2.1.2. Identificación de tráfico basada en flujos

Dentro de este ámbito se encuentran dos de las tres técnicas que se han presentado anteriormente.

- Identificación basada en los puertos de la capa de transporte. [9]

Esta técnica consiste en identificar el tráfico dependiendo de los puertos de la capa de transporte, fue diseñada por *IANA* [12]. Esto se debe a que *IANA* es la organización que se encarga de asignar los protocolos de Internet, algunos nombres de dominios y coordinar la asignación de direcciones IP.

Por lo tanto, *IANA* es quien asigna el número de puerto a los distintos protocolos, haciendo que se pueda identificar el tráfico por el número de puerto al que se envía. Este sistema tiene el inconveniente de que ahora mismo, se están desarrollando multitud de aplicaciones web. Estas aplicaciones se conectan mediante el puerto 80, es decir, mediante el protocolo HTTP. Pero esto es solo teoría, pues se puede hacer que cualquier aplicación envíe información por el puerto 80, sin que sea HTTP, lo cual daría como resultado una mala identificación y posterior clasificación.

- Identificación basada en la aplicación de técnicas de aprendizaje automático sobre estadísticas de tráfico [11].

Esta técnica de identificación es muy sencilla. Consta de algoritmos de aprendizaje automático, los cuales se dedican a realizar análisis estadísticos. De esta forma los algoritmos aprenden que características tienen los flujos y los identifican en función de estas. Por lo tanto cuanto más tiempo pase, los patrones serán mejores y se tendrán más aciertos a la hora de la identificación.

2.1.3. DPI

- Identificación basada en el contenido del paquete o DPI. [10]

Esta técnica de identificación consiste en analizar los paquetes, más allá de su cabecera. La Inspección Profunda de Paquetes, DPI por sus siglas en inglés, *Deep Packet Inspection*, esta técnica lo que hace es analizar todos los paquetes en busca de cadenas que se repitan, además de las IP's y los puertos. Una vez que tiene una cadena, la usa para buscar paquetes que tienen ese patrón, de esta forma establece un patrón de identificación.

Esta técnica es útil para buscar virus, intrusiones y demás. Suele ser usada por los proveedores de servicios de Internet, ISP, *Internet Service Provider* y grandes empresas. Se puede decir que es una técnica que no respeta la privacidad, pues entra dentro de la propia información del paquete. Aunque en el caso de una gran empresa si puede hacerse, pues si se está en la red propia no es delito mirar la información que contienen los paquetes.

En algunos países, como Alemania, donde es ilegal el uso de *torrents*, las autoridades pueden pedir información al proveedor de servicios sobre uno de sus clientes, si cree que este está infringiendo la red, por lo que esta técnica es ideal para este tipo de cometidos.

Además de para lo descrito, esta técnica también es ideal para saber si se está siendo atacado y por quien. [2]

2.2. BRO

Bro [5], es un analizador de tráfico de red de código abierto. Al ser de código abierto permite ser usado por quien lo desee sin necesidad de pagar licencias. Funciona sobre sistemas basados en Linux y Mac OS X. Una de sus principales características es la gran cantidad de información que devuelve con un solo escaneo. Otros monitores de red devuelven menos información, teniendo que ser el propio administrador el que analice después la información obtenida por estos. Por lo tanto tardará más en resolver los posibles problemas que encuentre en la red.



Figura 2.1: Logo de Bro.

No tiene interfaz gráfica, por lo que su gestión se realiza desde la línea de comandos. Cuando se analiza tráfico, Bro generará unos registros, los cuales están divididos en función de parámetros definidos por parte del equipo que lo ha desarrollado.

Bro incorpora la posibilidad de introducir funcionalidades nuevas, mediante la programación en su propio lenguaje de *scripting*, del mismo nombre. Esto se desarrollará más adelante y en la parte de implementación se podrá ver cómo se desarrolla y algunos ejemplos de código escrito en Bro.

El lenguaje de *scripting* está orientado a trabajar con eventos. Por lo tanto a la hora de añadir una nueva funcionalidad habrá que crearla a partir del uso de los eventos que puede gestionar el programa.

Bro está estructurado de forma que todos los flujos de paquetes que analiza son procesados por el motor de eventos. Este motor convierte los flujos de paquetes en eventos, de forma que es más sencillo trabajar con ellos. Una vez que los eventos son procesados se generan los registros correspondientes, los cuales podrán ser analizados posteriormente. [13]

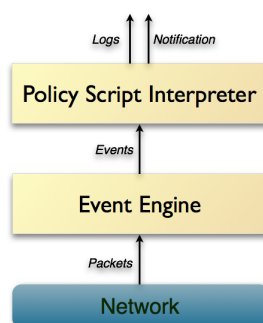


Figura 2.2: Arquitectura de Bro.

2.2.1. Funcionalidades básicas de BRO

La funcionalidad básica de Bro es la monitorización de la red en la que se ejecuta. Mientras que se encuentra en ejecución, genera registros o *logs* en texto plano que se podrán leer usando un editor de texto. Si se analiza un archivo *pcap* los *logs* no cambiarán. Sin embargo si se analiza tráfico a tiempo real, los registros se irán actualizando a medida que pase el tiempo. Algunos de los *logs* que se generarán son los siguientes.

- *dpd.log*. Consiste en un resumen de los protocolos encontrados en puertos que no son estándar.
- *dns.log*. Contendrá toda la actividad correspondiente al *DNS*.
- *ftp.log*. Un registro de la actividad a nivel de sesión de *FTP*.
- *files.log*. Un resumen con los archivos transferidos a través de una red. Incluye protocolos *HTTP*, *FTP* y *SMTP*.
- *http.log*. Registro de toda la actividad *HTTP* con sus réplicas.
- *ssl.log*. Un registro de las sesiones *SSL*, incluidos los certificados que se utilizan.
- *weird.log*. En este *log* se guarda la información correspondiente a actividad inesperada o rara a nivel de protocolo. Al analizar gran cantidad de tráfico no es muy útil, pues ocurren muchas cosas raras, pero a pequeña escala es bastante interesante para detectar intrusiones y demás.
- *conn.log*. Aquí se puede ver la información correspondiente a conexiones *TCP*, *UDP* e *ICMP*.

Para más información sobre *logs* generados por una monitorización de Bro lea [14].

Bro dispone además de varios *frameworks* que extienden su funcionalidad. Con ellos se podrá crear *scripts* muy potentes. Algunas de las utilidades de estos *frameworks* son las siguientes.

- *Geolocalización*. Se podrá encontrar la localización geográfica de una IP.
- *Análisis de ficheros*. El monitor de red tiene la capacidad de trabajar con ficheros.
- *Framework de loggins*. Con este *framework* se podrá extender los archivos de registro que se generan.

- *NetControl*. Este *framework* permitirá tener un control muy amplio y diverso del tráfico.

Para mayor conocimiento de estos *frameworks* y otros lea [15].

Ahora con toda esta información de los registros, el administrador del sistema podrá determinar si existen amenazas en la red o si hay algún componente defectuoso. Para ello deberá de hacer uso de los eventos con los que Bro trabaja, para sacar mayor partido a todo el trabajo que se realice sobre la red.

2.2.2. Eventos y trazas

La programación en Bro no es secuencial, no se escribe un *script* y se espera que se ejecute en el orden en el que está escrito. La programación aquí esta orientada a eventos. Por lo tanto el orden depende de cuando se ejecute un determinado evento en el análisis.

Un evento se da cuando Bro detecta un “comportamiento”, por ejemplo cuando detecta un paquete de respuesta *UDP*. Si se programa para que ese evento muestre un mensaje, este se mostrará cada vez que se detecte un paquete de respuesta *UDP*. Por lo tanto la programación deberá de estar pensada en función de los eventos disponibles.

Dentro de cada evento se podrá hacer lo que se desee, con la información que capture ese evento. Si el evento captura información de una conexión *TCP*, se tendrá que trabajar con esa información y las distintas variables globales que se hayan definido previamente.

Una traza es una captura del tráfico de una red. Con casi cualquier programa de diagnóstico se puede realizar capturas de red. En la propia web de Bro se encuentran algunos archivos de trazas de red, en formato *pcap*, de forma que se puedan realizar pruebas sobre ellos sin necesitar realizar una traza propia.

En el caso de Bro se podrá generar trazas de red mediante un *script* de *Python* [16], el cual se puede descargar desde el repositorio de *GitHub* de Bro. Con dicho script también se podrá trabajar sobre la traza, separando el tráfico entrante del saliente, desglosando en distintos registros el tipo de tráfico y demás.

2.2.3. Incorporación de funcionalidades

La incorporación de funcionalidades al monitoreo realizado por Bro es una característica muy llamativa. Gracias a esto se podrá realizar un análisis

muy personalizado usando un *script* creado por el administrador de redes. De esta forma podrá, por ejemplo, filtrar el tráfico de una determinada IP mientras se sigue analizando el tráfico de forma normal, con los registros que genera Bro de forma automática. Es una forma muy sencilla de comprobar si por ejemplo el servicio que administra está recibiendo demasiadas peticiones desde una misma IP. Lo cual sería un indicio de ataque de denegación de servicio.

A la hora de incorporar funcionalidades a Bro se puede hacer todo lo que se desee. Una búsqueda rápida por *GitHub* arrojará una gran cantidad de personas que contribuyen con una gran cantidad de nuevas funcionalidades [17]. Ahora lo ideal sería incorporar un módulo, de forma que si el resto de la comunidad lo desea pueda hacer uso de él de una forma sencilla.

2.3. Emparejamiento de flujos

La técnica de emparejamiento de flujos, fue planteada por el departamento de Teoría de la Señal, Telemática y Comunicaciones, *TSTC*, de la Universidad de Granada, en el año 2011 [3] [4].

En el emparejamiento de flujos, la idea de la que se parte es muy sencilla. Si dos paquetes comparten IP's de origen y destino, y puertos de origen y destino, se podrá decir que esos dos paquetes pertenecen a la misma clase. Una vez que están identificados como pertenecientes a la misma clase, lo que se debe de hacer es comparar los tiempos de esos paquetes, de forma que sean coherentes en el flujo. Si son muy lejanos en el tiempo, se podrá descartar que pertenezcan al mismo flujo. Esta idea se puede ver en [4] de una forma mucho más técnica.

En el mismo artículo [4] se encuentra la fórmula que se sigue para el emparejamiento.

$$F(x, y) = \begin{cases} G(x, y), & NIP(x, y) \geq 1 \\ -\infty, & \text{en otro caso} \end{cases}$$

Donde tenemos que $G(x, y)$ se corresponde con la siguiente función:

$$G(x, y) = |NIP(x, y) - 1| + 1 / (dp1(x, y) + k1) + 1 / (dp2(x, y) + k1) + 1 / (dt(x, y) + k2)$$

Las variables de la función son las siguientes:

- x : Primer paquete a comparar.

- y : Segundo paquete a comparar.
- $NIP(x,y)$: Es el número de paquetes analizados con las mismas características.
- $dp1(x,y)$: Se corresponde con los puertos de origen de los dos paquetes.
- $dp2(x,y)$: Serán los puertos de destino de los dos paquetes.
- $k1$: Es una variable definida previamente.
- $k2$: Es otra variable definida antes del comienzo del análisis. En trabajos anteriores esta variable y la anterior suelen estar definidas entre 1 y 10000.
- dt : Es la diferencia de tiempo existente entre los tiempos de inicio de los paquetes o *timestamps*.

Con esta fórmula se obtendrá como resultado un número, el cual será comparado con un umbral que se define previamente. Si el umbral definido es pequeño tendremos más flujos emparejados. Si el umbral es mayor serán menos los flujos que serán emparejados al no cumplir los requisitos.

El emparejamiento de flujos, por si mismo, sólo identifica el tráfico. Por lo tanto se necesitará otra técnica para clasificar el tráfico.

Capítulo 3

Diseño y arquitectura del sistema

En este capítulo se hablará cómo se va a estructurar el módulo, las diferentes funcionalidades que debe tener para que sea funcional y por último cómo Bro gestiona los flujos.

Todo esto se planteará desde la vista del diseño, se verán algunos tipos de datos interesantes para el desarrollo, aunque el código se verá más adelante en el capítulo de implementación. También se hablará de la arquitectura que va a seguir el sistema.

3.1. Arquitectura del sistema

Para describir la arquitectura del sistema hay que tener en cuenta la arquitectura que compone a Bro. Este monitor de red es un software modular, esto quiere decir que esta compuesto de diferentes módulos que al ser ejecutados funcionan como un único sistema.

Por lo tanto la arquitectura de este problema se debe de acoplar a la arquitectura propia de Bro. Para ello, se construirá un módulo que se acople con el monitor de red. Entre los requisitos del módulo a desarrollar se encuentra que sea ligero y eficiente. Por lo tanto se tendrá que usar los distintos eventos que proporciona Bro para crearlo de la forma deseada.

Se prescindirá del uso de los *frameworks* descritos en el capítulo anterior, pues su uso no aporta nada que no se pueda realizar exclusivamente con los eventos destinados a gestionar el tráfico de la capa de transporte. Se tendrá que usar los eventos para gestionar absolutamente todo lo que pasa en el módulo, a parte también se tendrá que usar una función propia. Esta función se encargará de aplicar la fórmula del emparejamiento de flujos vista en el

último apartado del capítulo anterior, esta función devolverá un número que será el que se compare con el umbral.

Por lo tanto, a modo de resumen, la arquitectura final del sistema que se plantea será un módulo compuesto por los distintos eventos destinados a la gestión del tráfico de la capa de transporte y una función auxiliar para el cálculo del umbral de emparejamiento.

3.2. Módulo y funciones

Dado que una mala programación del módulo puede dejar al monitor de red colgado con una traza sencilla. Se tratará que el módulo contenga solamente lo justo para llevar a cabo la identificación del tráfico. De esta forma se obtendrá un módulo cuya ejecución será óptima.

Las funcionalidades que se espera que tenga este módulo en esencia son dos, la detección y almacenado del tráfico y la aplicación de la fórmula para conocer si dos flujos son emparejables, a los distintos flujos que se han detectado. De una forma más amplia las funciones del módulo serán las siguientes.

- *Función que aplique la fórmula de emparejamiento.*

A esta función se le pasará dos flujos, de forma que se aplique la fórmula y devuelva un número, el cual será el que indique si los flujos son emparejables o no.

- *Funciones que detecten el tráfico.*

Esto se hará con los eventos de Bro. Los eventos detectarán el tipo de tráfico que se está analizando y aplicarán la función anterior.

Tras el uso de esta fórmula se almacenará o no el flujo que está siendo analizado. Por lo tanto será necesario el uso de algún tipo de contenedor para este cometido.

Lo que se espera es capturar el tráfico de la capa de transporte. Dicho tráfico se corresponde a los protocolos *TCP* y *UDP*. Por lo tanto será necesario ver que tipo de eventos son los que controlan el tráfico de estos dos protocolos. Esto se verá de forma más amplia en la siguiente sección.

Las entradas y salidas son de fácil gestión. Las entradas de tráfico podrán ser mediante archivos o analizando directamente el tráfico de la red. Las salidas, por su parte, serán mediante terminal. Esto puede suponer cierto inconveniente si se obtienen demasiadas salidas. También se pueden guardar las salidas en un fichero mediante el carácter *mayor que* . Con esto se guardará en un fichero en la ruta que se especifique, siendo el posterior análisis mucho más cómodo desde, por ejemplo, un editor de texto.

Se debe de tener en cuenta que siempre se podrá extender la funcionalidad del módulo. Pero de momento no resulta interesante. La posible extensión correspondería a posibles trabajos futuros.

Para realizar este módulo es necesario conocer como gestiona los flujos Bro y de que forma se mantendrán los que son emparejados y los que están activos.

3.3. Gestión de flujos

La gestión de flujos en Bro se realiza completamente con eventos. Por lo cual se tendrán que crear variables globales para el almacenamiento de los flujos que sean emparejables y los que estén activos.

El *nacimiento de un flujo* es controlado por un evento. Por lo tanto cuando se detecta un nuevo flujo se lanza un evento. Este evento se tendrá que controlar de forma que si se tiene ya un flujo activo con las mismas características, se compare y se almacene. Si por el contrario no se tiene ningún flujo con esas características se tendrá que almacenar directamente en el contenedor de flujos activos.

La *muerte de un flujo* también es controlada por un evento. Ahora lo importante es si es interesante a nivel del análisis seguir almacenando los flujos aunque estos hayan muerto. Si no se quiere tener almacenados flujos muertos se tendrá que eliminar de la estructura en la que está almacenado. Si se quiere seguir trabajando con ellos habrá que mantenerlos guardados en la estructura. Obviamente si el flujo que va a morir está emparejado con otro se borrará solo del contenedor de flujos activos. Manteniéndose en el contenedor de flujos emparejados. De lo contrario se perderá información. Todo esto habrá que decidirlo en el evento que gestiona la muerte del flujo.

Estos dos comportamientos de los flujos están controlados por eventos genéricos. Con un único evento se detecta que el flujo ha nacido y con otro evento si ha muerto, independientemente del tipo de protocolo al que pertenezca. No pasará esto con los distintos estados de los flujos que serán detectados. Pues no es lo mismo detectar un ACK de un flujo TCP que una

respuesta de un flujo UDP. Serán tratados en eventos distintos y tendrán que ser gestionados con eventos distintos, cada uno destinado a un protocolo distinto.

A continuación se verán las estructuras de datos necesarias para llevar a cabo el desarrollo del módulo.

3.4. Estructuras de datos

Las principales estructuras de datos que se necesitan para el desarrollo de este trabajo serán dos tablas, *table* [18], de vectores para el almacenamiento de los flujos activos y los emparejados. Los vectores son iguales que en cualquier otro lenguaje de programación, mientras que las tablas son parecidas a los *maps* de *C++*. Aunque esto se podrá ver con mejor detalle en el apartado de implementación.

Dichas estructuras de almacenamiento, deberán de ser capaces de estar ordenadas por las IP's y los puertos. Lo cual se obtiene juntando las tablas con los vectores para así conseguir una especie de matriz bidimensional, la cual está indexada. Por lo tanto se consigue que el acceso a los datos sea mucho más rápido. Incluso se podría prescindir de bucles, los cuales pueden acabar siendo un problema en cuanto a rendimiento si se llega a almacenar muchos flujos.

Bro proporciona cierto tipos de datos muy interesantes, los cuales además, incluyen mucha más información. Algunos de estos tipos de datos son.

- *connection*.

Este tipo de dato es el flujo en si. Por lo tanto será de vital importancia comprenderlo para poder trabajar como se desea.

Dentro del tipo de dato *connection* existe un registro llamado *id*, el cual esta compuesto por el tipo de dato *conn_id* [19]. Este dato sirve para identificar los flujos mediante una tupla formada por 4 datos. Estos datos son los que se precisan para indexar la matriz bidimensional, siendo pues las IP's y los puertos.

- *addr*. Este tipo de dato representa una IP. Reconoce tanto IPv4 como IPv6. Este tipo de dato puede ser comparado e incluso ordenado mediante operadores. [20]
- *port*. Este tipo es el usado para los puertos. Además del número de puerto también indica el protocolo de la capa de transporte que usa. Soporta la comparación y ordenación, pero no por el número, sino por el tipo de protocolo. [21]

Para obtener más información sobre el tipo de dato *connection* lea [22].

- *time*.

Este tipo de dato también es interesante. Aunque en otros lenguajes se puede obtener, en Bro es un tipo de dato por si mismo. Por lo tanto se podrá operar sobre él desde el principio, siendo una gran ventaja a la hora de calcular el tiempo de inicio de los flujos. Para leer más [23].

Es importante entender que para realizar el cálculo para el emparejamiento de flujos, se necesita el *timestamp* del primer paquete de cada flujo, pues será sobre este tiempo sobre el que se apoye el cálculo del emparejamiento.

Estos dos tipos de datos a parte de ser los más interesantes para el cálculo del emparejamiento, también serán los más utilizados junto a los contenedores para los flujos. Existen más tipos de datos e incluso los hay que extienden la información disponible sobre los flujos. Para leer más sobre esto [24].

Capítulo 4

Implementación

A continuación se contará cómo se ha implementado el módulo de Bro, resolviendo así el problema planteado.

La descripción del módulo será sin entrar en detalles de la programación. Se hará una descripción breve de los distintos eventos y funciones.

Lo primero que se va a detallar es la función para calcular el emparejamiento.

```
1 function emparejamiento(c1: connection , c2: connection ):double
```

Esta función recibe como entrada dos flujos y devuelve un número de tipo *double*. En esta función lo que se hace es aplicar la fórmula de emparejamiento a los dos flujos que entran.

Al querer hacer la comparación lo que se hace es sacar las IP's de origen y destino y los puertos de origen y destino de los flujos. También se obtendrán los *timestamps* de los flujos, de esta forma se conseguirá la diferencia de tiempo.

Para poder operar con los puertos habrá que pasarlos a tipo *count*, de esta forma se elimina la terminación con el tipo de protocolo del puerto. Se comentó anteriormente que se podía operar con esta terminación, pero al tener que operar con otros tipos de datos que no son puertos es mejor quitar la terminación, pues de lo contrario no se obtendrá un buen resultado.

Lo mismo que con los puertos pasa con el tipo *time*. Se puede operar con este tipo de dato pero no es recomendable hacerlo ya que se va a trabajar con más tipos de datos de carácter matemático.

El número de veces que se tiene un flujo con los mismos datos, o *Nip* en la fórmula, es de fácil cálculo. Gracias a la indexación basta con buscarlo en la tabla y calcular el tamaño del vector.

A continuación se procederá a la explicación de los distintos tipos de eventos usados y para que son usados.

```
1 event new_connection(c: connection)
```

Este evento recibe como entrada un flujo nuevo. Este flujo es una nueva conexión, la cual no está identificada previamente. Esto quiere decir o que es nueva o que ha sido borrada.

Este evento no devuelve nada, por lo tanto en el momento en el que es detectado el flujo se tendrá que crear un nuevo índice en la tabla de flujos activos o guardarlo para hacer la posterior comparación.

Este tipo de evento detecta las conexiones de tipo *TCP* y *UDP*.

Ahora se verá cómo se van a gestionar la muerte de los flujos.

```
1 event connection_state_remove(c: connection)
```

Este evento se activa cuando el flujo que entra como parámetro va a morir, o ser borrado de la memoria. Es un flujo que ya ha sido procesado por el módulo.

Lo que se realiza dentro de este evento es buscar en la tabla el índice el vector. De esta forma se borrará el primer flujo almacenado en el vector.

De ser el único flujo almacenado en el vector, se borrará el vector entero. Si hay más flujos almacenados se moverán los demás una posición hacia atrás. De esta forma se seguirá teniendo un rendimiento óptimo.

A continuación se verán los distintos eventos que van a detectar el diferente tipo de tráfico.

```
1 event connection_established(c: connection)
2
3 event connection_finished(c: connection)
4
5 event udp_request(u: connection)
6
7 event udp_reply(u: connection)
```

Los dos primeros eventos son los correspondientes al tráfico *TCP*. Los otros dos, como se puede ver en el nombre están destinados al tráfico *UDP*.

El primer evento relacionado con *TCP* se activa cuando se detecta un paquete *SYN-ACK* que responde al *handshake* que se realiza en las conexiones de este tipo.

El segundo evento detecta cuando la conexión *TCP* finaliza de forma normal.

Los dos eventos relacionados con *UDP* detectan paquetes de dos tipos distintos.

- *UDP request*. Se genera por cada paquete que es enviado por el creador del flujo.
- *UDP reply*. Este es generado por cada paquete que es enviado por el receptor del flujo.

Estos dos últimos eventos son bastante costosos, pero son absolutamente necesarios, son los únicos eventos que detectan conexiones de tipo *UDP*.

Dentro de todos estos eventos se realiza lo mismo. Primero se comprueba que los dos flujos no son el mismo, si son el mismo se termina el análisis. Si no son lo mismo se pasa a la función de emparejamiento ya descrita y se compara el número que se obtiene con el umbral que se ha definido antes de la ejecución. Si el resultado es mayor que el umbral son emparejables, por lo cual se tendrá que guardar en la tabla de emparejados y se informa mediante un mensaje en pantalla de que lo son. Si es menor que el umbral no son emparejables y mediante un mensaje en pantalla se informa de que no son emparejables.

Es necesario recordar que se puede dar que se puede estar usando dos eventos o más a la vez, por lo que los mensajes en pantalla pueden ser confusos. Se tiene que tener en cuenta que hay cierto retardo en los mensajes, por lo que lo interesante es el resultado final de las tablas.

Además de estos eventos se tienen los siguientes eventos genéricos.

```
1 event bro_init ()
2
3 event bro_done ()
```

El primero se lanza cuando Bro se inicia y mostrará el tiempo de inicio. El segundo se lanza cuando Bro finaliza, por lo tanto es el último evento que se lanzará y mostrará la hora de finalización. Con estos dos eventos se tendrá control de cuando fue lanzado el análisis y cuando finalizó.

Aparte de estos eventos, también existen dentro del módulo otros dos eventos destinados a detectar paquetes del protocolo *ICMP*. Este tipo de eventos son necesarios pues Bro los detecta igual que los de tipo *TCP* y *UDP*.

```

1 event icmp_echo_request(c: connection, icmp: icmp_conn, id: count, seq
  : count, payload: string)
2
3 event icmp_echo_reply(c: connection, icmp: icmp_conn, id: count, seq:
  count, payload: string)
```

El protocolo *ICMP* se usa para el control, enviando mensajes de error si un router o un host no son alcanzables.

Al igual que con los eventos de *UDP*, el *request* es enviado por el creador del flujo, siendo una petición. El *reply* es enviado por el receptor del *request*, por lo que se considera la respuesta del anterior. El funcionamiento es el mismo que en los anteriores eventos que gestionan el tráfico.

Aunque estos dos eventos tienen más parámetros de entrada que los eventos anteriores, para calcular el emparejamiento solo será usado el primer parámetro, el cual hace referencia al flujo. El trato del flujo dentro de los eventos es el mismo que el que se aplica a los eventos *TCP* y *UDP* anteriormente descrito.

A continuación se verá cómo se ha implementado las estructuras de almacenamiento de los flujos.

```

1 global collection: table[addr, addr, port, port] of vector of
  connection &synchronized;
2
3 global collection_added: table[addr, addr, port, port] of vector of
  connection;
```

Se trata de dos tablas globales, cuyo índice está constituido por las IP's de origen y destino y los puertos de origen y destino. Además como los índices son únicos, cada índice apunta a un vector, y dentro de dicho vector se almacenan los flujos ordenados dependiendo de cuando son detectados.

La primera tabla almacena los flujos que están activos. La segunda almacena los que ya están emparejados.

Capítulo 5

Evaluación y pruebas

Capítulo 6

Conclusiones y trabajo futuro

Bibliografía

- [1] James F. Kurose and Keith W. Rose. *Redes de computadores un enfoque descendente*, chapter 1. Pearson, 2010.
- [2] Dr. Thomas Porter. The perils of deep packet inspection. URL <https://www.symantec.com/connect/articles/perils-deep-packet-inspection>.
- [3] José Camacho, Pablo Padilla, F. Javier Salcedo-Campos, Pedro García-Teodoro, and Jesús Díaz-Verdejo. Pair-wise similarity criteria for flows identification in p2p/non-p2p traffic classification. 2011.
- [4] José Camacho, Pablo Padilla, Pedro García-Teodoro, and Jesús Díaz-Verdejo. A generalizable dynamic flow pairing method for traffic classification. 2013.
- [5] Bro Team. Bro indice, . URL <https://www.bro.org>.
- [6] Álvaro Maximino Linares Herrera. Repositorio del trabajo con bro. URL <https://github.com/Lynares/bro-flows>.
- [7] Bro Team. Descarga de bro, . URL <https://www.bro.org/download/index.html>.
- [8] Elvis Michael. La tarifa por hora de un programador. URL <http://pyme.lavoztx.com/la-tarifa-por-hora-de-un-programador-12286.html>.
- [9] Internet Assigned Numbers Authority (IANA). Service name and transport protocol port number registry. URL <https://www.iana.org/assignments/service-names-port-numbers/service-names-port-numbers.xhtml>.
- [10] Christian Fuchs. Implications of deep packet inspection (dpi) internet surveillance for society. URL <http://fuchs.uti.at/wp-content/uploads/DPI.pdf>.

- [11] Thuy T.T. Nguyen and Grenville Armitage. A survey of techniques for internet traffic classification using machine learning. URL <http://ieeexplore.ieee.org/document/4738466/?reload=true>.
- [12] Kim Davies. Una introducción a iana. URL <https://www.iana.org/about/presentations/davies-atlarge-iana101-paper-080929-es.pdf>.
- [13] Bro Team. Arquitectura de bro, . URL <https://www.bro.org/sphinx/intro/index.html#architecture>.
- [14] Bro Team. Logs de bro, . URL <https://www.bro.org/sphinx/script-reference/log-files.html>.
- [15] Bro Team. Frameworks de bro, . URL <https://www.bro.org/sphinx/frameworks/index.html>.
- [16] Bro Team. Trazas en bro, . URL <https://www.bro.org/sphinx/components/trace-summary/README.html>.
- [17] securitykitten. Finding beacons with bro. URL <https://gist.github.com/securitykitten/a7edcee0932c556d5e26>.
- [18] Bro Team. Tablas en bro, . URL <https://www.bro.org/sphinx/script-reference/types.html#type-table>.
- [19] Bro Team. Tipo conn-id de bro, . URL https://www.bro.org/sphinx-git/scripts/base/init-bare.bro.html#type-conn_id.
- [20] Bro Team. Tipo addr de bro, . URL <https://www.bro.org/sphinx-git/script-reference/types.html#type-addr>.
- [21] Bro Team. Tipo port de bro, . URL <https://www.bro.org/sphinx-git/script-reference/types.html#type-port>.
- [22] Bro Team. Tipo cennnection, . URL <https://www.bro.org/sphinx/scripts/base/init-bare.bro.html#type-connection>.
- [23] Bro Team. Tipo time, . URL <https://www.bro.org/sphinx/script-reference/types.html?highlight=time#type-time>.
- [24] Bro Team. Tipo conn, . URL <https://www.bro.org/sphinx/scripts/base/protocols/conn/main.bro.html#type-Conn::Info>.
- [25] Bro Team. Web de bro, . URL <https://www.bro.org/sphinx/intro/index.html>.
- [26] Bro Team. Instalación de bro, . URL <https://www.bro.org/sphinx/install/index.html>.

- [27] Bro Team. Función `get_port_transport_proto`, . URL https://www.bro.org/sphinx/scripts/base/bif/bro.bif.bro.html#id-get_port_transport_proto.
- [28] Bro Team. Analizadores de protocolos, . URL <https://www.bro.org/sphinx/script-reference/proto-analyzers.html>.
- [29] James F. Kurose and Keith W. Rose. *Redes de computadores un enfoque descendente*. Pearson, 2010.
- [30] James F. Kurose and Keith W. Rose. *Redes de computadores un enfoque descendente*, chapter 2. Pearson, 2010.
- [31] James F. Kurose and Keith W. Rose. *Redes de computadores un enfoque descendente*, chapter 3. Pearson, 2010.
- [32] James F. Kurose and Keith W. Rose. *Redes de computadores un enfoque descendente*, chapter 4. Pearson, 2010.
- [33] James F. Kurose and Keith W. Rose. *Redes de computadores un enfoque descendente*, page 55. Pearson, 2010.
- [34] Microsoft. Calidad de servicio. URL [https://msdn.microsoft.com/es-es/library/hh831679\(v=ws.11\).aspx](https://msdn.microsoft.com/es-es/library/hh831679(v=ws.11).aspx).

