

Sniff'eirb, sniffer intelligent

Rapport de projet

Nicolas Verdier Nicolas Retrain Gabrielle Schmidt
Benjamin Vandenberghe

Responsable : M. Aymeric Vincent

27 janvier 2013

Table des matières

Introduction	2
1 Les besoins et l'état de l'art	3
1.1 L'état de l'art	3
1.2 Nos objectifs	3
1.2.1 Premières étapes	3
1.2.2 La valeur ajoutée de Sniff'eirb	4
1.2.3 Les idées abandonnées ou non implémentées	4
2 L'architecture	5
2.1 L'architecture générale	5
2.2 Explication du rôle des différentes couches	6
2.2.1 Client	6
2.2.2 Serveur web	6
2.2.3 Analyse et reconstruction	6
2.2.4 Base de données	7
2.2.5 Acquisition des données	7
2.2.6 Paquets	7
3 L'implémentation	8
3.1 La base du sniffer : la récupération des paquets et leur stockage	8
3.2 La reconstruction	9
3.2.1 Reconstruction d'un flux	9
3.2.2 Reconstruction d'un document	10
3.3 L'affichage	10
Conclusion	11
Annexes	12

Introduction

Dans le cadre du projet de développement système et réseaux, nous avons le choix sur notre sujet. En parallèle de ce projet, nous travaillions sur un projet de sécurité dont le but était la réalisation d'un Sniffer intelligent¹. L'idée nous est alors venue de les fusionner afin d'approfondir d'avantage le sujet.

En effet, l'élaboration d'un Sniffer faisait partie des possibilités de sujet du projet de développement. Pour confirmer (ou infirmer) l'efficacité de ce Sniffer, nous avons dû réaliser plusieurs tests d'analyses de réseaux. Néanmoins, ces analyses ont uniquement été faites dans un but pédagogique².

Ce rapport décrit donc dans un premier temps les besoins et l'existant, puis l'architecture envisagée de notre Sniffer, et enfin son implémentation effectivement réalisée.

1. La description des besoins et fonctionnalités, faite plus loin dans le rapport, définit cette notion pour le moment ambiguë d'*intelligent*.

2. Nous certifions sur notre honneur ne jamais avoir eu d'idées malveillantes en testant notre Sniffer :).

Chapitre 1

Les besoins et l'état de l'art

1.1 L'état de l'art

Il existe un grand nombre de logiciels et de bibliothèques qui proposent des fonctionnalités de captures réseaux de bas niveau et d'analyse de paquets. Nous pouvons notamment citer Libpcap, Scapy, Winpcap ou encore Tcpdump. Ces outils ont été des modèles d'inspiration pour réaliser notre propre *Sniffer*¹ réseau.

1.2 Nos objectifs

1.2.1 Premières étapes

L'objectif premier était tout d'abord de retrouver les fonctionnalités principales de Wireshark, comme par exemple afficher une liste des paquets circulant sur le réseaux ainsi que les informations s'y rapportant (timestamp, adresse IP source, adresse IP destination, port, protocole...), tout en apportant une solution différente et plus simple d'utilisation à un utilisateur néophyte. Ainsi nous avons défini trois premiers objectifs :

- Récupérer les communications² en clair d'un réseau.
- Stocker les paquets de cette communication.
- Afficher les différents flux³ possibles de cette communication.

1. Terme anglicisé désignant un dispositif permettant d'analyser le trafic d'un réseau.

2. Une communication est définie par un échange réseau entre deux paires (Adresse IP ; Port).

3. Un flux est défini comme un assemblage possible des paquets d'une communication.

1.2.2 La valeur ajoutée de Sniff'eirb

Le but de ce projet n'étant pas de reproduire un Wireshark, nous devons y apporter des atouts significatifs. Plusieurs fonctionnalités ont donc été implémentées :

L'arbre des possibilités pour les flux TCP : Lors de la capture des paquets circulant sur le réseau pour une même communication, des événements inattendus, comme la répétition de paquets ou encore la perte de certains, pouvaient survenir. La capture de tous les paquets permet donc de connaître les différentes façons de les réassembler pour obtenir un flux décrivant la communication.

La gestion du protocole HTTP : La reconstruction des pages webs visitées par un utilisateur sur le réseau était le fil conducteur du projet. L'analyse des paquets circulant sur le réseau était la première étape afin d'arriver à cette fin.

1.2.3 Les idées abandonnées ou non implémentées

Beaucoup d'idées durent être abandonnée du fait d'un manque de temps pour les traiter. Voici une liste non exhaustive de ces idées :

Récupérer les cookies d'un utilisateur du réseau afin de les réinjecter dans le navigateur et pouvoir accéder aux données personnelles de cet utilisateur.

Établir un mode *ghost* qui permettrait au Sniffer d'être quasiment indétectable sur le réseau.

Permettre le choix de l'interface réseau pour pouvoir écouter sur un réseau choisi.

Exporter les données au format *pcap* afin de pouvoir réutiliser ces données avec d'autres analyseurs réseau.

Supporter l'analyse des paquets transmis en UDP pour récupérer toujours plus d'informations qui transitent sur le réseau.

Chapitre 2

L'architecture

2.1 L'architecture générale

Nous avons choisi dès le départ de suivre une architecture de type MVC, et nous nous y sommes tenu tout au long du projet. Afin de nous faciliter le découpage du code et de définir nos principales interfaces, nous avons conçu au démarrage du projet une architecture globale. Elle est représentée dans la figure 2.1.

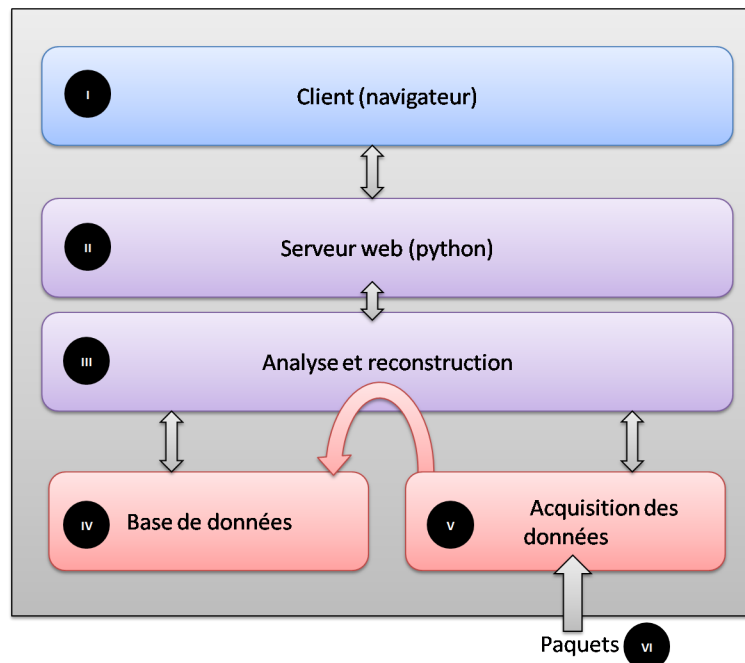


FIGURE 2.1 – Architecture de Sniffeirb

Nous avons choisi d'utiliser le langage Python pour notre projet, en partie pour ses nombreuses bibliothèques, y compris pour les captures réseaux. De plus, nous avons tous envie d'approfondir nos connaissances dans un langage de programmation autre que le C, C++ ou Java.

2.2 Explication du rôle des différentes couches

2.2.1 Client

Le client est soit le terminal dans lequel on lance le programme, soit un navigateur web. Nous n'offrons pas tout à fait les mêmes services sur ces deux interfaces, ceci est lié à notre volonté d'avoir une utilisation simple. De ce fait, le terminal permet seulement de lancer une capture, ou de lancer le programme avec différentes options (lancement du serveur web, choix d'une ancienne capture à charger, choix du nom de la nouvelle capture, chargement d'un fichier pcap..).

L'interface web permet de lancer une capture, d'afficher les flux des différentes communications, de reconstruire les documents et de changer certains paramètres.

Le navigateur utilise des requêtes AJAX afin d'interagir avec le programme. Dans le but de fournir un programme utilisable sur le plus grand nombre de navigateur possible, nous avons utilisé des frameworks reconnus dans le milieu. Nous avons donc choisi d'utiliser JQuery pour nos requêtes, Bootstrap pour la mise en page et le CSS, et DataTables pour la présentation des flux sous forme de tableau.

2.2.2 Serveur web

Le serveur web est l'interface entre les requêtes HTML provenant du client et les différentes actions possibles du sniffer. Après avoir lancé une action, le serveur renvoie si besoin est des données au navigateur au format JSON.

Nous n'avons pas utilisé de serveur déjà fourni dans un framework Python, trop lourd pour l'utilisation que nous en faisons. Nous avons choisi d'utiliser JSON pour plusieurs raisons. C'est un format de données structuré, moins verbeux que le XML, et qui est très bien intégré avec le JavaScript et le Python. De plus le format JSON est identique au format BSON utilisé par la base de données.

2.2.3 Analyse et reconstruction

Ce composant analyse une communication pour déterminer son protocole, reconstruit les flux possibles et reconstruit les documents qui transitent en donnant un pourcentage de validité. Nous assurons seulement la reconstruction des documents sur le protocole HTTP conformément à nos spécifications fonctionnelles.

2.2.4 Base de données

La base de données sert à stocker les paquets réseaux regroupés par communication. Stocker les données nous permet de rejouer des échanges de données (soit pour le débogage, soit suite à une amélioration algorithmique du projet).

Nous avons choisi d'utiliser la base NoSQL Mongo pour ses performances et sa souplesse. Elle s'utilise aisément en Python grâce au module Pymongo. Il faut donc installer mongodb et lancer le service mongod pour pouvoir utiliser le sniffer.

2.2.5 Acquisition des données

Ce composant récupère les données réseaux, ou lit un fichier pcap, et les met dans la base de données. Les paquets sont regroupés par communication. Nous utilisons le module Scapy de Python afin de récupérer les paquets. Il est de plus haut niveau que la bibliothèque Libpcap ou WinPcap, ce qui facilite la capture des paquets réseau.

2.2.6 Paquets

Ceci représente le réseau, et les paquets qui circulent. En paramétrant la carte réseau en mode monitor, l'utilisateur peut alors écouter tout le réseau et analyser les paquets sniffés.

Chapitre 3

L'implémentation

3.1 La base du sniffer : la récupération des paquets et leur stockage

La partie principale de notre projet est de capturer un flux réseau, nous nous sommes donc appuyés sur la bibliothèque Scapy. Cependant, la fonction de capture réseau fournie par Scapy permet de lancer une capture sur un nombre fixe de paquet ou un temps donné. C'est également à la fin de la capture que l'utilisateur peut récupérer les données.

Dans notre projet, nous voulions pouvoir sniffer le réseau un temps indéterminé, jusqu'à une interaction de l'utilisateur, et traiter les données au fur et à mesure. Nous avons donc repris le code de la fonction "sniff" de Scapy et modifié de telle sorte qu'une fonction de contrôle soit appelée régulièrement. Ainsi, nous pouvons couper et lancer le sniffer de façon événementielle, et récupérer les données en temps réel.

En parallèle de la capture, nous stockons les paquets capturés dans une base de données. Celle-ci est composée d'une unique collection. Cette collection comporte plusieurs sessions qui permettent d'identifier les différentes captures réalisées. Ensuite, chaque session est composée de communications, elles-mêmes composées de un ou plusieurs documents dont les champs peuvent varier selon le type de paquets. Un document correspond à une communication¹. Ce qui signifie que l'on ne prend pas en compte les éventuels paquets qui utilisent un autre protocole qu'IP.

Chaque document possède certains champs dans tous les cas. Il s'agit des champs :

- adresses IP source et destination
- ports sources et destination
- protocole (TCP, UDP, ICMP ou "other")
- type (IP)

1. Une communication est définie par un échange réseau entre deux paires (Adresse IP ; Port).

- session (un identifiant de type *sess_date*)
et du document paquets. Les documents paquets étant eux-mêmes composés de tous les champs utiles des paquets, le tout au format BSON.

Dans l'exemple situé en annexe, on peut voir un document de notre collection où l'on retrouve tous les éléments indiqués, classés par ordre alphabétique.

3.2 La reconstruction

Un des principaux objectifs de Sniff'eirb était de rendre l'exploitation des données capturées la plus évidente possible. Nous avons donc décidé de permettre l'affichage des échanges HTTP, notamment celui des pages HTML et des images. Pour ce faire, nous avons commencé par reconstruire les différents flux de données.

3.2.1 Recontruction d'un flux

Comme nous ne pouvions pas implémenter la reconstruction de flux de la même façon pour tous les protocoles de couche 4 (UDP et TCP) et que nous voulions présenter un résultat démonstratif même pour le néophyte, nous nous sommes consacrés au protocole TCP, puis au média HTTP.

Reconnaissance du protocole et du média

Il a donc tout d'abord fallu reconnaître le protocole qui nous intéresse et le média, avant d'en reconstruire le flux. Comme nous l'avons vu dans la partie précédente, le protocole est facilement trouvable directement dans la base de données (puisqu'il s'agit d'un champ du paquet IP). Pour le média, cela a été plus technique.

Au départ, nous l'avons simplement déterminé grâce au port (source ou destination). En effet, traditionnellement, certains ports sont utilisés pour certains médias (par exemple, 80 pour HTTP). Mais en ne nous concentrant que sur cette information, nous aurions manqué toutes les communications d'un média sur d'autres ports (comme 8080 pour HTTP), qui, d'un point de vue de la sécurité, pouvaient être les plus intéressantes.

Nous avons donc décidé de parser la charge utile du paquet TCP à la recherche d'expression régulière trahissant le HTTP(S), comme les en-têtes *HTTP/numero_de_version* et les mots clés *GET* ou *POST* ou des URLs.

La multiplicité des flux possibles

À ce moment-là, nous disposons donc de l'ensemble des paquets appartenant à la même communication ainsi que du média qui correspond. Comme le protocole est TCP, il paraît assez facile de reconstruire le flux : on classe les paquets par numéro de séquence croissant et on concatène les charges utiles.

Cependant, cette démarche n'est pas possible. En effet, il arrive parfois que plusieurs paquets arrivent avec le même numéro de séquence, voire qu'il manque

des paquets au milieu de la communication. C'est pourquoi nous avons dû prévoir de reconstruire tous les flux possibles.

La reconstruction de tous ces flux

Pour ce faire, nous avons créé une structure de liste doublement chaînée que nous avons ensuite parcouru pour en tirer tous les flux possibles. Pour remplir cette chaîne doublement chaînée, baptisée *LianaTree*, nous avons simplement parcouru la liste des paquets d'une communication en commençant par celui qui a le plus petit numéro de séquence et nous cherchions tous ses successeurs.

Le successeur d'un paquet ayant pour numéro de séquence la somme de la taille de sa charge utile et du numéro de séquence de son prédécesseur (modulo 2^{32}).

Cependant, nous avons un problème dans certains cas, car le numéro de séquence TCP est pris au hasard parmi 2^{32} possibilités et il était donc possible d'avoir la fin du flux avant le début (si le numéro de séquence du début de la communication était proche de $2^{32} - 1$). Finalement, nous recherchons simplement les successeurs et prédécesseurs de chaque paquet via le numéro de séquence et la taille de la charge utile.

Une fois cette liste créée, nous devons la parcourir en entier en partant du paquet qui n'a pas de prédécesseur. Si plusieurs paquets n'ont pas de prédécesseur, nous recommençons le parcours autant de fois que nécessaire. Nous réalisons un parcours en largeur de la liste qui découle de ce paquet et nous concaténons les charges utiles au fur et à mesure. À la fin, nous éliminons les doublons et renvoyons le résultat.

3.2.2 Reconstruction d'un document

Afin de reconstruire des documents, nous parsons chaque flux HTTP pour trouver les en-têtes et corps des requêtes. Après cette première étape, nous analysons le corps trouvé par rapport au content-type de l'en-tête. Plusieurs cas sont alors possibles. Si le content-type précise qu'il s'agit d'un document gzippé, nous décompressons le document grâce à la bibliothèque gzip.

Une autre solution est d'envoyer directement au navigateur le document reconstruit, et laissons le navigateur gérer la décompression. Dans le cas où il ne s'agit pas de gzip, nous créons des fichiers html pour tous les fichiers textes, et nous envoyons au navigateur les flux binaires dans le cas d'images.

3.3 L'affichage

Afin d'avoir un retour visuel des fonctionnalités implémentées, nous avons décidé d'afficher les résultats dans une page web interactive. Nous avons donc codé un serveur web minimaliste pour récupérer des requêtes HTTP du navigateur. Le serveur utilise un template HTML pour initialiser la première page et des variables globales javascript. Ceci permet de séparer le code HTML du code python pour respecter le modèle MVC.

Ensuite, après l'initialisation toutes les actions de l'utilisateur sont envoyées en AJAX via le framework JQuery.

Pour la mise en page nous nous sommes appuyés sur le framework Bootstrap, ainsi que sur Datatables pour la création du tableau. Datatables nous a permis par exemple d'intégrer les filtres, les tris, et l'actualisation en temps réel des données. Nous avons dû toutefois modifier quelques fichiers de Bootstrap et de Datatables afin de les faire fonctionner ensemble.

Vous trouverez en annexe des captures d'écran de l'interface graphique de Sniff'eirb.

Conclusion

Ce projet fut très intéressant dans le sens où il nous a permis de mettre en application nos connaissances sur l'écoute et l'analyse de paquets transitant sur un réseau.

Également, nous avons pu aiguïser nos compétences sur le langage qui nous a accompagné tout au long de ce projet : le *Python*. Enfin, malheureusement, le manque de temps nous a empêché de mettre en œuvre toutes les idées qui nous avaient traversé l'esprit et qui auraient non seulement élargi le panel de fonctionnalités de notre Sniffer, mais aussi gonflé davantage nos connaissances.

Vous trouverez le code de notre projet à l'adresse suivante : <https://github.com/bvandenberghes/Sniffeirb>.

Annexes

The screenshot shows the Sniffairb web application running in a Google Chrome browser. The interface includes a navigation bar with links for Launch, Parameters, Network map, Archive, and Help. Below the navigation bar, there is a search bar and a table displaying network traffic records. The table has five columns: Initialisation time, Source, Destination, Protocol, and Size. The records are listed in a table with alternating light blue and white rows.

Initialisation time	Source	Destination	Protocol	Size
11/1/2013 20:18:13.492	173.194.69.16.993	192.168.1.101.53010	TCP - imap sécurisé (ssl..	1861
11/1/2013 20:18:13.492	192.168.1.101.53010	173.194.69.16.993	TCP - imap sécurisé (ssl..	1551
11/1/2013 20:18:45.174	192.168.1.101.48016	173.194.78.125.5222	TCP - serveur Jabber	922
11/1/2013 20:18:45.198	173.194.78.125.5222	192.168.1.101.48016	TCP - serveur Jabber	222
11/1/2013 20:18:52.281	192.168.1.101.35668	173.194.34.7.443	TCP - https	7672
11/1/2013 20:18:52.301	173.194.34.7.443	192.168.1.101.35668	TCP - https	4763
11/1/2013 20:19:3.237	192.168.1.101.35140	147.210.19.227.993	TCP - imap sécurisé (ssl..	9714
11/1/2013 20:19:3.272	147.210.19.227.993	192.168.1.101.35140	TCP - imap sécurisé (ssl..	37026
11/1/2013 20:20:12.984	127.0.0.1.8080	127.0.0.1.34928	TCP -	1
11/1/2013 20:21:40.48	173.194.69.16.993	192.168.1.101.53021	TCP - imap sécurisé (ssl..	1346
11/1/2013 20:21:40.48	173.194.69.16.993	192.168.1.101.53022	TCP - imap sécurisé (ssl..	1321
11/1/2013 20:21:40.48	192.168.1.101.53021	173.194.69.16.993	TCP - imap sécurisé (ssl..	1032
11/1/2013 20:21:40.48	192.168.1.101.53022	173.194.69.16.993	TCP - imap sécurisé (ssl..	1032
11/1/2013 20:21:58.497	147.210.19.227.993	192.168.1.101.35158	TCP - imap sécurisé (ssl..	5070
11/1/2013 20:21:58.497	192.168.1.101.35158	147.210.19.227.993	TCP - imap sécurisé (ssl..	1665
11/1/2013 20:21:58.637	147.210.19.227.993	192.168.1.101.35155	TCP - imap sécurisé (ssl..	6028

FIGURE 3.1 – Page d'accueil

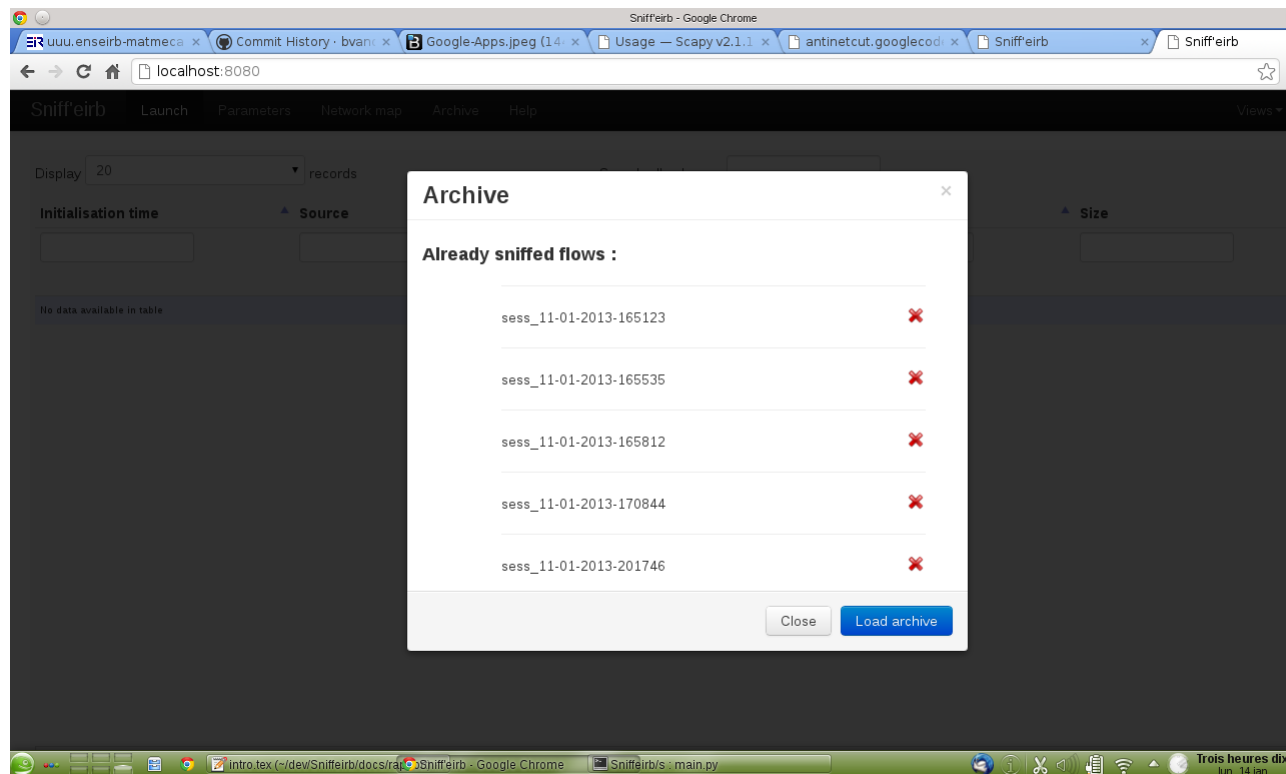


FIGURE 3.2 – Archive

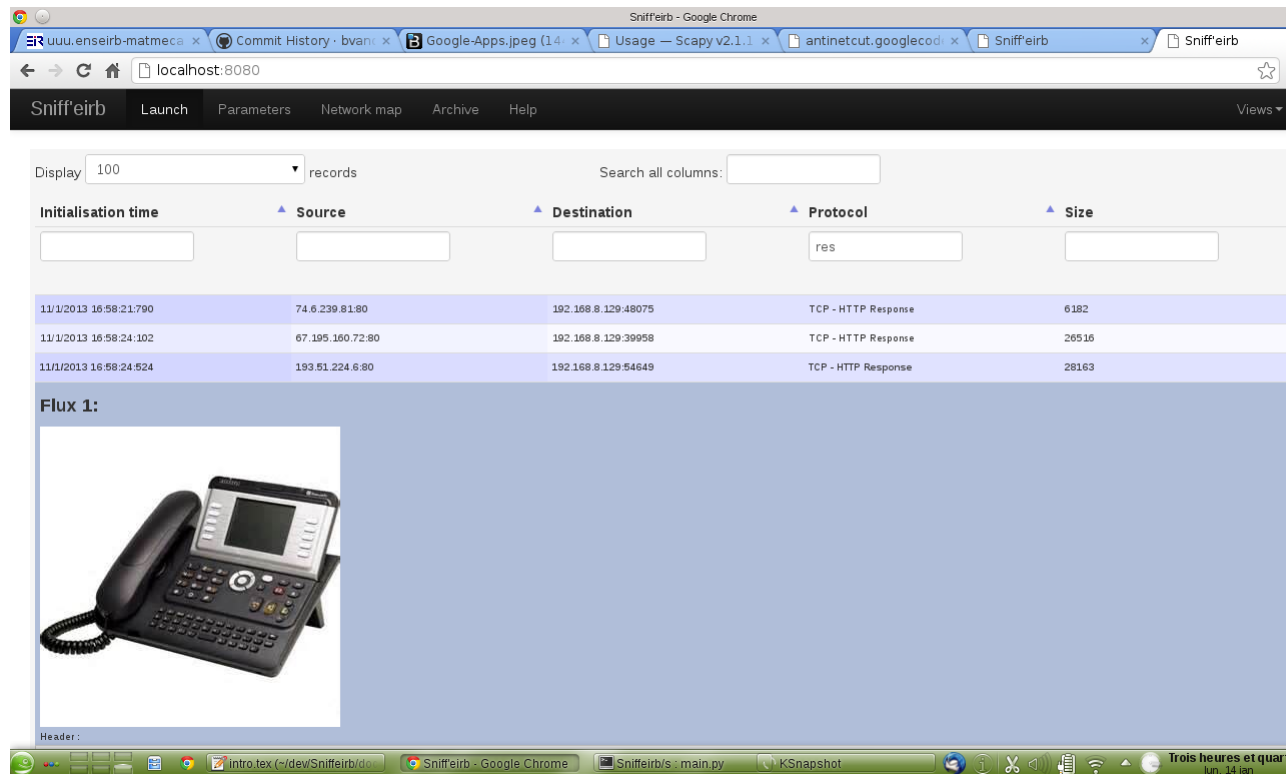


FIGURE 3.3 – Vu d'une image reconstruite


```
{
  "_id" : ObjectId("50edc0fe10b95715ceddf583"),
  "dport" : 80,
  "dst" : "50.57.168.107",
  "initTS" : 1357758718.609804,
  "packets" : [
    {
      "ack" : 0,
      "flags" : "S",
      "ts" : 1357758718.609804,
      "seq" : NumberLong("3047389275")
    },
    {
      "ack" : NumberLong("3368619615"),
      "flags" : "A",
      "ts" : 1357758718.736191,
      "seq" : NumberLong("3047389276")
    },
    {
      "ack" : NumberLong("3368619615"),
      "flags" : "A",
      "ts" : 1357758722.540476,
      "seq" : NumberLong("3047389276")
    }
  ],
  "proto" : "TCP",
  "session" : "sess_09-01-2013-201054",
  "sport" : 53879,
  "src" : "192.168.1.101",
  "type" : "IP"
}
```

FIGURE 3.4 – Extrait de la base Mongo : un document