# On detecting Internet-based criminal threats with XplicoAlerts: Current design and next steps

Carlos Gacimartín, José Alberto Hernández, Manuel Urueña, David Larrabeiti
Universidad Carlos III de Madrid, Spain
Avda. Universidad 30, E-28911 Leganés, Madrid, Spain
Contact author: {carlos.gacimartin}@uc3m.es

*Abstract*—**Criminals use more and more the Internet to plan their crimes. Hence, modern Police must be provided with powerfull threat detection tools to prevent crimes before they actually occur and as a means to provide evidence against criminals at court. In light of this, when monitoring suspect traffic generated by criminals, Deep Packet Inspection (DPI) tools must be combined with automatic threat detection techniques in order to filter out non-relevant information but trigger alarms when potential threats are detected.**

**This work shows an architecture and its implementation for such a combination between DPI tools with automatic thread detection techniques, and further proposes the next steps to follow in order to achieve a high-performance threat detection tool to be used by police officers.**

## I. INTRODUCTION

Criminals use more and more modern communication technologies and the Internet to plan their crimes. Thus, the police of the future will need to be provided with high-performance traffic monitoring tools by IT companies and experts. Essentially, such monitoring tools must provide a means to identify the potential threats that the traffic generated by suspects may carry. This poses a number of both technological challenges and legal issues. Concerning technological challenges, the traffic monitoring tools must operate at very high-speeds in some cases, and are also desired to must provide an easy interface to be used by IT non-experts, as it is often the case of police officers. Regarding legal issues, most EU countries require an authorisation by court to monitor the traffic of a given suspect.

A large number of open-source monitoring tools are available in the web, with wireshark being the most typical example. However, most of these tools only capture IP packets as they traverse a given link, and display their raw contents without any processing. Furthermore, for privacy reasons, such tools often do both source and destination address anonymisation, and even sometimes they remove the application layer contents. Hence, such tools are not valid for identifying potential threats and criminals, since the IP addresses or both source and destination computers are important parameters to be stored by the police.

Additionally, when the goal is to identify whether or not a given communication flow carries any potential threat, it is very necessary to take a look at the layer-five contents of all

packets from a given flow. Remark that a communications flow is the unidirectional set of packets characterised by a particular flow five-tuple (source address, source port, destination address, destination port, protocol).

In the literature, *Deep Packet Inspection* (DPI) tools are those which process the full contents of packets from a flow, from the link layer up to the application layer, and optionally extract the the application contents from the array of packets of a given flow. A number of open-source DPI tools are publically available in the Internet, however the number of application-layer protocols supported by them is currently very limited. For instance, MSNshadow [3] only provides support to msn traffic, whereas tcpick [7] does not support email decoding.

Essentially, the number of challenges in designing and developing DPI applications are many-fold:

- The number of protocols and applications in the Internet is so high and changes so quickly that is impossible to decode all captured traffic.
- The reconstruction of the application-layer contents split into a variable number of IP packets, each of them including layer-two to layer-five headers requires the decoder to keep track of a large number of communication aspects, such as sequence numbers, IP fragmentation, TCP error control, etc.
- Applying decoding rules for a particular traffic flow requires to accurately determine the flow's application-layer protocol. This is often related to port number identification (for instance, port 80 is often web traffic, port 21 is FTP traffic, and so on) but the relationship (port number, application protocol) is not always a one-to-one maping.
- The information bytes may belong to different media, i.e. text, images, video, audio, programming code, etc, and each media may have been encoding with different techniques, for instance, base64 for MIME attachments in an email, iso-8859-1 for text, etc.

In spite of the so-many challanges, the Xplico tool [1], currently in its 0.5.5 version, has shown good performance features and a wide range of layer-five protocols supported, while promised by the developers to support a wider number of applications in the forthcoming versions. In addition to this, Xplico also provides a very intuitive web-based interface very appropriate for their use and configuration/administration by

IT non-experts, as it is often the case of Police Officers.

However, the Xplico tool was designed to monitor and decode traffic, but not to generate any alarm when the decoded traffic matches a given set of parameters associated to criminal threats. Thus, when the amount of decoded traffic becomes large, the manual inspection of each decoded flow is impractical, and it becomes necessary to make the tasks of threat identification automatic. To this end, the authors of this work have developed an extension of the original Xplico, called XplicoAlerts, which automatically detects whether or not a given communication flow contains any suspicious file allocated in a database, say for instance a terrorist- or paedophilia-related image. This work reviews the XplicoAlerts extension to Xplico, its design criteria and performance operation concerning the detection of suspicious traffic, and its applicability to the FP7 INDECT project.

## II. DEEP PACKET INSPECTION TECHNIQUES AND XPLICO

### A. A comparison of DPI techniques

Table I shows a comparison of a number of popular DPI techniques together with their most important features, such as the application protocols that they support and the platforms on which they can operate. As shown, the current version of Xplico decodes many more protocols than any other, and is expected to further decode packets collected at 802.11 wireless networks (as long as the keys are provided) and IRC, XMPP in subsequent versions.

### B. Inside Xplico

Xplico has been developed modularly in order to increase its maintenance and ease cooperation with other developers. Xplico comprises four macro-components, building a modular architecture as shown in Fig. 1:

- A Decoder Manager, called DeMa. This module is in charge of organising the incoming packets into flows, and launch and control the execution of the IP decoder and manipulators.
- An IP/network decoder, comprising a set of data manipulators (dissectors), one per application protocol, to increase modularity of the system.
- A dispatcher which outputs the results to several output formats and storage systems (directory tree, SQLite, mySQL, etc).
- A web-based visualisation system, called XI, which displays the decoded data using a easy-to-use php-based web interface.

The operation of Xplico is as follows: Incoming traffic feeds the DeMa, which organises the individual packets into flows. The DeMa then stores the packets from each flow in a separate .pcap file per flow, identified by source and destination IP addresses, port numbers and protocol. The DeMa then creates an instance of the IP/network decoder to process each flow separately.

The IP/network decoder comprises a number of data manipulators or dissectors one per protocol. For instance, there are layer-2 dissectors, like PPP and Ethernet, IPv4 and IPv6
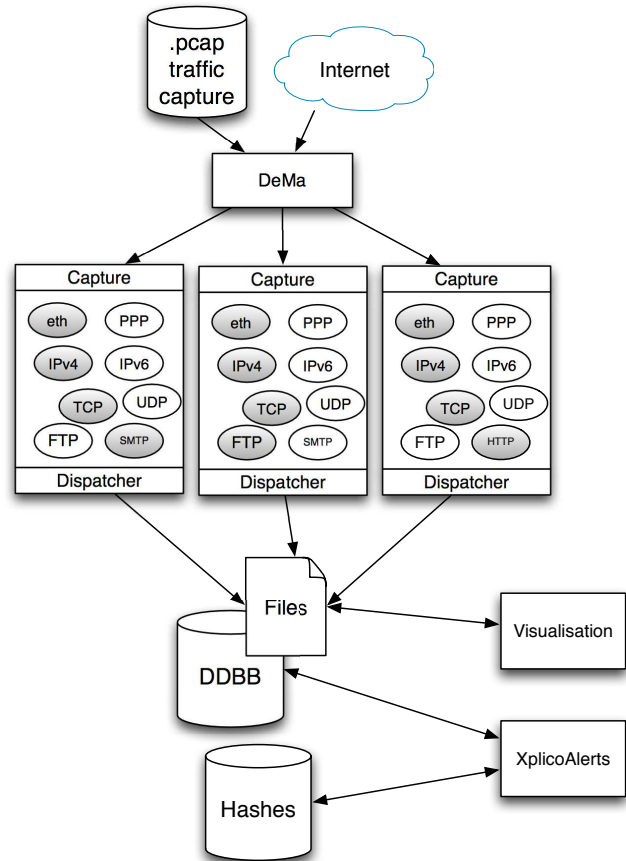


Fig. 1. Architecture of the Xplico tool and XplicoAlerts

dissectors, TCP and UDP dissectors, and finally a number of dissectors for the application-layer protocols, including SMTP, HTTP, telnet, etc. The IP/network decoder tries different dissectors from botton to top in the TCP/IP layer stack. For instance, in a traffic capture of web traffic collected in an Ethernet LAN, the first dissector tried is Ethernet, then IPv4, then TCP and finally HTTP to produce an output of the web page downloaded, which is then passed to the dispatcher. If none of the dissectors is suitable for a particular flow, then no output (other than the original .pcap flow) is passed to the dispatcher. For instance, in the example of Fig. 1, the DeMa generated three instances of the IP/network decoder, one for processing an email (thus using the eth, IPv4, TCP and SMTP dissectors), another one for processing a file transmitted through FTP, and a final one for processing a web page downloaded. After the flow is decoded, the dispatcher stores the results in a given previously-configured format (SQLite and directory tree, for instance).

Finally, the web-based visualisation system, which is based on php, displays the decoded data in a more friendly interface. This interface provides a main menu on the left-hand side with the decoded emails, SIP conversations, web sites visited and images transferred, among others. Clicking on each menu displays a list with information concerning the date and time

| Tool | Version | Platform | Language | Protocols |
|---|---|---|---|---|
| Xplico[1] | 0.5.5 | GNU/Linux | C, Python, php, javascript | HTTP, SMTP, POP3, IMAP, DNS, FTP, SIP, TFTP, IPP, PJL, MMSE, Telnet, NNTP, Facebook chat, Webmail AOL, hotmail, yahoo |
| PacketoMatic[2] | 20100227 | GNU/Linux, FreeBSD, Solaris | C | HTTP, POP3, MSN, IRC, RTP |
| MSNShadow[3] | 0.3beta | GNU/Linux | C | MSN |
| Pyflag[4] | 0.87 | GNU/Linux | Python, C | HTTP, SMTP, POP3, DNS, SIP, MSN, IRC, Yahoo Chat, Webmail Gmail, Google Image Search |
| tftpgrab[5] | 0.2 | GNU/Linux | C | TFTP |
| Chaosreader[6] | 0.94 | GNU/Linux, Windows98 | Perl | HTTP, SMTP, FTP, IRC, Telnet, VNC, ICMP, SSH |
| tcpick[7] | 0.21 | GNU/Linux, FreeBSD, Mac OSX | Ruby | HTTP, FTP |
| Yahsnarf[8] | | GNU/Linux | Ruby | Yahoo chat |

TABLE I

A COMPARISON OF OPEN-SOURCE DPI TOOLS

when the flow was captured, and other relevant information that summarises the contents of such flow. For instance, in Fig. 2, where all the emails captured and decoded are shown, such relevant information comprises the Subject of the email, the sender and the receiver, and the email size. Clicking on a given email displays its contents (see Fig. 3).
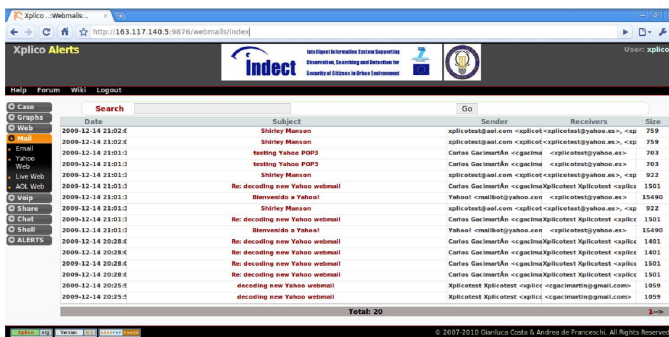


Fig. 2.  Xplico screenshot: Emails decoded



Fig. 3.  Xplico screenshot: Email contents

Next section explains how XplicoAlerts operate.

*C. XplicoAlerts*

The authors at Universidad Carlos III de Madrid have extended the original features of Xplico, providing an interface to quickly detect specific files within the flows collected from a given suspect. Recall from Fig. 1, XplicoAlerts combines information from two different databases: One database provides the monitored and decoded traffic from the suspect; the second database provides the hashes of a set of suspicious files. XplicoAlerts triggers an alarm when a decoded file matches

any suspect file. To do so, XplicoAlerts provides an interface to work with :

- Crime categories: New categories may be added by the Police Officer, etc.
- Include suspicious files in each category: Include the files themselves, the file hashes only or a directory tree with multiple files to be hashed and added.
- Alarms: Generate information regarding the alarm, including the flow that triggered the alarm, date and time and a description field for the police officer to type extra comments.

A screenshot of the index page of XplicoAlerts is given in Fig. 4. Fig. 5 shows how to add filehashes of new suspicious files.
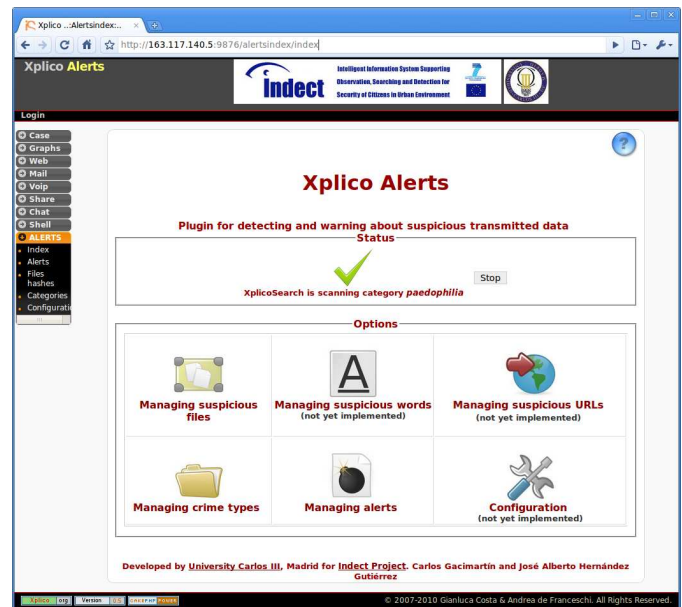


Fig. 4.  XplicoAlerts screenshot: Index

## III. A CASE EXAMPLE: DETECTING A THREAT

*A. Design requirements*

Before designing a tool for detecting criminal behaviour with Deep Packet Inspection techniques, it is necessary to take into account a number of requirements and constraints.

Fig. 5. XplicoAlerts screenshot: Adding filehashes

All these were taken under consideration in the design of XplicoAlerts:

- The analyser must store all information from the suspect. Basically, if the suspect sends an email which could be identified as a threat, he may then claim that after it he sent another email canceling the previous one. The Police must be able to check whether he sent a second email or not.
- The Police Officer does not need to check all decoded traffic manually, since this would be impractical when the amount of transferred information is large. Instead, the analyser must generate alarms only on those flows which it may consider suspicious. The Police Officer would then examine the flows that generated the alarm.
- Additionally, the Police Officer must have a means to type notes on the flows that generated alarms.
- Detecting suspicious flows shall be based on whether or not the flow contains a given suspicious file, e.g. an image or other type of object transmitted over an email, FTP, etc.
- The actual suspicious files do not need to be locally stored on the analyser, but only their hashes. The files trasferred by the suspect will be checked against the hash database. This increases the security of the analyser in case of loss.
- All decoded information must be stored locally, in an encrypted database, in order to increase security in case of loss.
- The suspicious hash database may be updated remotely. The analyser must be able to communicate safely and securely with the Police Office.

*B. Example of operation*

In a typical scenario, a Police Officer places the XplicoAlerts monitoring system somewhere near the house of a suspect from some criminal activity. Xplico then collects all traffic from the suspect and checks whether or not the decoded objects match any of the suspicious files, whose hashes are stored in a local (encrypted) database. If possitive, then an alert is sent to the Police Officer, who must check it manually.

In this scenario, the suspect is browsing a website which contains a suspect image (Fig. 6(a)), whose hash is included in the suspicious-file database. Xplico then decodes the full web session (Fig. 6(b) top-right). Then, XplicoAlerts hashes each decoded object and then compares the results with the hash database of suspicious-files. In case of a possitive, XplicoAlerts triggers an alarm (Fig. 6(c)), which must be analysed and anotated by the Police Officer (Fig. 6(d)).

## IV. NEXT STEPS

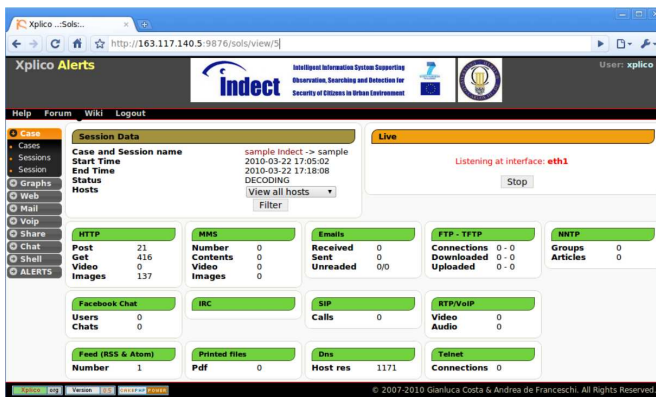The following comprises a set of development enhancements to be included in further versions of XplicoAlerts:

- Fuzzy hashing. With conventional hashing, when a file is slightly changed, then the hash changes completely. Thus, a suspect could slightly modify the files to make them undetectable by the analyser. Fuzzy hashing provides a mean to produce similar hashes of similar files (images), thus enabling to detect the transmission of images similar to those stored in the suspicious database.
- Natural language analysis. A further extension is to analyse not only the files transmitted but also to do some intelligent processing on the message contents (i.e. text of websites and emails, etc). To do so, it is necessary to store an array of words in the suspicious database such that, if some of them are identify within the same context (20 or 50 word distance), then an alarm is triggered to be manually inspected by the Police Officer.
- Host location. It is interesting to correlate the suspicious database with geographical information of the source and destination IP addresses, such that the Police Officer may identify the country or region that the suspect communicates with.
- Fast detection with Bloom Filters. As the number of suspicious items is expected to grow dramatically, an easy way to speed up the processing of files would include the use of Bloom filters. Bloom filters are compact data units that can be used to fast checking the membership of an item within a set.
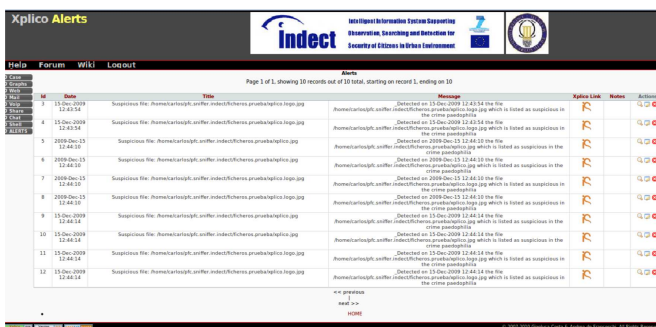
## V. SUMMARY AND DISCUSSION

This work has shown the architectural design of Xplico and XplicoAlerts and their potential to detect criminal activity by monitoring and decoding the traffic of suspects and then compare the results in a hash database of suspicious files. At present, the tool is capable of decoding a wide variety of protocols, but the detection features are still limited to the detection of suspicious files. However, a number of enhancements are planned for further versions, including natural language processing, threat location and fuzzy hashing -based detection.
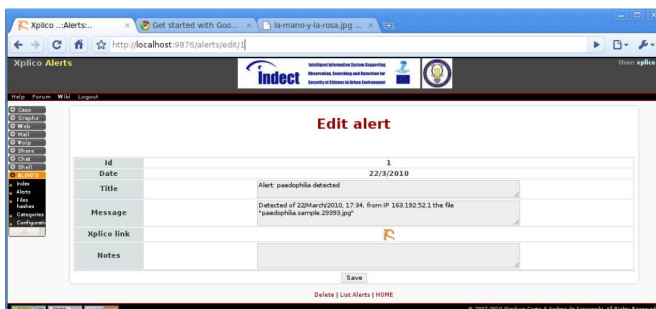
(a) Suspect browsing the web



(b) Traffic captured by Xplico



(c) Suspicious image detectec by XplicoAlerts



(d) Manual anotation by Police Officers

Fig. 6. A suspect browsing the web (a), all traffic captured by Xplico (b), suspicious image detected by XplicoAlerts (c) and manual anotation by the Police Officers (d)

REFERENCES

[1] "The Xplico Internet decoder", current version 0.5.5. Available at http://www.xplico.org.
[2] "PacketoMatic", current version 20100227. Available at http://www.packet-o-matic.org
[3] "MSN Shadow", current version 0.3beta. Available at http://msnshadow.blogspot.com/
[4] "Pyflag", current version 0.87. Available at http://www.pyflag.net
[5] "TFTPgrab", current version 0.2. Available at http://pseudo-flaw.net/content/tftpgrab/
[6] "Chaosreader", current version 0.94. Available at http://chaosreader.sourceforge.net/
[7] "Tcpick", current version 0.21. Available at http://tcpick.sourceforge.net
[8] "Yahsnarf". Available at http://writequit.org/projects/yahsnarf/