
RSA - Réseau Avancé

Projet MyAdBlock

RUCHOT Guillaume
MOLLARD Romaric

Avril 2017

Sommaire

1	Introduction	1
1.1	À propos de ce document	1
1.2	Remerciements	1
2	Partie I - Étude du fonctionnement	2
2.1	Observation du dialogue navigateur/serveur	2
2.1.1	Observation avec wireshark	2
2.1.2	Observation avec wireshark après mise en place du proxy	3
2.1.3	Observation avec un serveur TCP de base	3
2.2	Spécification du proxy MyAdBlock	3
2.2.1	Création d'un proxy HTTP transparent	3
2.2.2	Prise en compte des connexions HTTPS	4
2.2.3	Filtrage des URL	4
2.2.4	Résumé	4
2.3	Recherche de masques et filtres pour la publicité	4
3	Partie II - Création de l'outil	6
3.1	Fonctionnement général du programme	6
3.2	Difficultés rencontrées	6
4	Conclusion	7
4.1	Conclusion	7
4.2	Méthodes et outils	7
4.3	Répartition du temps	7

1 Introduction

1.1 À propos de ce document

Ce document a pour but de présenter le projet *MyAdBlock* et son élaboration. Le projet *MyAdBlock* consiste en la réalisation d'un proxy HTTP/HTTPS capable de bloquer les accès à certains sites en utilisant une liste spécifique de masques, dans le but par exemple de bloquer les publicités. Ce projet a été réalisé dans le cadre de l'étude de l'utilisation avancée du système et du réseau du niveau application au niveau transport. Le langage utilisé est le C, pour sa proximité avec le système.

Créer un serveur de ce type permet de comprendre pleinement le fonctionnement TCP pour le transport du protocole HTTP : les informations contenues dans les entêtes HTTP, la transformation des noms de serveurs en adresses, l'obligation pour le serveur d'avoir une gestion simultanée des clients...

1.2 Remerciements

Nous avons réalisé ce projet à l'aide des sources suivantes :

Documentation linux (ici `getaddrinfo(3)`)

<http://man7.org/linux/man-pages/man3/getaddrinfo.3.html>

Wireshark

<https://www.wireshark.org/>

Liste de filtres publicitaires

<http://easylist.to/>

Conversion nom de domaine vers adresse ip

<http://get-site-ip.com/>

2 Partie I - Étude du fonctionnement

2.1 Observation du dialogue navigateur/serveur

Nous utiliserons pour la suite deux navigateurs, Firefox qui sera configuré pour fonctionner avec un proxy local sur le port 80, et Chrome qui lui permettra de comparer les résultats.

2.1.1 Observation avec wireshark

Nous avons commencé par observer les transmissions TCP entre Chrome et www.telecomnancy.eu avec Wireshark. Comme www.telecomnancy.eu renvoie une redirection HTTP 301, nous avons filtré les résultats sur l'adresse ip 193.50.135.38, qui correspond à l'adresse ip du serveur telecomnancy.univ-lorraine.fr.

Ainsi le filtre WireShark devient celui-ci : `ip.dst_host == 193.50.135.38 or ip.src_host == 193.50.135.38`

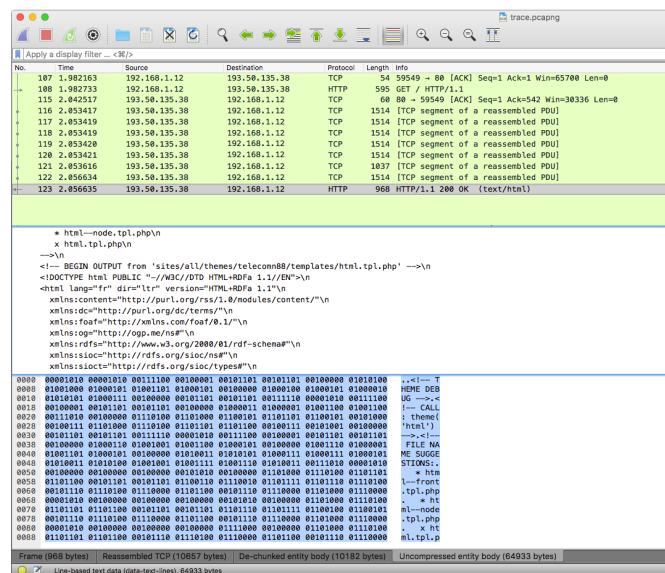


Figure 1 : Trace Wireshark contenant la communication vers telecomnancy.univ-lorraine.fr

Pour observer le contenu HTTP d'une transmission, il faut regarder le tout dernier élément de la transmission (ici l'élément 123, HTTP 200 OK, figure 1), et demander à Wireshark de décompresser le contenu. Le document html est compressé, et ainsi on ne peut pas travailler dessus tant qu'il n'est pas complet. Cependant, le header HTTP n'est jamais compressé, ni pour l'envoi ni pour la réponse, et ce pour permettre au navigateur d'obtenir des informations comme la taille du contenu en cours de réception (content-length).

Nous remarquons que la communication commence par l'envoi d'un paquet TCP contenant la requête HTTP "GET" (élément 108, figure 1, les lignes précédentes concernent la mise en place de la connexion TCP), suivi du contenu de la page envoyé en retour. Le tout se fait via TCP. En regardant les paquets TCP en détail, nous remarquons que le dernier paquet envoyé possède le flag FIN, et c'est le seul élément nous permettant à priori d'obtenir l'information de fin de transmission.

Nous pouvons noter que le protocole HTTP ne contient pas l'ip du serveur distant et seulement le nom du serveur.

2.1.2 Observation avec wireshark après mise en place du proxy

Nous avons configuré Firefox pour se connecter à un serveur proxy local qui n'est qu'un simple serveur TCP basique (qui accepte les connexions) afin d'observer le comportement de Firefox et surtout les informations envoyées au futur proxy.

Dans ce cas nous devons observer les transactions locales dans le mode "Loopback Io0".

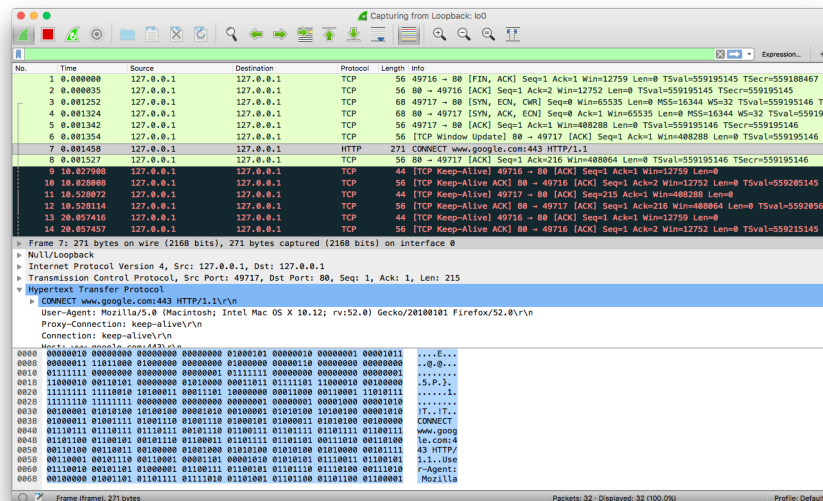


Figure 2 : Trace Wireshark contenant la communication vers telecomnancy.univ-lorraine.fr avec proxy

Nous observons ici qu'après la connexion TCP, Firefox envoie le header HTTP "CONNECT" (qui est un équivalent de GET dans le cas du https) au proxy local (ligne 7 de la figure 2). En observant le contenu de cette entête nous confirmons que nous n'avons aucun moyen de connaître directement par lecture l'adresse ip du serveur web demandé. Nous devrons donc utiliser des outils comme gethostbyname(1) ou getaddrinfo(3).

2.1.3 Observation avec un serveur TCP de base

En mettant en place un serveur TCP très basique étudié en cours, nous avons remarqué que les navigateurs ouvraient une nouvelle connexion au serveur pour chaque url voulue. La question de la gestion simultanée de plusieurs clients est donc primordiale si nous voulons qu'une page web ne charge pas chaque composant (js, css, iframes...) un par un.

2.2 Spécification du proxy MyAdBlock

Après les observations effectuées sur WireShark, nous avons mis en place le fonctionnement général de notre programme.

Nous avons séparé la réalisation du projet en plusieurs étapes.

2.2.1 Création d'un proxy HTTP transparent

Dans un premier temps, une grande partie du travail était la réalisation d'un proxy permettant l'accès à internet. Si ceci semble simple en apparence, c'est la partie la plus importante du travail dans ce projet. Internet évolue très rapidement et si nous développons un proxy limité à l'IPv4 et le protocole HTTP, nous ne pourrions pas afficher grand chose. Nous avons donc décidé de mettre en place notre proxy transparent en commençant de manière simple sur cette url :

<http://www.example.com/>

Son contenu est très court (tient en un seul paquet en temps normal), l'accès se fait via HTTP (pas de HTTPS) et enfin il n'y a pas de liens internes, c'est à dire que le navigateur ne va pas chercher d'autres liens pour charger des images par exemple ou des feuilles de style css.

Afin d'étendre notre proxy à la réception de plusieurs paquets réponse et la connexion de clients multiples, nous avons utilisé un autre site web simple :

<http://cheval.fr/>

Celui-ci contient une image, qui ne tient pas en un seul paquet, ce qui permet de vérifier un autre point du proxy.

2.2.2 Prise en compte des connexions HTTPS

Il est presque impossible d'utiliser internet sans le protocole HTTPS, car même les sites web HTTP utilisent des publicités externes ou bien des feuilles de styles générales (bootstrap) accessibles seulement via HTTPS.

Heureusement détecter les transactions SSL est simple puisqu'il suffit de détecter la méthode "CONNECT".

Nous testerons cette fonctionnalité simplement sur google et sur

<https://www.leboncoin.fr>

2.2.3 Filtrage des URL

Nous effectuerons le filtrage des accès sur le second argument des entêtes HTTP, soit le chemin demandé. Il contient à quelques détails près l'url visible dans le navigateur.

Le programme générera une LinkedList contenant l'intégralité des masques à utiliser. Lorsqu'une nouvelle url sera demandée, si l'un de ces masques est entièrement contenu dans cette url alors l'url est refusée.

2.2.4 Résumé

Voici un résumé graphique du fonctionnement du proxy MyAdBlock.

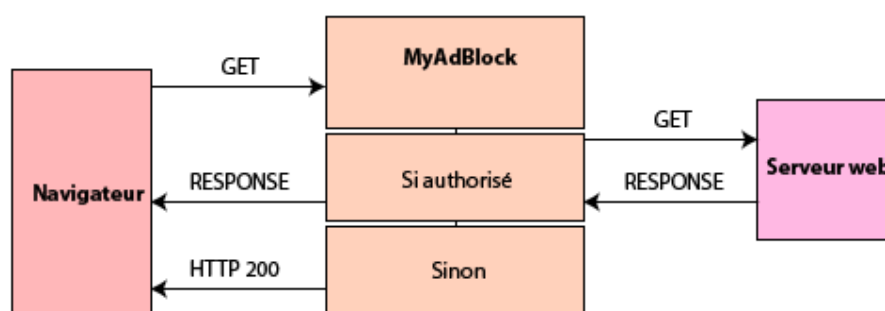


Figure 3 : Fonctionnement de MyAdBlock

2.3 Recherche de masques et filtres pour la publicité

Pour tester une des utilisations du proxy qui est le filtrage des publicité, nous avons utilisé les listes présentes sur le site :

<http://easylist.to/>

Une bonne partie de cette liste concerne des éléments du contenu HTML d'un site web, nous ne

pouvons pas utiliser ces éléments car nous n'avons à aucun moment accès au contenu brut des pages transférées.

3 Partie II - Création de l'outil

3.1 Fonctionnement général du programme

3.2 Difficultés rencontrées

4 Conclusion

4.1 Conclusion

4.2 Méthodes et outils

Programmation :

Nous utilisons plusieurs systèmes et logiciels.

La gestion des version s'est faite via un repository github et donc en utilisant l'outil git.

L'écriture de ce rapport a été effectuée sous le langage LaTeX.

Les systèmes d'exploitations utilisés ont été Mac OS X, Ubuntu et Windows 10 (sous-système Windows pour Linux).

Les éditeurs utilisés pour la programmation ont été Atom et Visual Studio Code.

Correction des erreurs :

Nous avons utilisé le compilateur GCC pour le développement, Valgrind pour la vérification des fuites mémoires et Valgrind + Helgrind pour la vérification des erreurs de multi-threading.

4.3 Répartition du temps

	Romaric MOLLARD	Guillaume RUCHOT
Partie 1 - Lecture basique	0h	3h
Partie 2 - Lecture des options	2h	0h
Partie 3 - Extraction	0h	4h
Partie 4 - Lecture et extraction avancées	5h	1h
Partie 5 - Ajout des threads	5h	4h
Partie 6 - Archives compressées	4h	0h
Partie 7 - Checksum	1h	0h
Corrections du code	10h	8h
Rédaction du rapport	4h	2h
Total	26h	22h