

D113 : une plateforme open-source dédiée à l'analyse des flux et à la détection des intrusions

David Pierrot, Nouria Harbi

Université Lumière Lyon 2, Laboratoire ERIC 69676 BRON Cedex, FRANCE
{david.pierrot1,nouria.harbi}@univ-lyon2.fr

Résumé. Ce travail se situe dans le domaine de la "Cybersécurité", le projet "D113" permet de visualiser en temps réel les flux transitant sur des équipements de filtrage sans avoir recours au traitement manuel des journaux d'événements. Nous centrerons notre démonstration sur la visualisation de grands "graphes" et l'exploitation d'analyses statiques des flux.

1 Introduction

L'explosion d'internet, couplée à l'effet de la mondialisation, a pour résultat d'interconnecter les personnes, les entreprises, les états. Le côté déplaisant de cette interconnexion mondiale des Systèmes d'Information réside dans un phénomène appelé "Cybercriminalité". Des personnes, des groupes mal intentionnés ont pour objectif de nuire dans un but pécuniaire ou pour une "cause", aux informations d'une entreprise, d'une personne voire d'un Etat. Il n'est pas rare que des faits de "cyber-attaques" soient relatés dans les médias envers des grandes sociétés comme "Google", "Visa", "Sony", "Apple". La sécurité d'un Système d'Information se doit d'être présente afin de garantir la confidentialité, l'intégrité, la disponibilité de l'information. De ce fait, il existe une multitude d'équipements de sécurité qui permettent de détecter les comportements anormaux. Un des principaux équipements de sécurité est le "Pare-Feu"¹ ou plus communément appelé "Firewall". Il a pour mission comme le décrit Al-Shaer et Hamed (2003) de filtrer, selon une politique fondée sur les flux autorisés à pénétrer dans un réseau selon leurs sources, leurs destinations et les services souhaités (navigation internet, transfert de fichiers, etc...). Par son positionnement, il donne une visibilité totale de l'ensemble des flux. Cet équipement offre aussi la possibilité "d'historiser" vers des journaux les flux ayant été autorisés ou interdits. L'exploitation et l'analyse des journaux d'événements liés aux équipements de sécurité sont devenues primordiales pour la maîtrise des flux et la détection des intrusions ainsi que pour la vérification du bien fondé de la politique de filtrage mise en place (Golnabi et al, 2006). Dans ce contexte, les constructeurs d'équipements de filtrage mettent à disposition des logiciels permettant d'analyser les flux. Ces derniers nécessitent un accès et une connaissance dudit équipement. La détection des anomalies et des comportements anormaux est conséquemment réservée à ces seuls utilisateurs. La problématique de la représentation des événements de sécurité est tellement répandue que plusieurs outils ont même été regroupés au

1. équipement de sécurité basé sur le filtrage des entrées/sorties des flux réseau

sein d'une distribution "Open Source" nommée "Davix"². Des outils comme Gephi³, NVisio-nIP de K.Lakkaraju et al. (2004) ou PicViz⁴ sont inclus dans cette dernière. Mais il s'avère que ces outils demandent une installation et ils n'ont pas été créés pour une interrogation via des clients légers. De facto, ils nécessitent un déploiement et une préparation des données plus ou moins importante pour permettre une visualisation de l'activité réseau.

A notre connaissance et à ce jour, il n'existe pas d'outil "packagé"⁵ permettant de réaliser ce type de tâches. Il convient donc de proposer une solution vulgarisant l'analyse des flux et d'être en mesure d'obtenir un résultat probant sans une quantité importante de moyens. La solution "D113" permet aux personnes en charge de la sécurité⁶ d'un système d'information de visualiser l'ensemble des flux transitant sur un équipement de filtrage en facilitant la lecture des événements et la prise en compte des anomalies. "D113" peut également être utilisé pour des diagnostics ou des levées de doute⁷. Le principe est de modéliser un système de "monitoring" et de visualisation des données réseau en temps réel permettant de détecter rapidement les tentatives d'intrusions.

2 Composition du projet "D113"

Le projet "D113" est composé de quatre phases qui s'inscrivent dans le cadre d'un travail de thèse en sécurité s'appuyant sur des données issues des différents équipements et outils de sécurité. Les différentes phases du projet se déclinent selon la liste suivante.

- Phase 1 : "Monitoring et visualisation" des données réseau, représentation graphique des activités d'un réseau informatique via un modèle de données.
- Phase 2 : "Extraction des profils d'attaque", phase qui s'appuiera sur des méthodes de Data Mining.
- Phase 3 : "Scoring" des risques et phase d'évaluation.
- Phase 4 : Détermination d'un plan d'actions.

Notre démonstration de logiciel portera uniquement sur la **phase 1** qui est le préambule à la "fouille de données" qui sera effectuée dans les phases suivantes. La première phase constitue un tout en soi dans la mesure où la visualisation des données pour les utilisateurs est un enjeu crucial en termes de prise de décisions sur les problématiques de sécurité.

3 Contexte d'application

L'architecture étudiée concerne une entreprise publique dans le domaine de la santé composée de 92 000 employés. Notre étude porte sur 3 réseaux interconnectés au sein d'un réseau étendu WAN⁸ distants géographiquement. Ces derniers sont tous pourvus d'équipements de

2. the Data Analysis & Visualization Linux : <http://www.secviz.org/node/89>
3. <https://gephi.github.io>
4. <http://www.picviz.com>
5. Client léger : navigateur web fondé sur une architecture Linux Apache Mysql Php
6. Responsables de la sécurité du système d'information, ingénieurs sécurité, administrateurs réseau
7. Vérification que l'équipement de filtrage n'est pas la source du problème lors d'un dysfonctionnement d'une application
8. Wide Area Network : réseau informatique étendu couvrant une grande zone géographique

filtrages, l'objectif est d'être en mesure d'analyser les événements liés aux règles de filtrage via une exportation vers un conteneur de données. Afin de simplifier les références aux différents réseaux, le nommage suivant sera utilisé :

- Site de production : **SP1**
- Site de qualification : **SQ1**
- Site d'administration distante et bureautique : **SAB1**

Le réseau SP1 propose des services à destination de **14 millions** de personnes. Les données peuvent être considérées comme sensibles et portent sur une quantité de 9.2 Teraoctets et plusieurs dizaines de millions d'euros par jour. Ces données sont hétérogènes et proviennent de plusieurs sources différentes. Les transactions réseaux filtrées par le "Firewall" représentent plus de 6000 lignes par minute en matière d'événements. Le second et le troisième sites sont respectivement utilisés pour des tâches dites de "bureautique-administration distante" et de qualification (reproduction des infrastructures de production).

Ces trois sites sont opérationnels, c'est à dire que les données traitées et analysées dans les sections suivantes correspondent à des données de production. Pour des raisons de **confidentialité** les adresses IP ont été anonymisées. Le réseau SP1 est doté de son propre conteneur de données qui est alimenté par les événements envoyés en temps réel par le "Firewall". Les réseaux SAB1 et SQ1 mutualisent un même conteneur. Les données brutes envoyées par l'ensemble des équipements filtrants sont traitées selon une extraction de motifs.

Description des données

Le contenu des variables listées ci-dessous sont exportées vers des conteneurs de données.

- adresse ip source
- adresse ip de destination
- port de destination
- protocole (udp et tcp)
- date et heure de la connexion
- numéro de la règle du pare feu correspondant aux flux

Le tableau 1 synthétise le volume en nombre de lignes traitées par les équipements de filtrage.

	flux traités par journée	moyenne par minute
SP1	9 886 928	6 865
SAB1	572 272	397
SQ1	20 670	14

TAB. 1: Flux traités par SP1, SQ1, SAB1 en nombre de lignes

4 Collecte et traitement des données

Afin de pouvoir produire un résultat graphique, il convient de détailler les opérations et configurations mises en place. Les différents moyens décrits dans cette section permettent de réaliser le pré-traitement. Ces derniers dispensent, à partir des événements bruts envoyés par

D113 : monitoring et analyse des flux

les équipements de filtrage, un format exploitable via le conteneur de données. La figure 1 présente en détail les différentes étapes de conversion des données.

Les différentes traces de connexions des équipements de filtrage (option LOG) sont envoyées au serveur SYSLOG-NG⁹ au format brut similaire à la figure 2. Ce dernier, via ses options intégrées de reconnaissance PCRE¹⁰ et de filtrage, dépose l'ensemble des flux traités sur un serveur de bases de données (MYSQL¹¹). Par la suite, un traitement est réalisé via un script Perl¹², créé par nos soins, ayant pour objectif de préparer le résultat des différentes requêtes. Enfin, les programmes issus de la suite Graphviz¹³ et du script Afterglow¹⁴ sont utilisés pour la création des graphiques comme démontré par Marty (2008). L'utilisateur n'a en fait besoin que d'un navigateur Internet pour être en mesure de visualiser les flux.

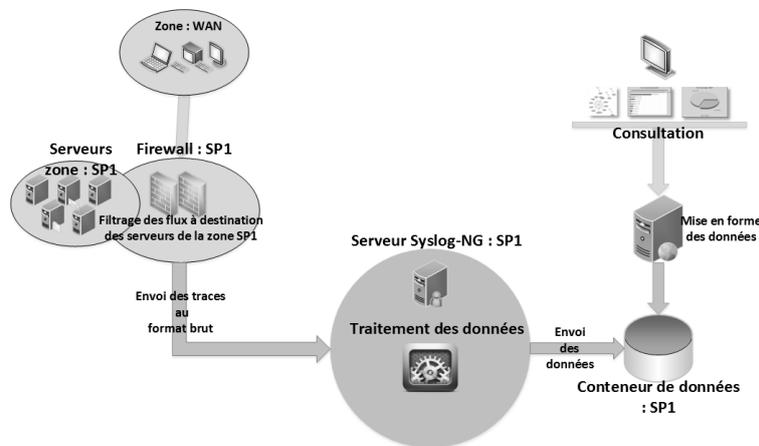


FIG. 1: Schéma traitement des événements

```
Aug 12 11:54:51 10.250.0.2 Site1b-FW1: NetScreen device_id=Site1b-FW1 [root]system-notification-00099(traffic): start_time="2014-08-12 11:54:51" duration=0 policy=16-26 service=SP proto=IP src_zone=Site1b-bureau@que-zone dst_zone=Site1b-interco-man-zone action=Permit sent=0 rcvd=0 src=172.250.8.12 dst=172.127.8.65 src_port=5060 dst_port=5060 src-xlated ip=172.250.8.12 port=5060 dst-xlated ip=172.127.8.65 port=5060 session_id=25145 reason=Creation
Aug 12 11:54:51 10.250.0.2 Site1b-FW1: NetScreen device_id=Site1b-FW1 [root]system-notification-00099(traffic): start_time="2014-08-12 11:54:51" duration=0 policy=16-26 service=TCP port=445 proto=6 src_zone=Site1b-bureau@que-zone dst_zone=Site1b-interco-man-zone action=Permit sent=0 rcvd=0 src=10.250.9.19 dst=10.11.9.142 src_port=1310 dst_port=445 src-xlated ip=10.250.3.19 port=1310 dst-xlated ip=10.11.9.142 port=445 session_id=28054 reason=Creation
Aug 12 11:54:51 10.250.0.2 Site1b-FW1: NetScreen device_id=Site1b-FW1 [root]system-notification-00099(traffic): start_time="2014-08-12 11:54:51" duration=0 policy=16-26 service=TCP port=5060 proto=6 src_zone=Site1b-bureau@que-zone dst_zone=Site1b-interco-man-zone action=Permit sent=0 rcvd=0 src=172.250.8.13 dst=172.143.12.62 src_port=5060 dst_port=5060 src-xlated ip=172.250.8.12 port=5060 dst-xlated ip=172.143.12.62 port=5060 session_id=34959 reason=Creation
Aug 12 11:54:51 10.250.0.2 Site1b-FW1: NetScreen device_id=Site1b-FW1 [root]system-notification-00099(traffic): start_time="2014-08-12 11:54:51" duration=0 policy=16-26 service=TCP port=8443 proto=6 src_zone=Site1b-bureau@que-zone dst_zone=Site1b-interco-man-zone action=Permit sent=0 rcvd=0 src=10.250.0.183 dst=10.11.45.128 src_port=62829 dst_port=8443 src-xlated ip=10.250.0.183 port=62829 dst-xlated ip=10.11.45.128 port=8443 session_id=36682 reason=Creation
```

FIG. 2: Flux bruts en provenance d'un équipement de filtrage

5 Scénario de visualisation

La représentation graphique de l'ensemble des flux autorisés selon la période souhaitée relève du problème de vision de grands graphes (voir figure 3), mais il est possible d'extraire

9. Gestion des journaux, <http://www.balabit.com/network-security/syslog-ng>
10. Perl Compatible Regular Expression, <http://www.pcre.org>
11. Base de données open Source, <http://www.mysql.com>
12. <https://www.perl.org/>
13. Logiciel de visualisation graphique, <http://www.graphviz.org>
14. AfterGlow, outil de génération graphique, <http://afterglow.sourceforge.net/>

des "sous graphes" basés du "requêtage" qui visent à sélectionner les modalités de certaines variables (adresses source et de destination ainsi que les services et protocoles). En revanche, l'analyse d'un graphique des flux rejetés (même agréés), comme le montre la figure 4, s'avère simple. Une adresse IP tente de se connecter à plusieurs autres adresses sur le port "135". Une recherche de répertoires partagés peut être à l'origine de ce type de comportement.

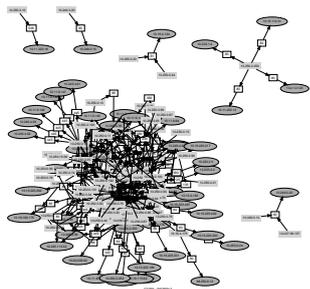


FIG. 3: Exemple de transactions

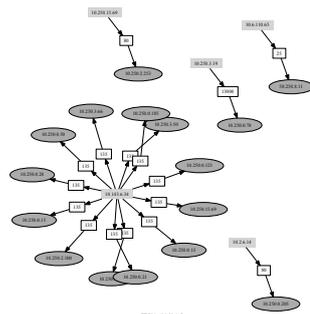


FIG. 4: Exemple de transactions rejetées

Si une activité comme le montre la figure 5 est constatée, un balayage de ports est certainement à l'origine de cet affichage. L'interface de visualisation propose plusieurs options d'exploration des données, qui permettent de corroborer ce constat comme illustré dans la figure 6. Les modalités généralement constatées ont été largement outrepassées, passant de 9 rejets par tranche de 30 secondes à 3287. Une surveillance peut être mise en place sur l'actif visé. Afin d'avoir

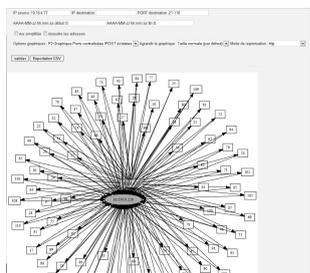


FIG. 5: Balayage de ports

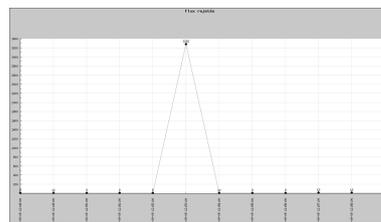


FIG. 6: Flux rejetés par tranche de 30 secondes

une vision critique et surtout externe par rapport aux travaux entrepris sur le projet "D113", un avis et une utilisation ont été demandés à plusieurs personnes occupant des postes relatifs à la sécurité des systèmes d'information dans plusieurs entreprises. L'installation du produit étant relativement aisée et l'utilisation simplifiée par l'utilisation d'un navigateur web, le constat est plus que positif. Il est facile de visualiser les événements et les diagnostics n'en sont que plus simples à déduire.

6 Conclusions et perspectives

A l'issue de la phase 1, l'ensemble des événements liés au filtrage est exporté en temps réel vers des conteneurs de données. Dans un souci de performance et compte-tenu de l'importance

de la volumétrie des données, nous avons opté pour une possibilité de "requêtage" rétroactif de deux jours. Ce choix de durée est aussi fondé sur les différents échanges réalisés avec les experts qui estiment qu'au delà de cette période, l'analyse n'est plus forcément nécessaire. La lecture graphique permet de visualiser rapidement les tentatives de connexion depuis plusieurs sources vers plusieurs destinations. Il peut s'agir d'une tentative d'intrusion ou d'une prise de renseignement ou encore d'une mauvaise configuration d'un script.

Il convient de poursuivre les phases 2 à 4 à savoir l'extraction des profils d'attaques, établir un "scoring" et la création d'un plan d'actions basé sur l'application des méthodes de Data Mining sur un espace de représentation graphique. Ces phases permettront de générer des règles d'association en fonction des actifs visés selon plusieurs vecteurs d'attaques. Une prise en compte de techniques issues de l'apprentissage supervisé et non-supervisé permettrait une meilleure considération des attaques avec une définition automatique des seuils de signalement. Il conviendrait de créer un système évolutif et adaptatif en temps réel permettant en fonction des changements intervenant sur un Système d'Information d'offrir une véritable aide à la décision.

Références

- Al-Shaer, E. et H. Hamed (2003). Firewall policy advisor for anomaly detection and rules. *Integrated Network Management, 2003. IFIP/IEEE Eighth International Symposium on.*
- Amanpreet, H. et al (2011). Survey on data mining techniques in intrusion detection. *International Journal of Scientific Engineering Research 2.*
- Bhruyan, H. et al (2011). Survey on incremental approaches for network anomaly detection. *International Journal of Communication Networks and Information Security 3.*
- Deepa, A. J. et V. Kavitha (2012). A comprehensive survey on approaches to intrusion detection system. *ICMOC-2012.*
- Golnabi, K. et al (2006). Analysis of firewall policy rules using data mining techniques. *Network Operations and Management Symposium, 2006. NOMS 2006. 10th IEEE/IFIP.*
- K.Lakkaraju et al. (2004). Nvisionip: netflow visualizations of system state for security situational awareness. *VizSEC/DMSEC '04 Proceedings of the 2004 ACM workshop on Visualization and data mining for computer security.*
- Marty, R. (2008). *Applied Security Visualization.*

Summary

Our present work is in the area of "Cybersecurity", the "D113 project" allows you to view real-time flow on the filtration equipment without the need for manual processing of event logs. We focus our demonstration on the visualization of large "graphs" and use static analysis of flow.