

# Entre urgence et exhaustivité : de quelles techniques dispose l'analyste pendant une réponse sur incident ?

Amaury Leroy  
amaury.leroy@airbus.com

Airbus Group CERT

**Résumé** Cet article présente quelques techniques d'analyse qui ont fait leurs preuves en situation réelle, dans le cadre d'une réponse sur incident ou d'une recherche de compromission, de type APT ou SPT (*Simple Persistent Threat* pour les intimes), sur tout type de parcs. Ce domaine étant vaste, cet article s'attachera à définir les techniques employées par un analyste en charge de l'étude des logs proxy et des flux réseaux, appelé en urgence sur une suspicion de compromission. Ces techniques seront présentées de manière chronologique au fur et à mesure de l'investigation, en commençant par l'urgence des premières heures, et en finissant avec les réflexions des dernières semaines.

## 1 Introduction

Appelé en urgence suite à la découverte d'un malware sur un serveur sensible ou la présence de fichiers d'archives sur Internet, l'analyste est très souvent livré à lui-même dans le cadre d'une réponse sur incident ou d'une recherche de compromission. L'obligation de résultat et le stress le plongent dans une course effrénée contre le temps. La recherche de comportements suspects sera rythmée par des analyses plus ou moins complexes et rentables. Ainsi, les premières semaines, l'analyste s'orientera en priorité sur une approche *low-hanging fruit* [49], avec des analyses rentables et faciles pour dégrossir la masse de travail et découvrir un maximum d'attaques. Pendant les semaines suivantes, ayant une meilleure idée de la compromission, il s'attachera à mettre en place des analyses et alertes pour surveiller l'attaquant. Pour finir, des analyses plus complexes seront étudiées pour s'assurer que rien n'a été oublié.

### 1.1 Préparer et connaître la guerre

Préparer et être prêt à la crise fait déjà partie de la réponse sur incident. Pour assurer une intervention rapide, l'analyste doit avoir revu et validé les

méthodes d'analyse, de travail et les outils<sup>1</sup>. Comme chaque intervention offrira son lot d'imprévus et de surprises, notre analyse doit pouvoir surmonter chaque nouveau problème rencontré. Ainsi il préférera se baser sur plusieurs briques libres ou faites maison, pour s'assurer de la flexibilité et du contrôle de ses solutions et outils. Concernant les méthodes d'analyse et de travail, il sera indispensable de respecter et appliquer les bonnes pratiques et méthodologies forensic : traçabilité des actions, main courante, modèles d'analyse et de rapport, structuration et consolidation des données dans une solution de partage et d'échange, comme un Wiki. Ceci dans l'objectif de structurer, consolider et partager les informations avec les équipes impliquées : du forensics système jusqu'aux administrateurs et managers.

La durée d'incident pouvant se compter en semaines, mois ou années, la masse d'informations à étudier peut devenir très importante et sa gestion complexe. Ainsi, notre analyste devra s'efforcer de :

- **réduire la quantité de données.** L'utilisation de listes blanches (sites connus, internes, partenaires, etc.) permet de réduire la quantité de logs, mais cette opération peut fausser l'analyse (pivots<sup>2</sup> internes ou d'entreprises partenaires, compromission d'un site légitime, etc.). Pour y remédier, une analyse des logs évincés par la *whitelist* sera possible dans un second temps ;
- **compresser.** Si possible compresser les données pour diminuer la taille et améliorer les performances ;
- **éviter les intermédiaires.** Connecter le support de données au plus proche du système d'analyse. Éviter les copies au travers du réseau et connecter le disque directement sur son système. Dans certains cas, l'envoi d'un disque dur, par courrier, est plus rapide que de le transfert via un VPN, même en tenant compte le temps pour chiffrer les données sur le disque et l'expédition de celui-ci. Dans le cas de gros fichiers, il sera toujours préférable de commencer à analyser les logs directement sur le support fourni, au lieu de faire la copie sur le système d'analyse et ensuite lancer l'indexation/analyse.

---

1. Il sera trop tard pour découvrir, apprendre ou inventer les techniques, outils ou méthodes pendant la crise.

2. Le terme « pivot » dans cet article définit une machine compromise, servant de relais entre 2 réseaux.

## 2 Les premières semaines

Les premières semaines sont primordiales à l'analyste pour détecter et esquisser l'ampleur d'une compromission. La connaissance de l'entreprise : son architecture, ses applications et ses pratiques, les bonnes et surtout les mauvaises, sera un atout essentiel pour détecter des comportements suspects. Une vérification des informations découvertes avec les équipes métiers sera régulièrement nécessaire lors de l'analyse de l'incident. Par exemple, l'envoi de 30 Go de données vers des systèmes hébergés dans un pays exotique pendant le week-end peut être tout à fait normal dans le cas d'un laboratoire qui envoie ses résultats à une équipe partenaire.

### 2.1 Bases d'IOCs privées ou publiques

Le terme IOC *Indicator of compromise* [39] désigne un *artéfact forensic*, détectable sur le réseau ou sur un système d'exploitation, qui indique la présence d'une infection ou attaque donnée. Généralement, les IOC sont des adresses IP, domaines, *hashs*, URLs, etc., qui indiquent la présence d'un malware spécifique ou l'utilisation de C&Cs (*Command and Control*) [38] donnés. Plusieurs formats existent pour consolider et échanger les IOCs, comme OpenIOC [26], STIX [35], IODEF [16] (RFC5070), XML, etc. Malheureusement, peu de bases publiques d'IOCs sont disponibles (Iocbucket [21]) et généralement, les IOCs sont à extraire directement des rapports proposés par les sociétés d'antivirus ou de *Cyber Intelligence*. Heureusement, plusieurs projets de partage d'IOCs, comme le *framework Open Source* MISP [11] et des communautés d'échange [14] ont vu le jour pour permettre le partage et la consolidation des IOCs.

Face à l'urgence, une comparaison entre les IOCs privées/publiques et les informations à notre disposition (logs proxy, captures réseau, etc.) sera une première étape pour débroussailler la quantité d'informations et offrir peut-être les premières pistes.

*La quantité de données* La compromission pouvant s'étaler sur plusieurs semaines, mois ou années, les recherches sur plusieurs téraoctets de logs ou captures réseau peuvent être longues. Pour parvenir à optimiser les recherches et son précieux temps, l'analyste doit comprendre les goulots d'étranglement et les minimiser. Pour des traitements de recherche de motifs ou statistiques légers tels que les IOCs, les entrées-sorties sont le facteur limitant. Rien ne sert d'investir dans des machines contenant des centaines de cœurs et To de RAM car le traitement par le CPU ou la place utilisée dans la RAM seront négligeables devant l'attente des informations

lues sur le disque. Pour cette raison, il faut éviter de lire plusieurs fois la même donnée depuis le disque dur. Ainsi, il faudra utiliser de préférence des logs compressés pour maximiser la quantité d'informations utiles récupérées par chaque accès au support et d'autre part utiliser un moteur d'indexation pour minimiser le nombre d'accès aux fichiers.

Il existe de nombreuses solutions d'indexation et de recherche pour le traitement en masse, comme Splunk ou RSA SA. Toutefois il est également possible de bricoler, pour des analyses ponctuelles, avec la solution libre ELKS<sup>3</sup> qui est disponible pour Windows et Linux.

## 2.2 Maximums et minimums

Les premiers débroussaillages étant effectués, l'analyste continuera la recherche d'une compromission par le calcul de métadonnées sur certaines valeurs, telles que les domaines/IPs, volumétries, etc. pour en étudier les variations sur les valeurs extrêmes (maximums et minimums) et mettre en évidence des comportements illégitimes. Ces recherches de maximums et minimums permettront de lever une alerte si les résultats changent dans le temps, et donneront à l'analyste les premières pistes de recherche.

*Machines infectées* Comme les *polling intervals*<sup>4</sup> des malwares sont généralement faibles, inférieurs à 5 minutes, il n'est pas rare de trouver un C&C dans les domaines les plus accédés depuis le réseau. Ainsi, étudier les domaines et IPs les plus accédés, pendant une période d'inactivité comme le week-end, sera intéressant pour détecter des machines infectées. Cette analyse peut être aussi appliquée sur les comptes<sup>5</sup> et machines qui accèdent ou envoient le plus d'informations, à l'extérieur.

*Possible exfiltrations* L'échange d'une grosse quantité d'informations est toujours une action suspecte. Lorsque les informations sont transmises à l'extérieur (volumétrie montante), la suspicion est autrement grande car il pourrait s'agir d'une exfiltrations de données. Ainsi, la recherche du top 10 des domaines/IPs réalisant le plus de volumétrie montante, pendant les weekends ou par mois, sera un bon exercice pour découvrir d'éventuelles exfiltrations. Au même titre, étudier la volumétrie montante

---

3. *Elasticsearch Logstash Kibana Stack*

4. Correspond à l'intervalle de temps entre 2 requêtes, effectués par le malware, pour maintenir le contact avec le C&C.

5. Le terme « compte », dans cet article, fait référence au compte d'authentification utilisé sur les proxy, VPN, etc.

selon les domaines/IPs pour les méthodes HTTP CONNECT et POST sera intéressant. Plus globalement, l'utilisation du ratio :

$$\frac{\text{quantité d'informations transmise}}{\text{quantité d'informations transmise} + \text{quantité d'informations reçue}}$$

servira à l'étude de volumétries dites « asynchrones » (échange de grandes quantités de données émises ou reçues) en analysant les valeurs extrêmes.

Chaque découverte demandera une vérification avec les équipes métiers pour s'assurer que le comportement est bien anormal, et ainsi se familiariser avec les applications et pratiques de l'entreprise (les bonnes, et surtout les mauvaises). Ces études permettent de trouver les grandes étapes et indices de la compromission pour commencer à remonter le fil de l'investigation. Les informations récoltées permettront aussi de déterminer les seuils et informations pertinentes à surveiller durant les prochaines semaines.

### 2.3 Tirer parti des informations extérieures

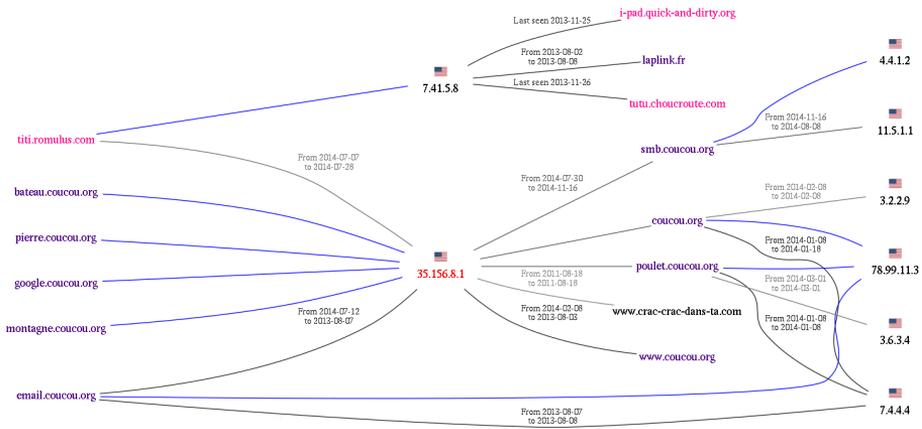
Quand cela est possible, les éléments découverts lors du traitement de l'incident peuvent être recherchés et corrélés sur plusieurs bases d'informations privées ou publiques. Pour des raisons d'image et de marketing, cette opération de collecte et de corrélation est souvent appelée *Cyber Intelligence*. Les objectifs de la *Cyber Int* (pour les intimes) sont d'obtenir des informations nouvelles ou historiques sur un IOC recherché, et confirmer ou infirmer certaines hypothèses, telles que la famille du malware, le groupe d'attaquants, etc. Ces informations récoltées seront ensuite réutilisées pour découvrir de nouveaux éléments de compromission.

Cette section se limitera à une rapide introduction des outils et techniques utilisables pendant les premières semaines d'intervention :

- **évolutions DNS.** Les bases de *passive DNS* [5,8,3,7,6,4] conservent les modifications dans le temps des enregistrements DNS liant nom de domaine et adresses IP. Dans l'exemple figure 1, les arcs marqués d'une date représentent les résolutions DNS effectuées dans le passé entre un domaine et une ou plusieurs adresses IP et les plus foncés et sans date<sup>6</sup>, indiquent l'IP courante du domaine. Sur le terrain, en partant d'un ou plusieurs C&C déjà identifiés dans l'investigation, l'analyste va pouvoir découvrir de nouveaux C&Cs liés à l'attaque, estimer la date de début de l'utilisation de ces C&C ou encore détecter la réutilisation de vieux domaines par le ou les attaquants. Ces nouveaux éléments seront

---

6. En couleur bleue sur la version pdf.



**Figure 1.** Exemple de relations entre C&Cs dans le temps

ensuite ajoutés dans les sondes ou dans l’analyse des logs proxy pour augmenter l’étendue de détection de l’attaque et trouver de nouvelles machines compromises ;

- **enregistrements des domaines.** Les bases Whois stockent des informations concernant l’enregistrement des noms de domaines et offrent la possibilité de rechercher dans ces informations et d’en suivre l’évolution. Dans l’investigation, extraire l’adresse email ou les coordonnées de la *personne* qui a enregistré le C&C, pour les recouper avec des enregistrements d’autres domaines pouvant faire partie de la même attaque ou du même groupe d’attaquants<sup>7</sup> sera utile. Pour illustrer l’intérêt de ce type d’informations, prenons l’exemple du domaine `aeroconf2014.org`, enregistré avec l’email `info@usa.gov.us`. Ce domaine fut utilisé dans l’opération Saffron rose [23,22] comme faux site de la conférence IEEE sur l’aérospatiale. En regardant de plus près l’enregistrement du domaine, on découvre que l’email `info@usa.gov.us` a enregistré 12 domaines, dont :
  - `aps-facebook.com`
  - `home-paypal.com`
  - `linckedln.com`
  - `boeing.com`
  - `www-mail-google.com`

7. Ce service Reverse Whois est très souvent payant bien que les informations y sont parfois partielles. En effet, beaucoup de registraires anonymisent ou ne renseignent qu’une partie des informations Whois.

Ces domaines nous renseignent sur une éventuelle « attaque ciblée » ou « attaque de point d'eau » (*watering hole attack* [40]) utilisant des noms de domaines proches de celui d'entreprises existantes ;

- **géolocalisation des IP.** Plusieurs sites ou services proposent, moyennant finances, de fournir ce type d'informations. Toutefois pour des raisons de performance, et réduire la divulgation d'informations à un tiers il est conseillé de télécharger des bases de géolocalisation comme Maxmind [2] ou Hostip.info [1], et de réaliser le traitement localement.
- **type de domaine.** Déterminer comment les domaines contactés par les malwares sont générés peut renseigner sur le type et le niveau de l'attaquant : *cherche-t-il à éviter la détection ou fait-il le strict minimum ?* Plusieurs techniques peuvent être utilisées par l'attaquant, comme par exemple les *Dynamic DNS* (DDNS), le *Fast Flux* ou encore les *Domain Generation Algorithm* (DGA). Il sera utile de garder à jour une liste [12,13] de l'ensemble des DDNS et ne pas hésiter à étudier, pour chaque domaine de second niveau, l'ensemble des domaines de niveaux suivants. Le calcul de l'entropie de Shannon [48] sera une bonne métrique pour découvrir les domaines générés par des DGA ;
- **domaines connexes.** Étudier les domaines résolus par un système compromis avant et après une requête vers un C&C, par exemple sur une durée de 60 secondes, peut donner à l'analyste d'autres éléments sur le fonctionnement de l'attaque. Ces domaines nommés *Co-occurrences* et *Related Domains* [36] ne sont pas forcément malveillants, mais ils peuvent aider à comprendre le fonctionnement des connexions du malware vers le C&C. En effet, il ne sera pas étonnant de détecter, dans ces domaines connexes, l'enchaînement des différentes étapes de la compromission ou découvrir des requêtes a priori légitimes mais effectuées par le malware pour tester la connectivité à Internet ;
- **bases de recherches privées ou publiques.** De multiples bases publiques de malwares sont consultables sur Internet : Virus Total Intelligence, malwr.com, Totalhash, Clean-mx, Virusshare, ainsi que l'utilisation d'opérateurs dans les moteurs de recherche classiques [24] permettent de collecter de nouveaux éléments. Certaines bases, comme Virus Total offrent, attachés aux analyses, une mine d'informations, notamment les commentaires des utilisateurs. Pour certaines campagnes ou familles de malwares<sup>8</sup>, plusieurs analyses

---

8. Une recherche par *imphash* [25] sera utile pour identifier des malwares similaires.

sont déjà disponibles sur ces bases. Extraire les résultats, et notamment les commentaires, sera un plus dans la collecte d'informations.

### 3 Les semaines suivantes

Les résultats des premières semaines révéleront les grandes étapes et indices de la compromission, comme les C&C, les machines compromises, des patterns réseau ou les malwares utilisés. Avec cette vision plus claire de la compromission, notre analyste perspicace continuera à détecter une éventuelle nouvelle attaque<sup>9</sup>, tout en surveillant la première. Le but ultime est de suivre les actions de l'attaquant<sup>10</sup> pour réagir au bon moment et, quand cela est possible couper une éventuelle exfiltration. Bien sûr, il n'est pas exclu que l'analyste puisse revenir sur les actions effectuées dans les premières semaines, pour ajuster ou récupérer de nouvelles informations.

#### 3.1 Surveillance des incohérences

Surveiller l'attaque, en parallèle de l'investigation, est le subtil équilibre que l'analyste doit trouver durant cette phase. Cette surveillance est généralement orchestrée selon 2 axes :

- surveillance à l'aide des IOCs trouvés dans les logs proxy, captures réseau, analyses de reverse engineering, etc, jusqu'à présent. Le défi sera de compléter et maintenir à jour cette liste d'IOCs pour détecter les nouveaux comportements, par l'analyse des flux réseaux et, dans le passé, au travers des logs.
- analyse des incohérences sur les points sensibles. Étudier les logs et les flux réseaux aux points sensibles du SI, comme par exemple les serveurs de fichiers, messageries, les accès VPN, etc., fournira à l'analyste les clés pour détecter et suivre les actions ou phases de l'attaquant.

Sans préciser l'ensemble des outils ou solutions disponibles, cette sous-section présente des idées et axes de recherche pour découvrir des incohérences sur des points clés du SI :

- **géolocalisation des IP sur les accès VPN.** Vérifier la géolocalisation des IP sources pour chaque compte VPN permettra de détecter rapidement qu'un agent d'une filiale angevine se connecte également tous les soirs depuis Santiago du Chili. Cette technique

---

9. Une attaque peut en cacher une autre.

10. Sans aborder l'ensemble des solutions, on notera l'existence du *framework Open Source ChopShop* [9] pour *décoder* le trafic de plusieurs RAT connus.

- est aussi utilisée par des banques pour détecter des fraudes, suite à des retraits ou paiements dans plusieurs endroits dans le monde ;
- **comptes de service/administrateur utilisés sur des accès distants (RDP, VPN, etc.).** Ces types de comptes sont souvent employées par les attaquants pour prendre le contrôle à distance des postes, serveurs, ou accéder à des ressources. Une revue de l'activité de ces comptes sera intéressante pour détecter d'éventuels comptes compromis et utile pour faire le nettoyage.
  - **échange de fichiers d'archives sur les ports HTTPS.** La détection d'en-têtes de fichiers compressés peut être utile pour découvrir une exfiltration de données et réagir si cela est possible. Attention aux nombreux faux positifs que ce type d'alertes peut générer. Si l'analyste utilise une sonde d'enregistrement des flux réseaux et si l'en-tête n'est pas chiffré, un traitement sur les fichiers pcap permettra *a posteriori* l'archive, ainsi que la liste des fichiers contenus. Comme un grand nombre d'équipements réseau n'analysent pas (ou peu) les flux sur les ports HTTPS, les malwares utilisent fréquemment les ports HTTPS, avec un chiffrement maison ou simplement en clair. Une sonde d'analyse de trafic réseau sera un grand « plus » pour l'analyste qui pourra suivre les flux réseaux et être alerté. Malheureusement, peu d'équipements d'analyse de flux réseaux offrent la flexibilité et le contrôle nécessaire pour répondre aux besoins d'une réponse sur incident. Ainsi, il n'est pas rare que notre analyste dégourdi réalise sa propre sonde basée sur plusieurs briques existantes ou maison, s'appuyant sur des NIDS/NIPS, comme Snort [19], Suricata [20] ou Bro [18] , des FPC [17] ou des outils classiques, comme tcpdump, ngrep, faup [10], etc. ;
  - **métadonnées des protocoles d'échange de fichiers (SMB<sup>11</sup>, FTP, etc.).** Les captures de ce type de trafic sur le LAN et/ou la DMZ sont souvent très volumineuses, mais offrent la possibilité de détecter les comportements suspects suivants :
    - **quantité de données et nombre de fichiers accédés.** La lecture ou l'accès à un grand nombre de fichiers sur les partages SMB permettent de détecter une tentative de collecte des fichiers avant exfiltration. L'utilisation de seuils selon les types de comptes, l'heure de l'action, et par machine seront des idées intéressantes pour détecter ce type de comportements ;

---

11. Le terme « SMB » dans cet article, fait référence aux différentes versions du protocole SMB, ainsi qu'au protocole CIFS.

- **comptes et machines utilisés.** Accéder à, et lister, tous les fichiers présents sur les partages SMB avec un seul compte ou depuis une seule machine sera intéressant à étudier pour découvrir les phases de *scan* des partages SMB et de collecte des fichiers durant une attaque.

L'utilisation des *dissectors* SMB [51], SMB2 [52] et ntlmssp [50] de Tshark<sup>12</sup> sera utile pour bricoler sur les études des volumétries et d'accès aux fichiers des partages SMB. Une description plus détaillée, ainsi que des exemples pratiques sur l'analyse des flux SMB, est disponible dans l'annexe A.

Les analystes les plus combattifs n'hésiteront pas à aller plus loin dans l'étude pour :

- extraire certains ou la totalité des fichiers échangés. L'extraction des binaires, DLL ou encore scripts (VBS, VBA, etc.) circulant sur le réseau sera une première étape pour alimenter automatiquement une solution de *Sandboxing* et faire rapidement le tri. Une approche plus simple passera par le calcul des *hashs* pour les comparer avec les bases d'IOCs ;
- inspecter les commandes des tâches planifiées avec le filtre `atsvc.JobInfo.command`<sup>13</sup> ;
- être alerté de l'accès à un dossier, partage ou fichier sensible. Cela s'avère utile pendant des crises, pour s'assurer que certains projets ne soient pas récupérés ;
- ou encore rechercher des mots ou commandes clés dans le flux SMB, comme par exemple `net use`, `at`, `Command completed successfully`, des noms de fichiers utilisés par l'attaquant, ou encore l'activité d'un compte compromis durant l'attaque.

### 3.2 Comportements inhabituels ou automatisés

En parallèle de la surveillance, le travail d'investigation de notre chercheur forensic est toujours d'actualité. La course pour découvrir un maximum d'indices, détails et preuves continue, et la découverte d'une éventuelle nouvelle attaque est toujours possible. Ainsi, en plus du travail des premières semaines, il continuera à affiner ses recherches de compromission au travers d'études comme celles-ci :

---

12. Ne pas oublier d'utiliser *seccomp* quand on agit avec ce type de programmes et *dissectors*, connus pour de nombreux bugs et vulnérabilités.

13. Il est conseillé d'utiliser le filtre `atsvc.opnum eq 0 and not atsvc.job_id eq 0 and atsvc.status eq 0x00000000` pour analyser uniquement les tâches planifiées valides et correctement créées.

- **navigation automatisée.** La détection d'une navigation Web scriptée n'est pas une facile, mais il existe plusieurs idées pour découvrir rapidement des connexions suspectes. Par exemple, l'absence de téléchargement des scripts JavaScript, CSS ou encore du fichier *favicon.ico* permet de déceler des connexions automatiques, possiblement issues de malwares ou de scripts malveillants ;
- **common name incorrect.** Plusieurs architectures de serveurs C&C ou malwares utilisent des certificats avec un CN invalide, comme par exemple lorsque le CN ne correspond pas aux domaines demandés ou acceptant l'ensemble des domaines possibles CN=\*. ;
- **contenu non chiffré sur les ports HTTPS.** Bien que cette métrique puisse générer des faux positifs, il n'est pas rare de trouver des malwares communiquant en clair sur les ports véhiculant habituellement des flux chiffrés, comme le port 443. De manière plus rare, on trouve aussi des malwares utilisant une connexion SSL avec le NULL cipher ;
- informations non standards contenues dans les champs de protocoles réseau. Certains malwares détournent les champs des protocoles pour leurs usages. Les valeurs utilisées peuvent être statiques ou bien dynamiques. Par exemple, beaucoup de malwares utilisent un *user agent* spécifique, souvent configuré en dur dans le code. Plus rarement, le champ *user agent* sera employé comme canal de contrôle par le malware<sup>14</sup>. Dans ce contexte, le calcul d'une simple distribution de cette valeur selon les domaines contactés permettra de découvrir ce type de méthode de communication ;
- **méthode HTTP non classique.** L'analyse de certaines méthodes HTTP inhabituelles peut fournir des pistes. Par exemple, l'utilisation de la méthode PUT pour exfiltrer une archive vers un serveur Web distant. Ainsi, notre analyste passionné se plongera dans l'analyse des méthodes HTTP (PUT, CONNECT, POST, GET, etc.) selon les domaines/IPs, la volumétrie et le temps (week-end, mois, etc.) pour découvrir un ou plusieurs canaux de contrôle et d'exfiltration.

## 4 Les dernières semaines

Après quelques semaines passées sur l'incident, notre analyste aura trouvé la ou les attaques, compris l'étendue de la compromission, surveillé,

---

14. Malheureusement, certains firewall ou proxy n'enregistrent pas le champ *user agent* par défaut.

et dans certains cas sonn  l'alarme pendant les phases critiques, telles que la r cup ration massive de documents ou les tentatives d'exfiltrations.

En regard du manque de temps et du *coup de feu* des premi res semaines, la fin de l'investigation laisse place   des recherches plus abstraites et plus longues – d corr l es de l'urgence – au travers d'une approche un peu plus math matique. L'objectif de ces  tudes est d'analyser la compromission sous un angle diff rent, pour s'assurer de n'avoir pas oubli  ou n glig  des pistes dans l'urgence des semaines pr c dentes. Beaucoup d'axes d' tudes sont possibles mais les calculs statistiques et de p riodicit s offrent les meilleurs r sultats sur les logs proxy et captures r seau.

## 4.1 Clustering

L'approche statistique par *regroupement* (*clustering*) permet de fournir les premiers  l ments/outils pour aider   la d couverte de comportements suspicieux. Cette technique consiste   grouper les informations qui partagent des m triques similaires pour former un groupe (*cluster*), l'objectif  tant de les analyser individuellement et conjointement. Le *clustering* rend possible de traiter une grande masse de donn es et d'offrir une visualisation adapt e   la recherche d'anomalies par des  tre humains. Cette visualisation permet de se focaliser non pas sur les nuages de points (comportements classiques), mais sur les points ext rieurs aux nuages de points (comportements anormaux). Cette technique est tr s souvent exploit e dans l'exploration de donn es (*data mining*) et l'apprentissage automatique (*machine learning*).

Un grand nombre de m thodes de calcul et d'algorithmes sont disponibles, comme OPTICS [45], DBSCAN [41], MDS [44], LOF [43]. En pratique il est pr f rable de partir sur les algorithmes DBSCAN [41] et MDS [44] pour les raisons suivantes :

- param tres de calcul limit s et simples. Deux param tres sont   d finir pour l'algorithme DBSCAN [41]. Une explication des param tres est disponible dans l'annexe B ;
- aucune limite dans le nombre de *clusters*<sup>15</sup> ;
- impl mentation simple et disponible en Python [37], Ruby et C.

L'ensemble des analyses de cette partie repose sur le choix de la m trique (*distance function*) utilis e. Dans le cas pr sent, la m trique correspond   une fonction de combinaison des champs, des logs proxy ou des protocoles que l'on souhaite  tudier conjointement. Par exemple, pour  tudier la corr lation de l'activit  web des utilisateurs dans la journ e,

---

15. Contrairement   certains algorithmes, comme *K-means* [42].

une métrique intéressante sera la combinaison entre les domaines/IPs, le *timestamps* et le compte ou la machine utilisé(e). Des exemples et idées de métriques applicables sur les logs proxy vont être présentés ci-dessous.

**URIs variables** L'objectif est de collecter, pour un domaine donné un ensemble de métriques sur les URIs pour déterminer si elles sont suffisamment variables et ressemblent à un site web classique (page web, CSS, JavaScript, etc.). En effet, les URIs générées par les malwares comportent, le plus souvent, beaucoup ou peu de variations. La répartition des variations est souvent définie selon une fonction de répartition, où les URIs malveillantes sont réparties dans les valeurs extrêmes et les sites web classiques autour de la moyenne. Ainsi, pour étudier les variations il est intéressant de déterminer pour chaque domaine de deuxième et troisième niveau les statistiques suivantes, selon le temps et le nombre :

- d'extensions différentes (TLD, gTLD, ccTLD, sTLD) ;
- de profondeurs de chemins différents ;
- de noms de fichiers différents ;
- de clefs, valeurs, (clef, valeur) différentes pour l'ensemble des *path* et *query* [15] de l'URI (par exemple : `foo.php?cle1=valeur1`).

La liste de ces métriques n'est pas exhaustive, mais constitue un moyen rapide de découvrir des domaines ou adresses IP suspect(e)s ;

**Périodes** Comme un grand nombre d'actions légitimes ou malveillantes sont prédictibles (mise à jour des logiciels périodiquement) ou *scénarisables* (consultation des emails tous les matins en arrivant au travail), il est intéressant d'étudier les périodes ou fréquences d'une métrique, dans le but d'analyser les déviations de la métrique. Le calcul des fréquences étant complexe et parfois impossible<sup>16</sup>, l'utilisation d'un simple algorithme d'estimation de période est préconisé. Estimer la période d'une valeur de la métrique, suppose de déterminer :

1. les intervalles de temps entre chaque valeur ;
2. compter combien de fois les intervalles de temps apparaissent dans la recherche ;
3. et prendre l'intervalle qui apparaît le plus comme période estimée.

---

16. En raison de limitations de la transformation de Fourier discrète (DFT) [46]. Par exemple, les fréquences, exprimées en Hz sont très faibles lorsque les périodes sont supérieures à plusieurs jours. Dans ce contexte, l'utilisation de la DFT avec des valeurs très faibles donne des résultats hasardeux avec des bibliothèques de fonctions habituelles.

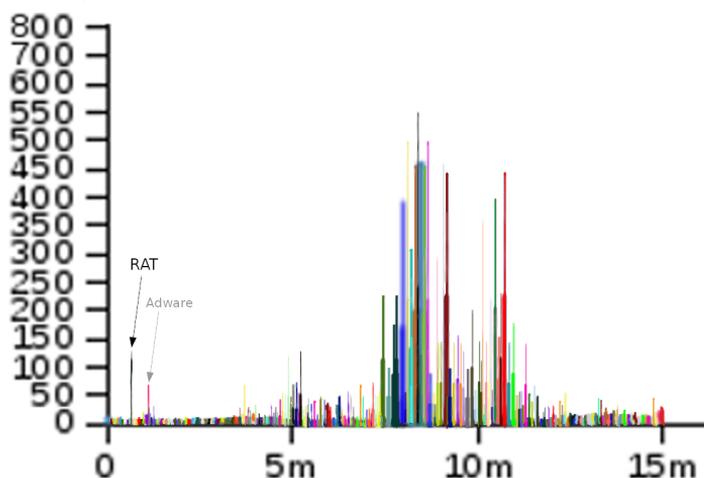
Par exemple sur une semaine, si on observe que le domaine `youtube.com` est accédé tout les jours, à 9h00, 12h00, 14h00 et 16h00. Les intervalles de temps entre les valeurs sont : 2, 3, 4, 5, 7 et 25 heures. Si on compte combien de fois les intervalles de temps apparaissent dans notre exemple, nous obtenons :

- 2 fois pour 2 heures et
- 1 fois pour 3, 4, 5, 7 et 25 heures

Comme l'intervalle 2 heures est celui qui apparaît le plus de fois, notre estimation de période est de 2 heures.

Des exemples et axes de recherches basés sur cette idée sont présentés ci-dessous :

- **Périodes des domaines/IPs.** L'analyse de la périodicité des requêtes est intéressante pour regrouper les domaines/IPs de mêmes types ou offrant les mêmes contenus (flux RSS, site de news ou de publicités, webmails, etc.). Notre chercheur pourra se focaliser sur les domaines ou adresses IP ne répondant pas à ces regroupements et découvrir des domaines ou adresses IP suspect(e)s. Pour cela, l'analyste peut s'aider d'un graphique représentant le nombre de hits sur la période estimée des domaines/IPs. L'étude permet de



**Figure 2.** Nombre de requêtes pour chaque domaine/IP selon leurs périodes estimées (1 semaine de logs proxy)

distinguer les domaines et adresses IPs liés<sup>17</sup> aux :

- sites de microblogage, messagerie instantanée, chats et flux RSS, entre 1 et 10 minutes ;
- sites d'actualités/news, proche des 10 minutes ;
- annonces publicitaires, autour de 10 minutes. Comme les annonces publicitaires sont très souvent intégrées dans les sites d'actualités, ces contenus se trouvent « synchronisés » ;
- webmail, entre 6 et 12 minutes ;
- sites de mises à jours d'antivirus et logiciels, autour d'une heure ;
- sauvegardes, toutes les 24 heures environ ;

Avec ce rapide classement en tête, l'analyste pourra se lancer dans l'étude des valeurs importantes dans les zones non conventionnelles, comme par exemple en dessous de 3 minutes. En regardant plus attentivement, deux valeurs importantes inférieures à 1 minute ressortent de la figure 2. La première, flèche noire à gauche, représente l'activité d'un RAT étalée sur une semaine. La périodicité de cette valeur s'explique par un faible *polling interval* dans la configuration du RAT (50 secondes). La seconde valeur, flèche grise, est due à la présence d'un *adware* qui essaye de joindre un site pour y télécharger les publicités à afficher. Pour continuer l'étude, il est envisageable d'attribuer une couleur à tous les domaines/IPs d'une même catégorie et d'analyser les domaines et IPs ne répondant pas aux catégories prédéfinies. Pour finir, il sera aussi intéressant de comparer les domaines/IPs catégorisés par nos soins avec des listes ou sites web dédiés à la catégorisation ;

- **fréquence et volume d'envois selon un couple d'email (From et To).** La répartition des fréquences et volumes d'envois sont des métriques qui permettent de trouver des canaux de contrôle utilisant les emails ou des boîtes mails compromises utilisées pour envoyer du spear phishing ;
- **temps de réponse d'un email pour une adresse extérieure.** Les temps de réponse fixes indiquent des comportements automatiques (robots, emails d'absence, emails d'erreur, etc.). Les canaux de contrôle utilisant les emails peuvent être détectés en analysant par exemple des temps de réponse courts et fixes.

Ces recherches plus abstraite et plus mathématique, restent encore expérimentales, difficiles et longues à calculer. Cependant, elles offrent les premiers pas dans l'expérimentation du *data mining forensic* pour

---

17. L'utilisation massive de backlist et sites de catégorisation est indispensable pour arriver à regrouper les domaines et IP selon les catégories ci-dessous.

l'analyste imaginatif, qui n'hésitera pas à investir dans plusieurs *clusters* Hadoop pour subvenir à ses besoins.

## 5 Conclusion

Les semaines passées sur l'incident auront permis de trouver les méthodes, outils et actions de l'attaquant. Le plan de reconstruction sonne la fin de l'investigation et le repos pour l'analyste. En espérant que toutes les mesures de la bascule, boutent définitivement l'attaquant hors du SI, le chercheur forensic épuisé s'assurera de faire le point avec l'ensemble des équipes. Malheureusement, le *Retour d'Expérience* (RETEX) est trop souvent rapidement expédié ou même complètement oublié, alors qu'il reste fondamental une fois l'investigation finie<sup>18</sup>. On s'assurera que tout le monde, interne et externe, puisse s'exprimer librement pour échanger et proposer des idées d'améliorations pour les prochaines réponses sur incident.

Depuis plusieurs années, l'analyse d'un système à chaud (*live forensic*) a fait son apparition. Cette nouvelle vision du forensic n'échappe pas aux analyses des logs et flux réseaux : alors qu'une analyse rapide de quelques centaines de Mo de logs proxy était suffisante il y a quelques années, l'analyste doit aujourd'hui trouver une aiguille dans des teraoctets de logs proxy et firewall, tout en jonglant avec les différents formats de flux réseaux et sous la pression de ses supérieurs qui attendent une analyse complète et sans faute, « pour hier ». Cette évolution montre bien le besoin de professionnalisation des outils et des analystes qui doivent arriver déjà préparés et outillés pour mener à bien ces nouvelles missions. Cette préparation doit également passer par des discussions dans la communauté sécurité pour échanger sur les méthodes d'analyse et partager des outils. Nous espérons que cet article aura donné quelques pistes aux lecteurs confrontés à ces problématiques.

## Annexes

### A Détails et exemples sur l'analyse des flux SMB avec Tshark

Comme vu à la section 2.3, l'analyse des flux SMB peut être très utile pour détecter les collectes des fichiers avant exfiltration ou encore découvrir les phases de *scan* des partages SMB et de collecte des fichiers durant une

---

18. Il n'est jamais trop tard pour faire un RETEX.

attaque. Ainsi, cette annexe détaillera les différents *fields* à utiliser pour manipuler les flux SMB avec Tshark et présentera des exemples utiles sur le terrain d'une réponse sur incident.

### A.1 *Dissectors et fields* SMB

Les *dissectors* SMB [51], SMB2 [52] et ntlmssp [50] de Tshark nous permettent d'extraire l'ensemble des informations nécessaires pour réaliser l'étude des flux SMB. Voici les *fields* les plus utiles dans nos recherches :

- smb.cmd. Indique le type de commandes SMB effectuées [34]. Dans notre contexte, les commandes intéressantes sont :
  - SMB\_COM\_WRITE\_ANDX (0x2F) [33] pour obtenir les écritures ;
  - SMB\_COM\_MOVE (0x2A) [29] indique le déplacement de fichiers ;
  - SMB\_COM\_READ\_ANDX (0x2E) [31] renseigne sur les opérations de lecture ;
  - SMB\_COM\_NT\_CREATE\_ANDX (0xA2) [30] est utilisé pour créer et ouvrir un nouveau fichier, ouvrir un nouveau fichier, ouvrir un fichier existant, ouvrir et tronquer un fichier existant à 0 et créer un répertoire, une connexion ou encore un *named pipe* ;
  - SMB\_COM\_TREE\_CONNECT\_ANDX (0x75) [32] indique l'établissement d'un montage réseau ;
  - SMB\_COM\_DELETE\_DIRECTORY (0x01) [28] et SMB\_COM\_DELETE (0x06) [27] renseignent sur la suppression d'un fichier ou d'un répertoire ;
- smb.account, smb.primary\_domain, ntlmssp.auth.username, ntlmssp.auth.domain et ntlmssp.auth.hostname spécifient le compte, domaine et poste utilisés ;
- smb.file et smb.path indiquent le fichier et le *path* manipulés ;

### A.2 Exemples pratiques d'utilisation

La complexité du protocole SMB et les nombreux faux positifs engendrés durant l'analyse, peuvent rapidement dérouter un analyste. Néanmoins, sur le terrain, des exemples simples peuvent être exploités. Par exemple, il sera envisageable de :

- surveiller les énumérations de répertoires et de fichiers ;
- détecter la copie d'un fichier.

```

$ tshark -r smbtoriture.cap -E separator="," -Y "smb.find_first2.
  flags.continue and smb.nt_status eq 0x0" -T fields -e smb.path -
  e smb.file -e smb.search_pattern
....
\\192.168.114.129\TEST,,\testsfileinfo\*
\\192.168.114.129\TEST,,\torture_search.txt
\\192.168.114.129\TEST,t599-599.txt,
\\192.168.114.129\TEST,t699-699.txt,
\\192.168.114.129\TEST,,\torture_search-NOTEXIST.txt
\\192.168.114.129\TEST,,\torture_search.txt
\\192.168.114.129\TEST,,\testsearch\*.*
\\192.168.114.129\TEST,t099-99.txt,
....

```

**Listing 1.** Répertoires et fichiers accédés

Attention, au fichier gpt.ini (stratégie de groupe) qui peut générer un grand nombre de faux positifs, avec le filtre utilisé ci-dessus.

**Copie d'un fichier** Le filtre ci-dessous affiche l'ensemble des fichiers copiés.

```

$ tshark -r smbtoriture.cap -E separator="," -Y "smb.trans2.cmd eq 0
  x0008 and smb.nt_status eq 0x0" -T fields -e smb.path -e smb.
  file -e smb.alloc_size | sort -u
....
\\192.168.114.129\TEST,\mkdirtest\mkdir.dir,
\\192.168.114.129\TEST,\rawchkpath\nt\VB98\vb6.exe,
\\192.168.114.129\TEST,\rawioct1\test.dat,
\\192.168.114.129\TEST,\rawopen\torture_chained.txt,
\\192.168.114.129\TEST,\rawopen\torture_open.txt,
\\192.168.114.129\TEST,\rawopen\torture_openx.exe,
\\192.168.114.129\TEST,\rawopen\torture_t2open_yes.txt,
\\192.168.114.129\TEST,\testeas\ea_max.txt,
\\192.168.114.129\TEST,\testsearch\T013-13.txt.3,
....

```

**Listing 2.** Fichiers copiés

Si l'exécutable PsExec n'est pas renommé par l'attaquant, l'exemple précédent listera toutes les exécutions de PsExec. En effet, l'exécution de PsExec entraîne la copie du binaire sur Admin\$ du poste distant. S'il a été renommé, une recherche du smb.file :

psexecsvc avec la commande SMB : SMB\_COM\_NT\_CREATE\_ANDX sera nécessaire pour détecter la présence de PsExec sur le réseau (filtre : smb.cmd eq 0xa2 and smb.nt\_status eq 0x0 and smb.file contains psexecsvc).

## B Détails des paramètres pour l'algorithme DBSCAN

L'utilisation de l'algorithme DBSCAN [41] demande 2 paramètres :

- La distance  $\epsilon$ .
- Le nombre minimum de points « MinPts ». « MinPts » définit le nombre minimum de points, devant se trouver dans le rayon  $\epsilon$  pour qu'ils soient considérés comme un cluster.

Malheureusement, il n'y a pas de valeur type pour « MinPts » et  $\epsilon$  et leurs valeurs changent en fonction de ce que vous recherchez. Plus « MinPts » est petit, plus il y aura de clusters et donc une augmentation des *clusters* de bruit. Pour  $\epsilon$ , si la valeur est trop grande, une grande partie des informations seront vues comme appartenant au même *cluster*. Pour un  $\epsilon$  trop faible, une majorité des informations ne seront pas regroupées en *cluster*. Pour les plus coriace, la valeur  $\epsilon$  peut être déterminée par le biais d'un graphe de K-distance [47]

## Références

1. Base de géolocalisation Hostip.info. <http://www.hostip.info/dl/index.html>.
2. Base de géolocalisation Maxmind. [https://www.maxmind.com/en/geolocation\\_landing](https://www.maxmind.com/en/geolocation_landing).
3. Base de passive DNS de Dnsparse . <https://dnsparse.insec.auckland.ac.nz/dns/query.php>.
4. Base de passive DNS de Exposure. <http://exposure.iseclab.org/>.
5. Base de passive DNS de la société BFK. [https://www.bfk.de/bfk\\_dnslogger.html](https://www.bfk.de/bfk_dnslogger.html).
6. Base de passive DNS de société Farsight Security (DNSDB). <https://www.dnsdb.info/>.
7. Base de passive DNS de société VirusTotal. <https://www.virustotal.com/en/#search>.
8. Base de passive DNS du CERT-EE . `sim.cert.ee:43`.
9. Github du projet ChopShop. <https://github.com/MITRECNDC/chopshop>.
10. Github du projet Faup. <https://github.com/stricaud/faup>.
11. Github du projet MISP. <https://github.com/MISP/MISP>.
12. Liste de DDNS (dnslookup). <http://dnslookup.me/dynamic-dns/>.
13. Liste de DDNS (freedns.afraid.org). <https://freedns.afraid.org/domain/registry/>.
14. Plateforme MISP du CIRCL. <https://www.circl.lu/services/misp-malware-information-sharing-platform/>.
15. RFC 3986. <http://tools.ietf.org/html/rfc3986>.
16. RFC 5070. <https://www.ietf.org/rfc/rfc5070.txt>.

17. Site officiel du FPC Open Source : OpenFPC. <http://www.openfpc.org/>.
18. Site officiel du NIDS Bro. <https://www.bro.org/>.
19. Site officiel du NIDS Snort. <https://www.snort.org/>.
20. Site officiel du NIDS Suricata. <http://suricata-ids.org/>.
21. Site web d'Iocbucket. <https://www.iocbucket.com>.
22. CrowdStrike. Document de CrowdStrike sur l'opération Saffron rose. <http://blog.crowdstrike.com/cat-scratch-fever-crowdstrike-tracks-newly-reported-iranian-actor-flying-kitten/>.
23. Fireeye. Document de Fireeye sur l'opération Saffron rose. <https://www.fireeye.com/resources/pdfs/fireeye-operation-saffron-rose.pdf>.
24. Google. Définition et explication des advanced-operators de Google. <https://sites.google.com/site/gwebsearcheducation/advanced-operators>.
25. Mandiant. Définition et explication du imphash par Mandiant. <https://www.mandiant.com/blog/tracking-malware-import-hashing/>.
26. Mandiant. Site officiel sur le format OpenIOC. <http://www.openioc.org/>.
27. Microsoft. Documentation sur la commande SMB\_COM\_DELETE (0x06). <https://msdn.microsoft.com/en-us/library/ee442133.aspx>.
28. Microsoft. Documentation sur la commande SMB\_COM\_DELETE\_DIRECTORY (0x01). <https://msdn.microsoft.com/en-us/library/ee441479.aspx>.
29. Microsoft. Documentation sur la commande SMB\_COM\_MOVE (0x2A). <https://msdn.microsoft.com/en-us/library/ee441847.aspx>.
30. Microsoft. Documentation sur la commande SMB\_COM\_NT\_CREATE\_ANDX (0xA2). <https://msdn.microsoft.com/en-us/library/ee442091.aspx>.
31. Microsoft. Documentation sur la commande SMB\_COM\_READ\_ANDX (0x2E). <https://msdn.microsoft.com/en-us/library/ee441503.aspx>.
32. Microsoft. Documentation sur la commande SMB\_COM\_TREE\_CONNECT\_ANDX (0x75). <https://msdn.microsoft.com/en-us/library/ee441940.aspx>.
33. Microsoft. Documentation sur la commande SMB\_COM\_WRITE\_ANDX (0x2F). <https://msdn.microsoft.com/en-us/library/ee441848.aspx>.
34. Microsoft. Documentation sur les commandes SMB. <https://msdn.microsoft.com/en-us/library/ee441741.aspx>.
35. MITRE. Site officiel sur le format STIX. <http://stix.mitre.org/>.
36. OpenDNS. Définition et explication des termes co-occurrences et Related Domains. <http://labs.opendns.com/2013/07/24/co-occurrences/>.
37. Scikit-learn. Site officiel de la bibliothèque python Scikit learn. <http://scikit-learn.org/stable/index.html>.
38. Trendmicro. Définition du terme Command and Control. <https://www.trendmicro.com/vinfo/us/security/definition/command-and-control-%28c-c%29-server>.
39. Wikipedia. Définition du terme IOC (Indicator of compromise). [https://en.wikipedia.org/wiki/Indicator\\_of\\_compromise](https://en.wikipedia.org/wiki/Indicator_of_compromise).
40. Wikipedia. Définition du terme watering hole attack (attaque de point d'eau). [https://en.wikipedia.org/wiki/Watering\\_Hole](https://en.wikipedia.org/wiki/Watering_Hole).

41. Wikipedia. Explication de l'algorithme DBSCAN. <https://en.wikipedia.org/wiki/DBSCAN>.
42. Wikipedia. Explication de l'algorithme K-means. [https://en.wikipedia.org/wiki/K-means\\_clustering](https://en.wikipedia.org/wiki/K-means_clustering).
43. Wikipedia. Explication de l'algorithme LOF. [https://en.wikipedia.org/wiki/Local\\_outlier\\_factor](https://en.wikipedia.org/wiki/Local_outlier_factor).
44. Wikipedia. Explication de l'algorithme MDS (Multidimensional scaling). [https://en.wikipedia.org/wiki/Multidimensional\\_scaling](https://en.wikipedia.org/wiki/Multidimensional_scaling).
45. Wikipedia. Explication de l'algorithme OPTICS. <https://en.wikipedia.org/wiki/OPTICS>.
46. Wikipedia. Explication de transformation de Fourier discrète. [https://fr.wikipedia.org/wiki/Transformation\\_de\\_Fourier\\_discr%C3%A8te](https://fr.wikipedia.org/wiki/Transformation_de_Fourier_discr%C3%A8te).
47. Wikipedia. Explication du graphe de K-distance. [https://en.wikipedia.org/wiki/K-distance\\_graph](https://en.wikipedia.org/wiki/K-distance_graph).
48. Wikipedia. Formule de l'entropie de Shannon. [https://en.wikipedia.org/wiki/Information\\_theory#Quantities\\_of\\_information](https://en.wikipedia.org/wiki/Information_theory#Quantities_of_information).
49. Wiktionary. Définition du terme low-hanging fruit. [https://en.wiktionary.org/wiki/low-hanging\\_fruit](https://en.wiktionary.org/wiki/low-hanging_fruit).
50. Wireshark. Documentation sur le dissector ntlmssp. <https://www.wireshark.org/docs/dfref/n/ntlmssp.html>.
51. Wireshark. Documentation sur le dissector SMB. <https://www.wireshark.org/docs/dfref/s/smb.html>.
52. Wireshark. Documentation sur le dissector SMB2. <https://www.wireshark.org/docs/dfref/s/smb2.html>.