

# Hierarchical object log format for normalisation of security events

Andrey Sapegin, David Jaeger, Amir Azodi, Marian Gawron, Feng Cheng, Christoph Meinel

Hasso Plattner Institute (HPI), University of Potsdam

P.O.Box 900460, 14440 Potsdam, Germany

{andrey.sapegin, david.jaeger, amir.azodi, marian.gawron, feng.cheng, christoph.meinel}@hpi.uni-potsdam.de

**Abstract**—The differences in log file formats employed in a variety of services and applications remain to be a problem for security analysts and developers of intrusion detection systems. The proposed solution, i.e. the usage of common log formats, has a limited utilization within existing solutions for security management. In our paper, we reveal the reasons for this limitation. We show disadvantages of existing common log formats for normalisation of security events. To deal with it we have created a new log format that fits for intrusion detection purposes and can be extended easily. Taking previous work into account, we would like to propose a new format as an extension to existing common log formats, rather than a standalone specification.

**Keywords**—log normalisation, intrusion detection, common log format

## I. INTRODUCTION

Nowadays, every application has its own log file format. Such a variety of formats complicates log analysis [1] and causes problems for system administrators, developers of intrusion detection systems (IDS) and security analysts [2], [3]. Therefore, several solutions [4] were introduced aiming to solve this problem: CEE [3], CEF [5], IODEF [6].

However, there are many developers who avoid using common log formats in their projects. This statement particularly applies to intrusion detection systems. Many software vendors and open source projects still use their own log format for IDSs [7], [8], [9]. The reasons for this situation hide in the IDS architecture. To operate rapidly, such systems need to analyse huge amount of logs from hundreds of devices and services on the fly. Therefore, the log format should be lightweight and able to encapsulate all network and host logs from large datasets.

Log file standards that aim to handle all possible log output of any service, do not fit for these purposes. Moreover, in some cases, the variety of standardised fields, offered by such a common log file format, does not provide a suitable schema for normalisation and correlation of security events. During the development of our own IDS —Security Analytics Lab [10]— we have faced the same challenges and came to the same solution: a log format that is specially designed for IDSs. However, the object log format we present in this paper noticeably differs from the existing solutions.

Our goal was to experiment with a flexible lightweight format to optimise attack detection. First, we analysed a variety of large log files and came up with a list of common fields that can be found in the logged events. Then, this list was reduced

to the fields that we considered most important for the security analysis of occurred events.

As a result, the developed log format facilitates the detection of attacks and simplifies their correlation. We demonstrate the ability of our format to simplify attack detection within a case study. Besides this particular use case, the proposed log format retains to be adjustable and could be utilised for generic normalisation concepts.

The remainder of the paper has the following structure: Section II describes the data set used, Section III evaluates and compares existing log formats, Section IV provides details on the log format we developed and Section V tells about results of a case study. We discuss our results and conclude in the Section VI.

## II. DATA

To prove our concepts, we utilised data from “Scan of the Month” HoneyNet Challenge: “Scan 34 - Analyze real honeynet logs for attacks and activity” [11]. The HoneyNet Project [12] aims to help security specialists to sharpen their forensics and attack detection skills. To achieve it, the project shares real attack data collected from honeynets from all over the world. The “Scan of the Month” challenge provides archives with log files from different hardware and software systems for the postmortem analysis. The “Scan 34” challenge contains log files being collected between 30 January and 31 March 2005. We chose this particular challenge because it implies analysis of multiple log files for different services and includes various heterogeneous attack traces. We describe the log files in Table I.

The “Scan 34” challenge also provides the description on the attacks and important events in the log files. See the details of the 15 most important events in Table II.

In total, four attacks led to the server’s compromise. Two of four attackers used the vulnerability of ‘awstats.pl’ installed after the server reboot, probably by a system administrator. Another two attackers successfully brute-forced the SSH password. The first attacker has also installed an IRC bot right after the intrusion. Other events listed in the table were not harmful for the server; they represent suspicious activity that did not result into a successful intrusion.

## III. EXISTING LOG FORMATS FOR NORMALISATION OF SECURITY EVENTS

Before we decided to use our own log format, we have evaluated existing solutions, trying to estimate if they could

Server	Service	Date	Total number of lines
n/a	HTTP Server	Jan 30 - Mar 16	3925
bridge	iptables firewall	Feb 25 - Mar 31	179752
bastion	snort IDS	Feb 25 - Mar 31	69039
combo	syslog	Jan 30 - Mar 17	7620

TABLE I. LOG FILES AVAILABLE TO SOLVE THE CHALLENGE

N	Event	Date	Details
1	Reboot	Feb 11, 2005	n/a
2	Software installation	Feb 25, 2005	AWSTATS installed
3	Server compromised	Feb 26, 2005	Code injection through awstats.pl
4	Server compromised	Mar 04, 2005	Code injection through awstats.pl
5	Server compromised	Mar 06, 2005	ssh brute-force successful
6	Server compromised	Mar 13, 2005	ssh brute-force successful
7	Software installation	Feb 26, 2005	IRC bot installed by an attacker
8	ICMP alert	n/a	ICMP Destination Unreachable
9	Slammer worm	n/a	Worm propagation attempt
10	IIS attacks	n/a	WebDAV search access, cmd.exe access, etc.
11	SMTP scan	n/a	POLICY SMTP relaying denied
12	Typot trojan	n/a	trojan traffic
13	RPC scan	n/a	RPC portmap status request
14	Port scan	n/a	NMAP -sA (ACK scan)
15	Slapper worm	n/a	Worm propagation attempt

TABLE II. SECURITY RELATED EVENTS THAT OCCURRED DURING THE MONITORING FOR THE CHALLENGE

fit for our purposes of normalisation of security events. Please see the overview of proposed log formats in Table III.

The selected formats have different structure, size and even usage purpose. CEE is a generic format for logging any types of events [3], IDMEF was developed for exchange of security messages [13]. IODEF is specially designed for computer security incidents [6], as well as CEF developed as a part of intrusion detection system [5]. Although the described formats have different purposes, each of them could be used for log normalisation.

However, not every part of the log message could be normalised to the standard fields defined in every format. To deal with it, each format allows to use extra fields or extend its structure. Therefore we decided to check how many changes we need to fully normalise log messages we have into each of 4 formats. The Table IV shows, that some parts of log messages have no corresponding field in the log format. This small observation reveals a deeper problem, if investigated. The reviewed log formats offer limited number of options to parse the textual event description precisely. For example, using the IDMEF, messages like “authentication failure” and “Relaying denied. IP name lookup failed” are supposed to be written into the Classification class as whole. Such single field available in the log format does not allow to effectively normalise descriptions like “BLEEDING-EDGE WORM Mydoom.ah/i Infection IRC Activity [Classification: A Network Trojan was detected]”. CEE format demonstrates a similar issue, as all descriptions and other details are supposed to be stored in the “message” field. Since the information still could be stored in such generic fields, we do not mention it in the table as parts without corresponding field.

Another common issue affecting all log formats implies complications while normalising the event details like user name, rule/action applied, method (e.g., GET) and response

(404 or 550). In other words, the existing formats offer a limited number of fields to specify, what has happened. CEF deals with it better than other formats, due to the presence of keys like ‘act’ (Action mentioned in the event), ‘app’ (Application protocol), ‘request’, ‘requestMethod’ and others.

Finally, all formats have too few – or do not have any – fields for expression of properties specific for intrusion detection, like relation between events, security metrics or classifiers such as CVE [14] and CWE [15] identifiers, object and subject of the event. These properties usually could not be extracted from the log lines itself, but are expected to be filled by IDS.

Thus, not only missing fields from the table should be added to each of formats, but also the fields to improve parsing of message parts and allow to store intrusion detection properties. That’s why we decided to select and re-engineer one of the formats. However, the structure, number of fields/attributes and possible values or contents are different for all formats, which make them hard to compare with each other. Therefore, we chose the following additional criteria to examine log formats:

- *scalability*. There should be an ability to add custom fields, if some log formats could not be normalised into standard ones.
- *light weight*. The log format should have reasonable number of fields/attributes to avoid redundancy and fit for high-speed normalisation purposes.
- *multilevel schema*. Intrusion detection implies correlation of the normalised logs. The hierarchical or other connected structure is preferable as it will allow to categorise different fields into classes.

Applying the first criterion, all selected formats are scalable to some extent. IDMEF [13] and IODEF contain an *Additional*

Format name	Organisation	Size estimation	Format structure
CEE (Common Event Expression)	MITRE	58 fields, 7 objects	two levels, hierarchical, object-based
IDMEF (Intrusion Detection Message Exchange Format)	IETF	118 elements, 5 core classes, 53 attributes	multi-level, class-based
IODEF (Incident Object Description Exchange Format)	IETF	53 elements, 19 top-level classes, 83 attributes	multi-level, class-based
CEF (ArcSight Common Event Format)	HP	104 keys	one level, key-value pairs

TABLE III. EXISTING STANDARDS FOR LOG FORMATS

Log line	CEE	IDMEF	IODEF	CEF
81.181.146.13 - - [15/Mar/2005:05:06:53 -0500] "GET //cgi-bin/awstats/awstats.pl? configdir= —%20id%20— HTTP/1.1" 404 1050 "-" "Mozilla/4.0 (compatible; MSIE 6.0; Windows 98)	HTTP/1.0, GET, 404	Mozilla/4.0, 404	81.181.146.13, GET, //cgi-bin/awstats/awstats.pl, 404, Mozilla/4.0	-
Mar 15 13:38:03 combo sshd(pam_unix)[14490]: authentication failure; logname= uid=0 euid=0 tty=NODEV ssh ruser= rhost=202.68.93.5.dts.net.nz user=root	-	-	14490, 202.68.93.5.dts.net.nz, user=root	14490
Mar 1 20:45:12 bastion snort: [1:2001439:3] BLEEDING-EDGE WORM Mydoom.ah/i Infection IRC Activity [Classification: A Network Trojan was detected] [Priority: 1]: TCP 11.11.79.67:2568 -> 129.27.9.248:6667	TCP	-	Priority: 1, 11.11.79.67, 129.27.9.248	Priority: 1
Mar 24 19:46:50 bridge kernel: INBOUND ICMP: IN=br0 PHYSIN=eth0 OUT=br0 PHYSOUT=eth1 SRC=63.197.49.61 DST=11.11.79.100 LEN=32 TOS=0x00 PREC=0x00 TTL=111 ID=1053 PROTO=ICMP TYPE=8 CODE=0 ID=512 SEQ=29421	ICMP, eth0, eth1, br0	-	SRC = 63.197.49.61, DST = 11.11.79.100, eth0, eth1, br0	br0
Feb 1 10:08:32 combo sendmail[32433]: j11F8FP0032433: ruleset=check_rcpt, arg1 = <china9988@21cn.com>, relay=[61.73.94.162], reject=550 5.7.1 <china9988@21cn.com>... Relaying denied. IP name lookup failed [61.73.94.162]	ruleset = check_rcpt, 550, china9988@21cn.com	ruleset = check_rcpt, 550, china9988@21cn.com	check_rcpt, relay=[61.73.94.162], 550, 61.73.94.162	32433, ruleset = check_rcpt

TABLE IV. PARTS OF LOG MESSAGES WITHOUT CORRESPONDING FIELD IN THE LOG FORMAT.

Data class, CEF supports *Custom Dictionary Extensions* and CEE allows to create custom events and fields.

Analysing the second criterion, i.e. the number of fields/attributes, the CEE format is much more compact than IODEF, IDMEF and CEF. We also would like to mention, that IDMEF is designed for data exchange for intrusion detection and security management systems only. So we could expect a lot of additional attributes being added to IDMEF to adopt it for normalisation purposes.

Finally the last criterion is examined. The CEF has a flat structure; IODEF and IDMEF have similar multi-level schemas with connections between different classes; and CEE has a clear object-based hierarchical structure.

Based on the criteria defined, we selected CEE to be extended as the log format that fits better for our purposes.

#### IV. OBJECT LOG FORMAT

As mentioned in the previous Section, the decision to develop our own log format for IDS came with the attempt to utilise the CEE format in our experimental intrusion detection system [10].

We started with modifying the CEE format and adding custom fields to cover all possible log message variants from

the Honeynet challenge. However, this draft log format still had several problems. First, we have used less than half of fields (24 of 58) defined in the standard. Compared to the number of custom fields added – 59 – and the fact that standard CEE fields do not always present the key properties of the log message, it becomes hard to argue for using the resulting format inside the IDS system. Second, the object-field hierarchy defined in the Field Dictionary [16] contains only one abstraction level. Taking into account overlapping semantics (e.g. *appname* and *app.name* fields) in the CEE notation, this relatively flat structure contains many similar fields and adds a lot of confusion for a developer. To handle with our goals, we made changes to the CEE structure as well. We extended the hierarchical structure of CEE to three levels to achieve flexibility and improve clarity.

We present the proposed format in Figure 1. The first level (marked with bold) describes global parameters or classes of parameters, such as *network* or *original\_event*. On the second level of our hierarchy (written in normal font) we describe the most significant properties. And on the third level (marked with italics) we place specific information such as network protocol fields.

Compared to the formats examined in the Section III, the proposed format offers multiple fields to store event details

## Object Log Format

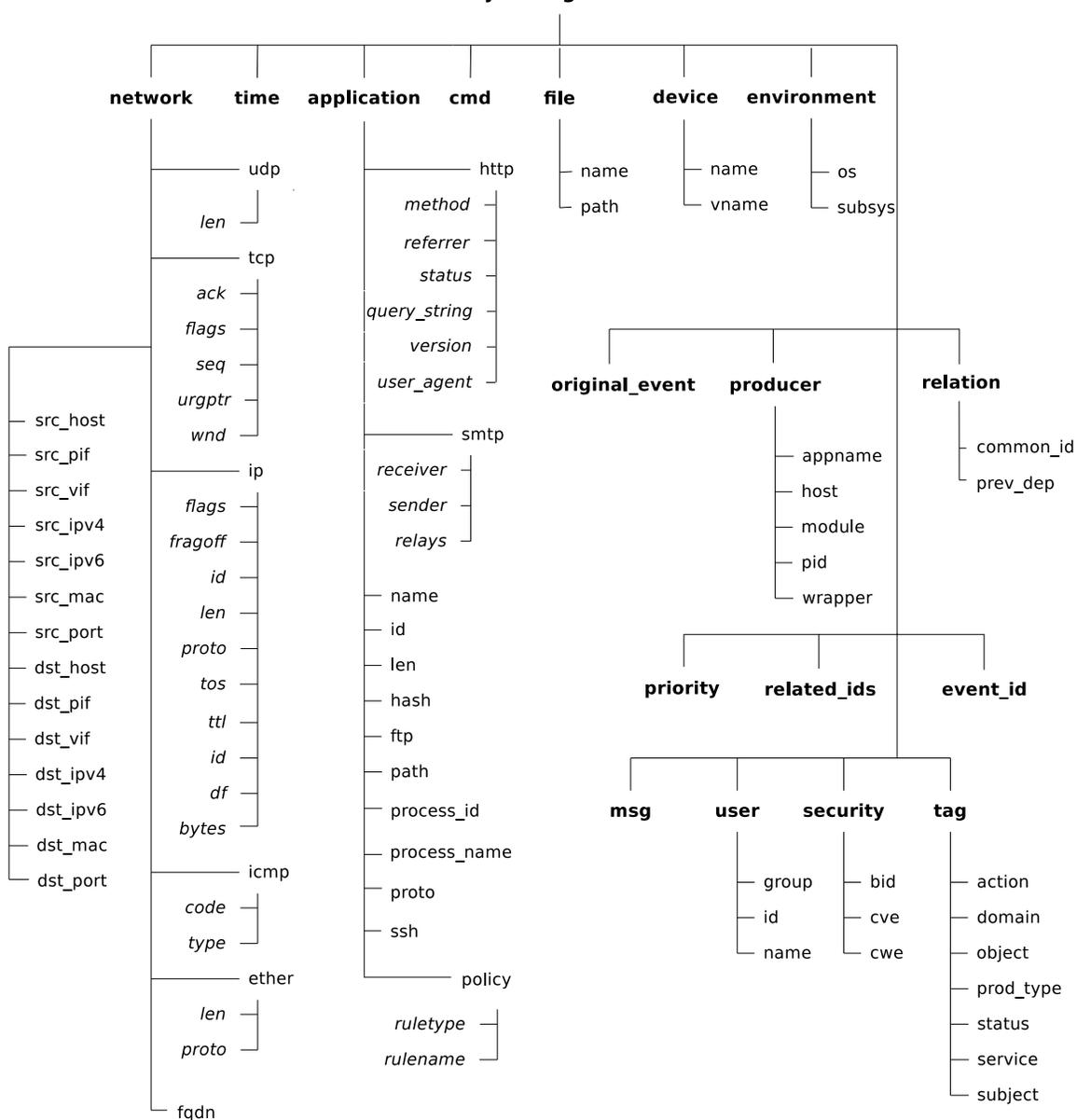


Fig. 1. Tree of properties in our log format

describing what has happened (*tag* and *application classes*), as well as fields, which are highly relevant for intrusion detection (mainly, classes *related\_ids*, *relation* and *security*).

Let's consider the example from the Table IV to show how to parse real log data into the proposed format:

```

81.181.146.13 - - [15/Mar/2005:05:06:53 -0500] ``GET //
cgi-bin/awstats/awstats.pl? configdir= \\%20id%20|
HTTP/1.1'' 404 1050 ``-'' "Mozilla/4.0 (compatible
; MSIE 6.0; Windows 98)
  
```

The log line listed above is taken from the 'access\_log' file on the HTTP server. Using it, we now describe the most significant elements and how the information from our example log entry should be distributed over them.

- network** This class covers all properties around the lower network layers, i.e. the link, network and transport layer of the TCP/IP protocol stack. It describes information about the source and destination endpoints of network events, including their MAC address, IP address and ports. The protocol fields in captured network packets are organised in subclasses, such as *ether*, *ip*, *icmp*, *tcp* and *udp*. These details can then be used to analyse the exact workings of a network communication. It should be noted, that since we are parsing security events and not only the network packets, this class could provide details for multiple network packets associated with an event, even if they were reported in different logs. From the log line sample, we extract only

the IP address—81.181.146.13—to fill the ‘network.src\_ip\_v4’ field.

- **application** This class covers all properties of applications and services involved in the event. The application class can represent different kind of involved applications in an event. In a network connection, such an application could be the client application, which initiates a connection or sends a packet to the server, or the service application, which is the target of a sent packet or initiated connection. To cover the various special characteristics in service communications, the *application* class also allows to further specify parameters in the context of HTTP, FTP and more. In a host-based event, this application is usually the application that initiates an action, e.g. an application that writes a file. The semantics of the specified application can be obtained from the *subject* and *object* field of the *tag* class.

Analysing the example log line, it is possible to notice, that the application is a user agent, so the ‘application.http.user\_agent’ field should be filled with “Mozilla/4.0 (compatible; MSIE 6.0; Windows 98)” line and ‘application.http.query\_string’ field – with the following line:

```
“GET //cgi-bin/awstats/awstats.pl? configdir=
—%20id%20— HTTP/1.1”.
```

- **producer** This class gives information on the application that observed and eventually persisted an event. It should describe the first application that persisted the events and should not be changed to one of the intermediate processing applications. In our example, a producer is a HTTP server, namely “Apache HTTP Server”, which should be written into the ‘producer.appname’. If we would know some details about the host from the other log lines, we could also fill the ‘producer.host’ field, e.g. with *http\_server*.
- **file** This class describes the files that were involved in an event. A file can appear in different contexts, i.e. mainly as a data source and target of access operations, such as read and write. In the case of an FTP or HTTP connection, this parameter could give information on the accessed resource. Similar to the *application* class, the concrete context for one event is defined by the *subject* and *object* field of the *tag* class. Considering the sample log line, the ‘file.name’ field will contain ‘awstats.pl’ and ‘file.path’ – ‘//cgi-bin/awstats/awstats.pl’.
- **original\_event** This field keeps the original log as found in the log source. This is usually a string containing full log line or fields from log database. The whole sample log line should be filled into the ‘original\_event’ in our case.
- **relation** This class serves the cases when several events are related to each other. In the case of multiple events sharing a common identifier for correlation, this identifier is stored in the *common\_id* field. The *prev\_dep* property indicates if the current event depends on a previous one. For the example considered, these fields will be empty,

because the logged event is standalone and not related to any other events.

- **tag** This class provides abstract information on top of the message details, mainly to categorise and tag events. As this information is not always directly represented in the log data, this is required to be set by the user. In some cases the action (and other fields as well) cannot be explicitly identified with a single term. Therefore, we propose all fields within this class as multi-value. E.g. the ‘tag.action’ could be {*get,access*}, ‘tag.subject’ – *host* or *user\_client*, ‘tag.object’ – *file* or *web\_document*, ‘tag.prod\_type’ – *web\_server* and ‘tag.service’ could be *web*.
- **security** This class provides links to identifier of vulnerability, related to the logged event. For example, ‘security.cve’ could be “CVE-2010-4369”.
- **event\_id** This field contains a unique internal id of the event. This should be unique among all generated events in a management system. For our example, we used an SQL database for automatic id generation.

The developed format structure is easy to present, extend and map into database relations. These features simplify the developer’s tasks and clarify the semantics of fields with similar names. For example, now the former *appname* and *app.name* fields are easier to distinguish by using *producer.appname* and *application.name* as names. Finally, with 107 fields used, we were able to effectively normalise every log message from the dataset. After the normalisation step, most of the attacks could be discovered using simple search queries, as shown in the next section.

## V. THE ROLE OF THE COMMON LOG FORMAT IN ATTACK DETECTION

As mentioned earlier, we normalise files before searching for attacks. This pre-processing step includes parsing of logs into described log format with regular expressions and inserting them into SQL database.

The proposed log format allows to detect attacks described in the Table II using simple queries, without the usage of a correlation engine or other advanced intrusion detection techniques. The challenge winners on the other part have used self-written scripts and manual analysis [17], [18]. We now provide several use cases to demonstrate the benefits of the proposed log format. E.g. to check for SMTP scans, we use the following query:

```
select * from event where application_protocol = 'smtp'
and tag_status = 'failure'
```

All 220 log lines returned correspond to event 11 from Table II. These lines include both *sendmail* messages from the mail server (“combo”) and *snort* alerts from the logging server (“bastion”).

Next, to detect the code injection attempts, we suggest another simple SQL query:

