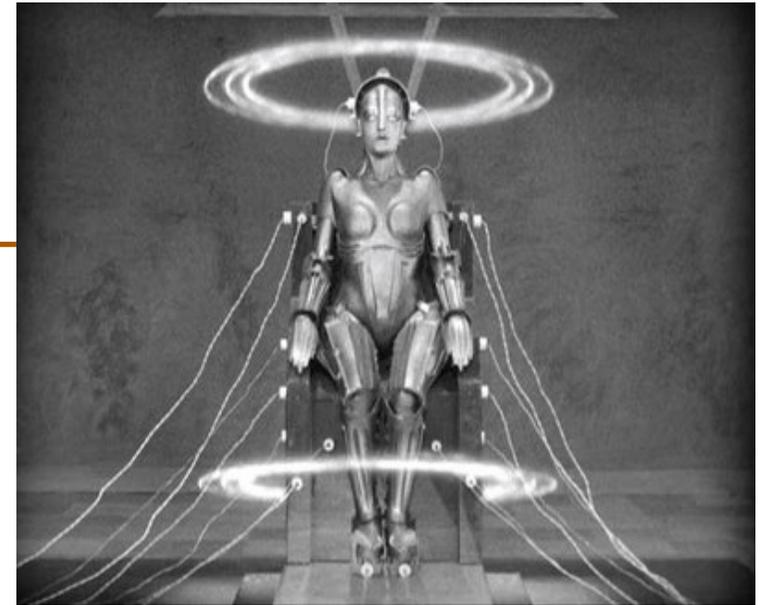


 Observatoire
de la vie littéraire



LIP6

Des *Big Data* aux *Big Brothers*

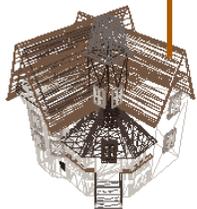
Jean-Gabriel Ganascia

Equipe ACASA – LIP6 – Université Pierre and Marie Curie

Labex OBVIL – PRES “Sorbonne Université”

4, place Jussieu, 75252 Paris Cedex 05, FRANCE

Jean-Gabriel.Ganascia@lip6.fr




SORBONNE UNIVERSITÉS

Synoptique

- *Big Data*
 - *Déluge informationnel*
 - *Les 3 V*
 - *La logique des Big Data*
- *« Big Brother »*
 - *Données personnelles*
 - *Protection?*
- *Qui sont les « Big Brothers »?*
 - *Utilisation des Big Data*
 - *Nouvelle économie*
- *Conséquences politiques*



Big Data – masses de données

Déluge informationnel



- Livres
- Gazouillis (tweets, ...)
- Traces
 - Caméras de surveillances
 - Empreintes digitales
 - Requêtes google
- Archives individuelles
- Science
 - Génome, séquence ADN
 - Astronomie
- ...



Unité de mémoire: le nœud de mouchoir

1 nœud de mouchoir = 1 « bit »



1 caractère typographique = 8 mouchoirs



Accroissement des quantités Masses de données - « Big Data »

Quantité d'information

- *1 livre, 1 million de caractères*
→ $1000 \text{ Ko} = 1\text{Mo} (10^6 \text{ octets}) = 8 \text{ millions de mouchoirs}$
- *Catalogue des livres et imprimés de la BNF:*
 $14 \text{ millions d'ouvrages} = 14\text{To}$
 $1 \text{ Go} = 10^3\text{Mo} = 10^9\text{o} (1 \text{ milliard})$
 $1 \text{ To} = 10^{12}\text{o} (1000 \text{ milliards})$



Production des données en nombre de BNF

- **Twitter** produit 7 To/jour
c'est-à-dire ½ BNF!
- **Facebook**: 500 To/jour
- Radio télescope *Murchison
Widefield Array* (Australie):
 - **données brutes** 7000 To/
minutes, c'est-à-dire 500
BNF/minute
 - **données analysées**: 50To
par jour (4 BNF/jour)
- **Volume total du web en 2015**:
7 Zeta-octets = $7 \cdot 10^{21} \text{o}$

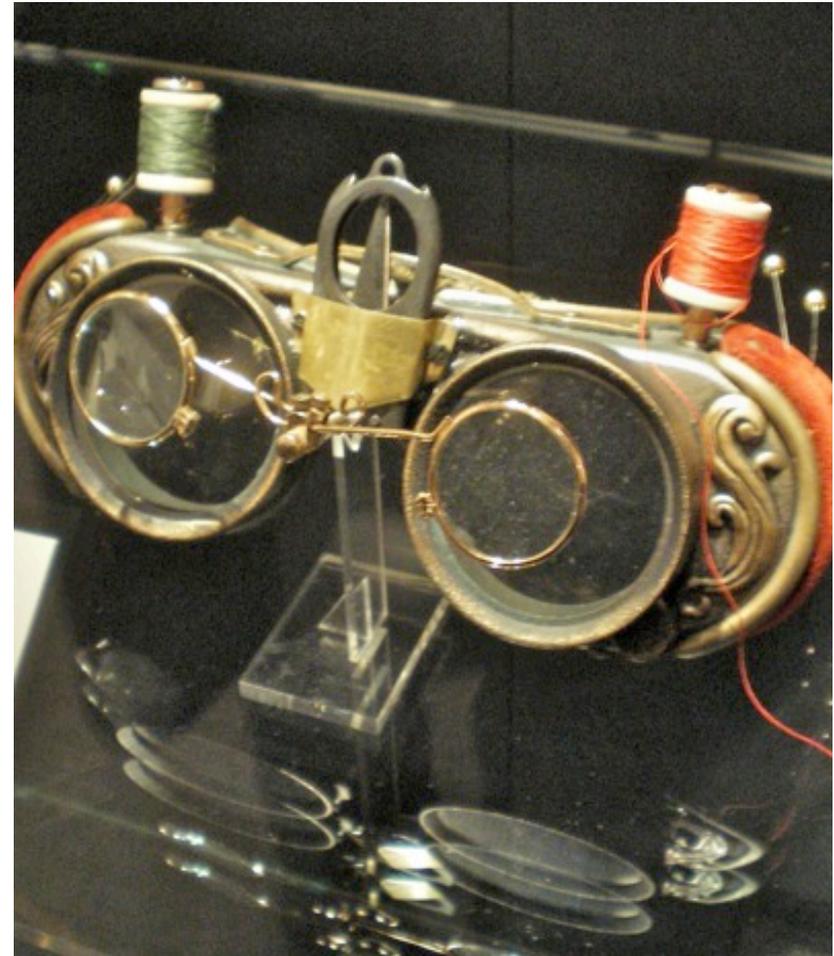


½ milliard de BNF

Big Data

Les 3 V

- **Volume:** To – Po (de 10^{12} à 10^{15})
- **Vélocité:** données changeantes
- **Variabilité:** non organisées, toutes les langues, mélange image, son, ...

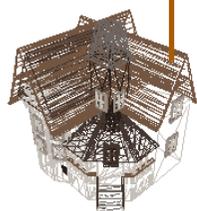


La logique des masses de données



Trois caractéristiques:

- **Les 3V**
 - Volume
 - Variété
 - Vélocité
- **Acquisition continue**
 - « Crowdsourcing »
 - requêtes, ...
- **Absence d'échantillonnage**
 - Recherche de d'anomalies

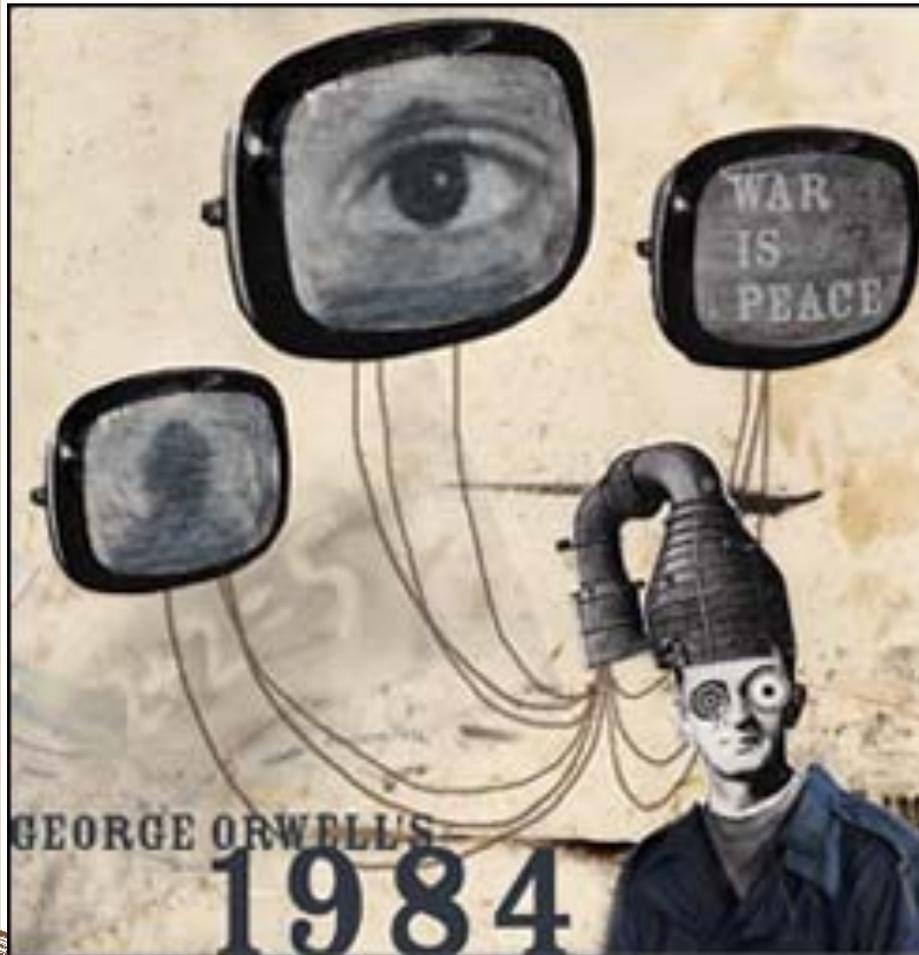


Synoptique

- *Big Data*
 - *Déluge informationnel*
 - *Les 3 V*
 - *La logique des Big Data*
- « Big Brother »
 - *Données personnelles*
 - *Protection?*
- *Qui sont les « Big Brothers »?*
 - *Utilisation des Big Data*
 - *Nouvelle économie*
- *Conséquences politiques*



1984 est il devant ou derrière nous?



- Progrès des dispositif de prise d'information personnelle
- Nous sommes suivis à la trace



Ces traces que nous laissons!



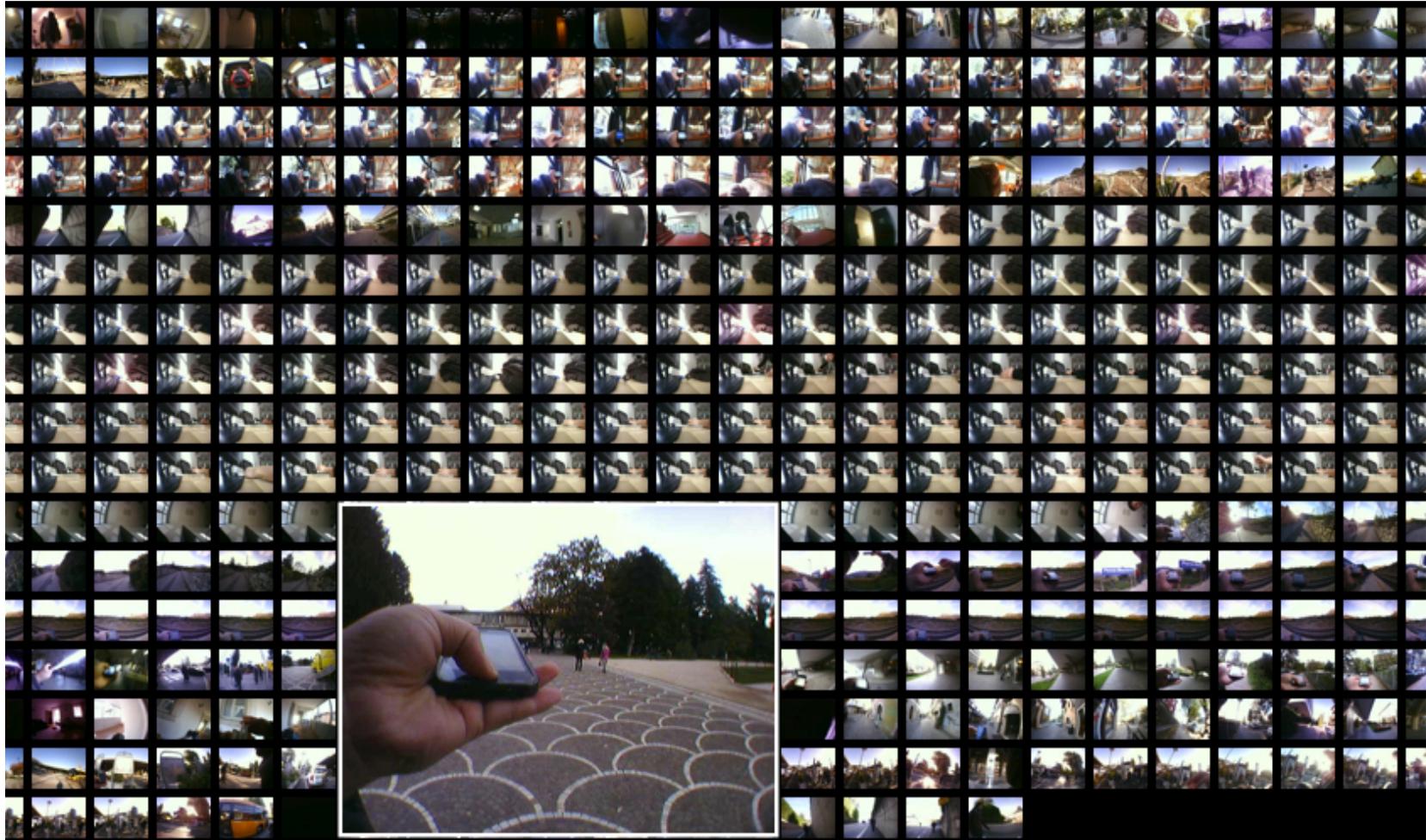
- Conscientes
 - Réseaux sociaux
 - Blogs, écrits, ...



Archives individuelles



« Lifelog » – mémoire de vie



Kiwi: traces physiologiques



[Home](#)

[Product](#)

[Developers](#)

[About Us](#)

[Contact](#)

PRODUCT INFO



Temperature



Heart Rate



Blood Pressure



Location (GPS)



Activity/Motion



Time

Body Vitals

The Kiwi accurately and continuously measures your body vitals using the latest in sensor technology. Using real time data from different sensors the Kiwi is able to derive blood pressure, stress, and a host of other signals from your body



Context

Combining rich sensor data with context is where the real magic happens. Without user intervention the Kiwi is able to distinguish between sleep and a variety of different modes (running, walking, sitting). With this data, the number of potential applications is endless!



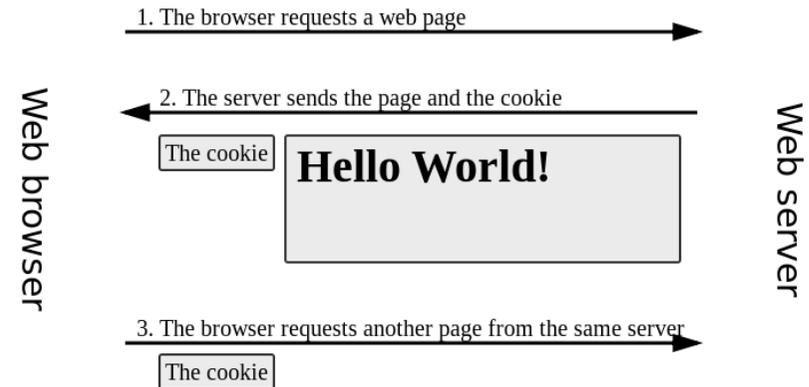
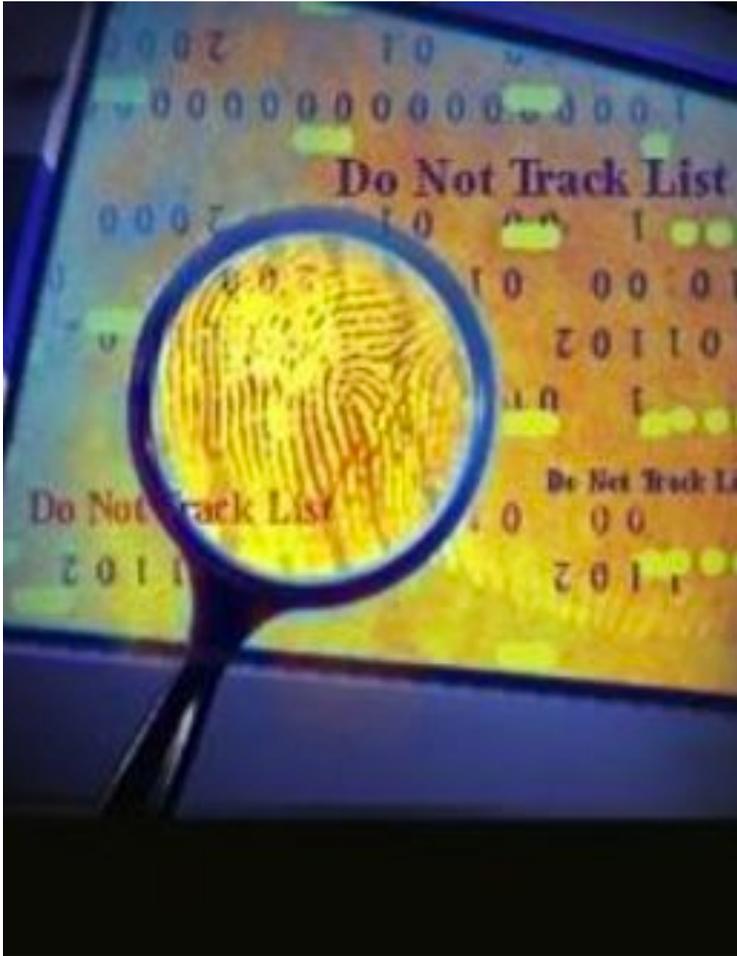
Ces traces que nous laissons!



- Conscientes
 - Réseaux sociaux
 - Blogs, écrits, ...
- Par devers nous
 - Paiements bancaires
 - Cartes de santé
 - Passe navigo, ...
- Inconscientes
 - Coups de fils, mails...
 - Cookies



Cookies

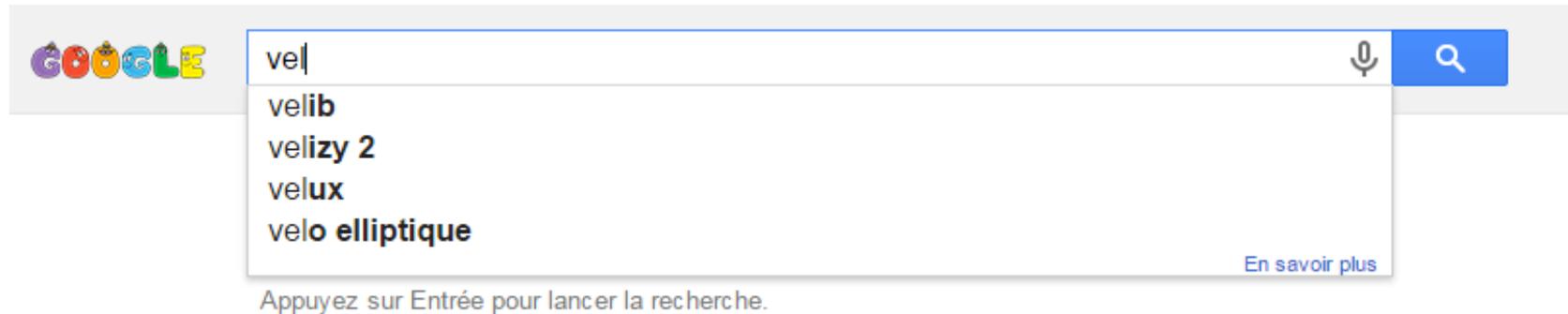


```
HTTP/1.1 200 OK
Cache-Control: private
Content-Type: text/html

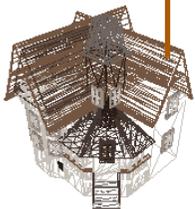
Set-Cookie: PREF=ID=5e66ffd215b4c5e6:
TM=1147099841:LM=1147099841:S=Of69MpW
Bs23xeSv0; expires=Sun, 17-Jan-2038 1
9:14:07 GMT; path=/; domain=.google.c
om
```



Requêtes sur moteurs de recherche



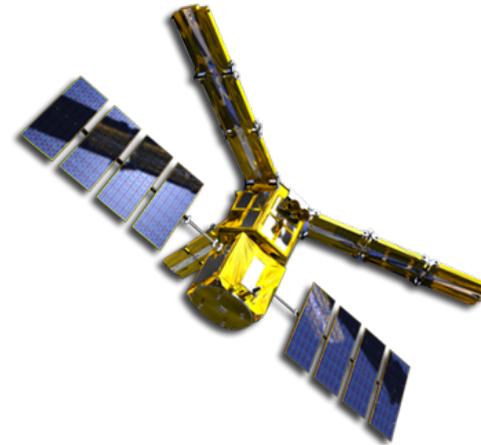
- Auto-complétion des requêtes
- Suivi et profilage
 - Publicité ciblée
 - Sécurité...
- Exploitation massive de ces requêtes – *Big Data*



Echapper?



- Peut-on l'espérer en se contentant de cultiver son jardin?
- Pas pour longtemps...

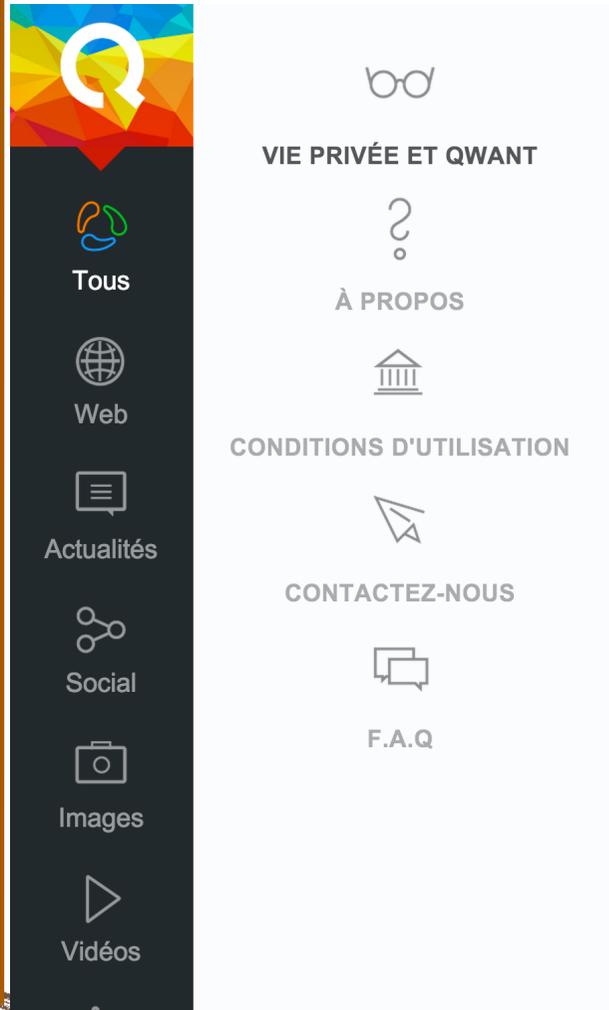


Synoptique

- *Big Data*
 - *Déluge informationnel*
 - *Les 3 V*
 - *La logique des Big Data*
- *« Big Brother »*
 - *Données personnelles*
 - *Protection?*
- *Qui sont les « Big Brothers »?*
 - *Utilisation des Big Data*
 - *Nouvelle économie*
- *Conséquences politiques*



Changer de moteur de recherche



Connexion

FR



Vie privée et Qwant

Qu'est-ce qu'un cookie ?



Les cookies sont des informations enregistrées par le navigateur servant à stocker diverses données telles que votre identifiant de session ou vos préférences.

La philosophie de Qwant repose sur 2 principes : ne pas tracer les utilisateurs et ne pas filtrer le contenu d'internet. Nous faisons notre possible pour respecter la vie privée des internautes tout en garantissant un environnement sécurisé et des résultats pertinents.

Principe de finalité de la CNIL

- **Big Data:** recueil systématique, sans hypothèse *a priori*
- Comment respecter ce principe dans le contexte des *Big Data*?

La finalité des traitements

Un fichier doit avoir un objectif précis.

Les informations exploitées dans un fichier doivent être **cohérentes par rapport à son objectif**.

Les informations **ne peuvent pas être réutilisées de manière incompatible avec la finalité** pour laquelle elles ont été collectées.

Tout détournement de finalité est passible de 5 ans d'emprisonnement et de 300 000 € d'amende.

art. 226.21 du code pénal



Principe de proportionnalité de la CNIL

- **Big Data:** conservation, partage et accumulation de toutes les données
- Comment respecter ce principe dans le contexte des *Big Data*?

La durée de conservation des informations

Les données personnelles ont une **date de péremption**.

Le responsable d'un fichier fixe **une durée de conservation raisonnable** en fonction de l'objectif du fichier.

Le code pénal sanctionne la conservation des données pour une durée supérieure à celle qui a été déclarée de 5 ans d'emprisonnement et de 300 000 € d'amende.

art. 226-20 du code pénal



Peut-on oublier sur Internet?

- Fonctionnement moteurs de recherche:
 - Aspiration du web avec des bots, puis analyse
 - Indexation des sites:

- Principe: dictionnaire inversé:

De « collectionner »

et « monnaies »,

retrouver

« numismate » et

exclure

« philatéliste » parce

qu'il y a « monnaies »

philatéliste

n. m. et f.

Personne qui collectionne les timbres-poste.

numismate

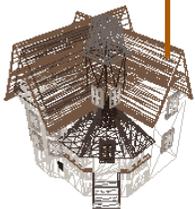
n. m. et f.

Personne qui collectionne les monnaies anciennes et des médailles.



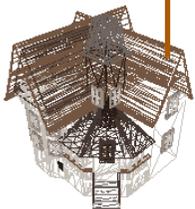
Peut-on oublier sur Internet?

- Fonctionnement moteurs de recherche:
 - Aspiration du web avec des bots, puis analyse
 - Indexation des sites:
 - Principe: dictionnaire inversé
 - Ordonnancement des sites:
 - Élimination des mots sans intérêt (TFIDF)
 - Mots trop fréquents
 - Mots trop rares (hapax)
 - Algorithme de préférences (ex. Google)



Indexation moteurs de recherche

- **Mot₁**: url_{1,1}, url_{1,2}, url_{1,3}, ... url_{1,n}, url_{1,n+1}, ...
- **Mot₂**: url_{2,1}, url_{2,2}, url_{2,3}, ... url_{2,m}, url_{2,m+1}, ...
- ...
- **Mot_p**: url_{p,1}, url_{p,2}, url_{p,3}, ... url_{p,l}, url_{p,l+1}, ...
- ...



Oubli sur Internet $\text{Mot}_p \text{ url}_{p,1}, \text{ url}_{p,2}$

Sociétés commerciales de déréférencement:

- Mot_1 : $\text{url}_{1,1}, \text{url}_{1,2}, \text{url}_{1,3}, \dots, \text{url}_{1,n}, \text{url}_{1,n+1}, \dots$
- Mot_2 : $\text{url}_{2,1}, \text{url}_{2,2}, \text{url}_{2,3}, \dots, \text{url}_{2,m}, \text{url}_{2,m+1}, \dots$
- ...
- Mot_p : $\text{url}_{p,q}, \text{url}_{p,q+1}, \text{url}_{p,q+2}, \dots, \text{url}_{p,r}, \dots$
 $\text{url}_{p,1}, \text{url}_{p,2}, \text{url}_{p,3}, \dots, \text{url}_{p,l}, \text{url}_{p,l+1}, \dots$
- ...

Action: ajout information neutre

But: enfouir l'information gênante



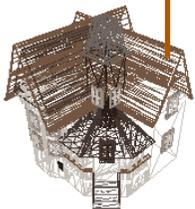
Oubli sur Internet $\text{Mot}_p \text{ url}_{p,1}, \text{ url}_{p,2}$

« Droit à l'oubli » Google

- $\text{Mot}_1: \text{url}_{1,1}, \text{url}_{1,2}, \text{url}_{1,3}, \dots, \text{url}_{1,n}, \text{url}_{1,n+1}, \dots$
- $\text{Mot}_2: \text{url}_{2,1}, \text{url}_{2,2}, \text{url}_{2,3}, \dots, \text{url}_{2,m}, \text{url}_{2,m+1}, \dots$
- ...
- $\text{Mot}_p: \text{url}_{p,1}, \text{url}_{p,2}, \text{url}_{p,3}, \dots, \text{url}_{p,l}, \text{url}_{p,l+1}, \dots$
- ...

Suppression $\text{Mot}_p: \text{url}_{p,1}, \dots, \text{url}_{p,l}, \text{url}_{p,l+1},$

Les URLs $\text{url}_{p,1}, \text{url}_{p,2},$ restent présentes, mais associées à d'autres noms...



Synoptique

- *Big Data*
 - *Déluge informationnel*
 - *Les 3 V*
 - *La logique des Big Data*
- *« Big Brother »*
 - *Données personnelles*
 - *Protection?*
- *Qui sont les « Big Brothers »?*
 - *Utilisation des Big Data*
 - *Nouvelle économie*
- *Conséquences politiques*

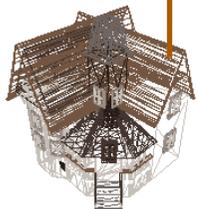


Nouvelle révolution scientifique

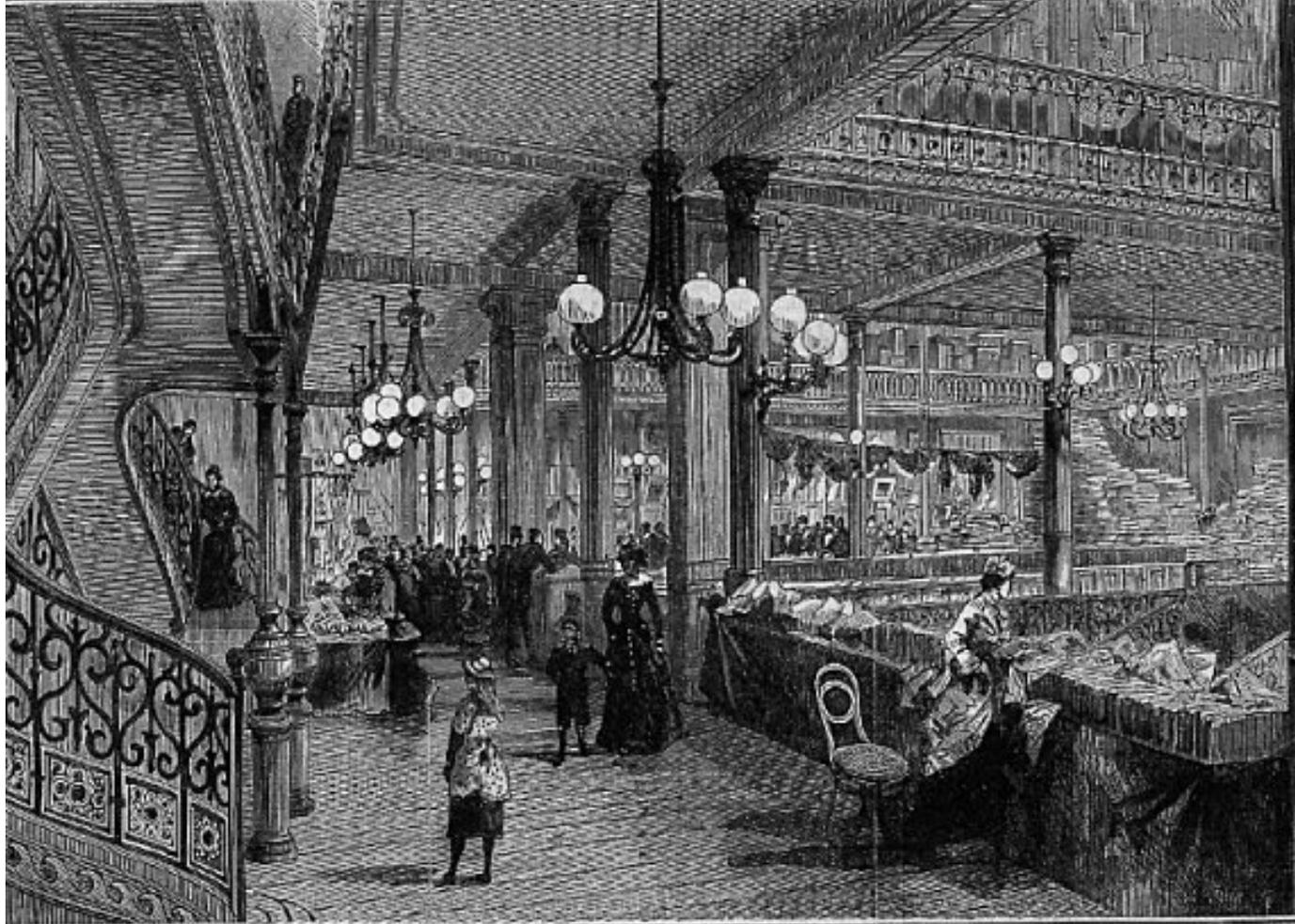
- ◆ Science moderne
expérimentation *in vivo* et *in vitro* Intervention
sur le monde extérieur



- ◆ Science *in silico*
Expérimentation sur données



« Au bonheur des dames » *Emile Zola*



Marchandisation des identités

*De la
« femme sans tête »
à la
« femme 100 têtes »*

**Importance du
profilage...**

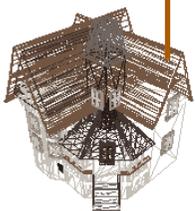


Qui sont les Big Brothers?

- Gouvernements
 - PRISM (Etats-Unis)
 - ...



**BIG BROTHER IS
WATCHING YOU**



PRISM

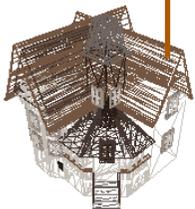
Differents programmes de surveillance

- « UPSTREAM »: Aspiration données sur cables sous-marins
 - Mails
 - Téléphones
- Requêtes auprès des acteurs de l'internet
 - AOL, Skype, Microsoft, Facebook, Yahoo, Apple, Google, YouTube, ...



Qui sont les Big Brothers?

- Gouvernements
 - PRISM (Etats-Unis)
 - ...
- Sociétés privées
 - GAFA (*Goole, Apple, Amazon, Facebook*)
 - Nouvelle économie fondée sur le profilage



La prise de pouvoir des GAFA

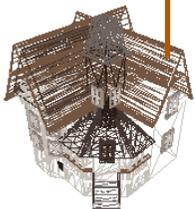
Prérogatives régaliennes des Etats

- Battre monnaie
- Organiser la santé
- Gérer l'état civil
- S'occuper de la formation
- Dresser le cadastre
- ...

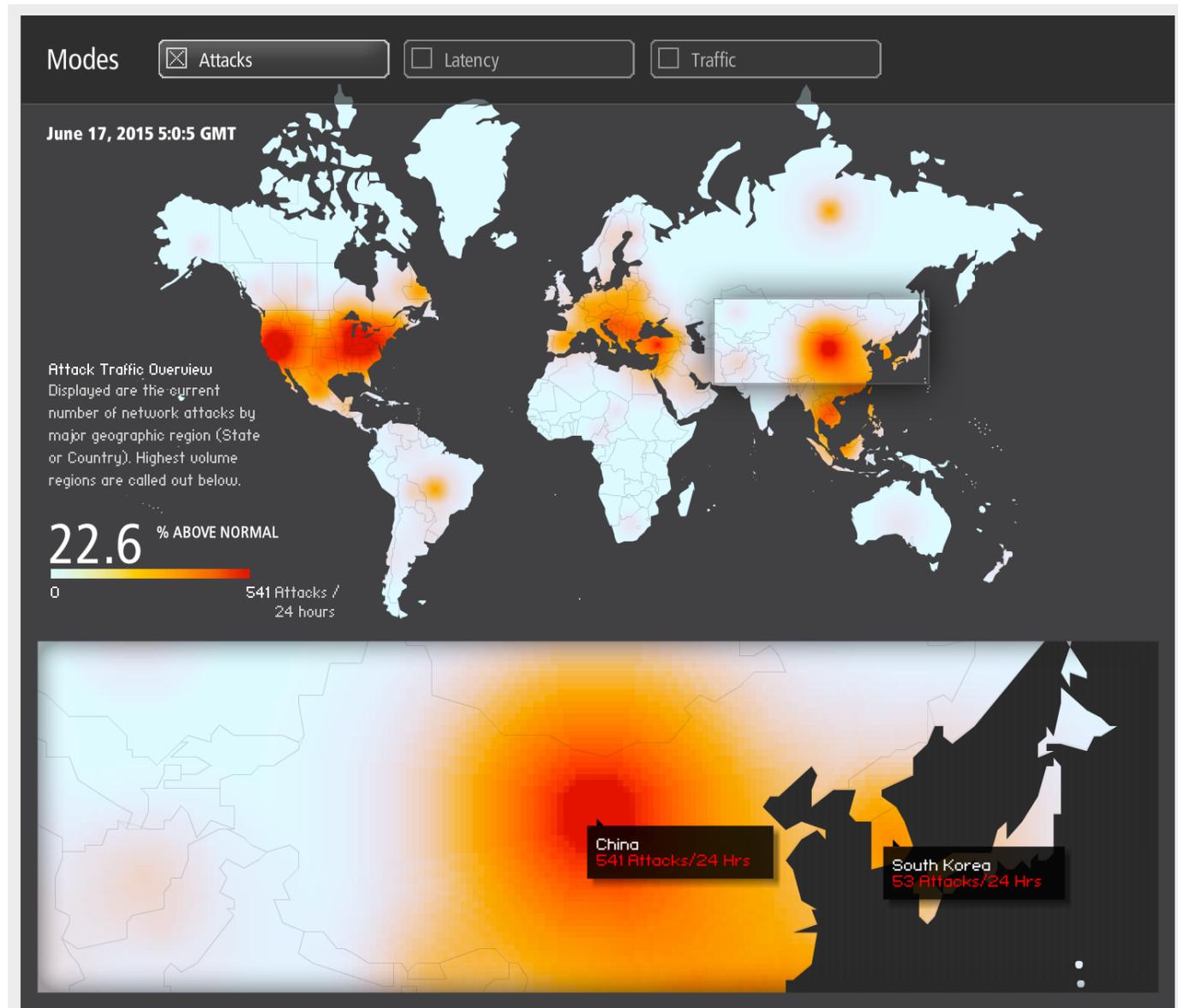


Qui sont les Big Brothers?

- Gouvernements
 - PRISM (Etats-Unis)
 - ...
- Sociétés privées
 - GAFA (*Goole, Apple, Amazon, Facebook*)
 - Nouvelle économie fondée sur le profilage
- Nouveaux acteurs
 - Hommes politiques
 - Maffias
 - ONG
 - ...



Cyber-attaques



Cyber-attaques



LIFE IS FOR SHARING.

OVERVIEW

STATISTICS

INFO

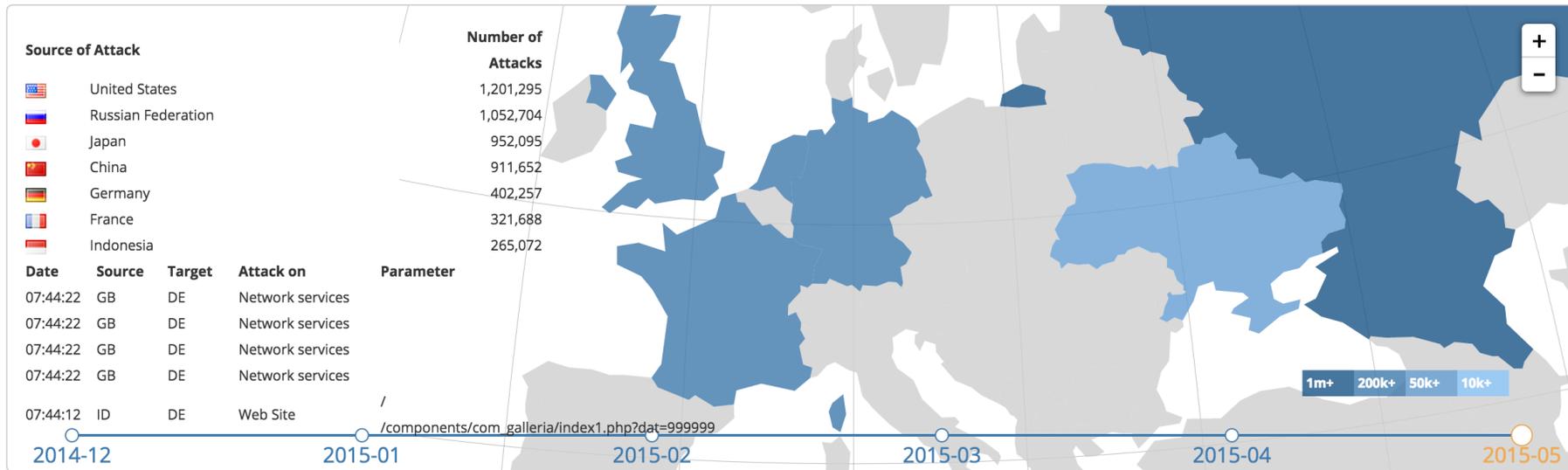
DOWNLOAD

IMPRINT

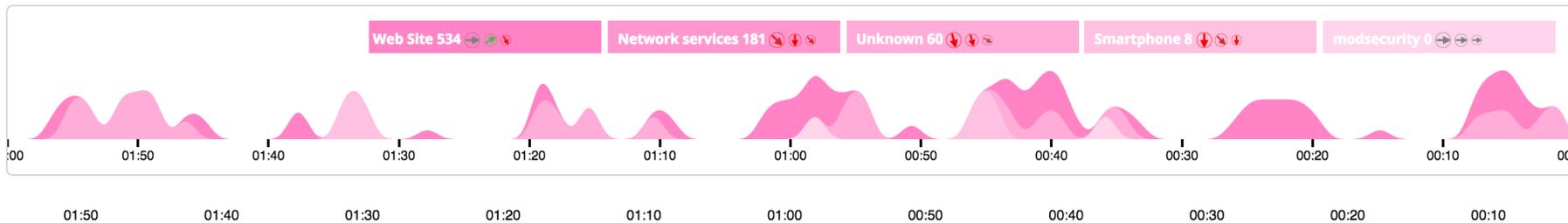
DTAG



Overview of current cyber attacks on DTAG sensors (logged by 180 Sensors)



Trend Analyse



Nombre d'attaques



LIFE IS FOR SHARING.

OVERVIEW

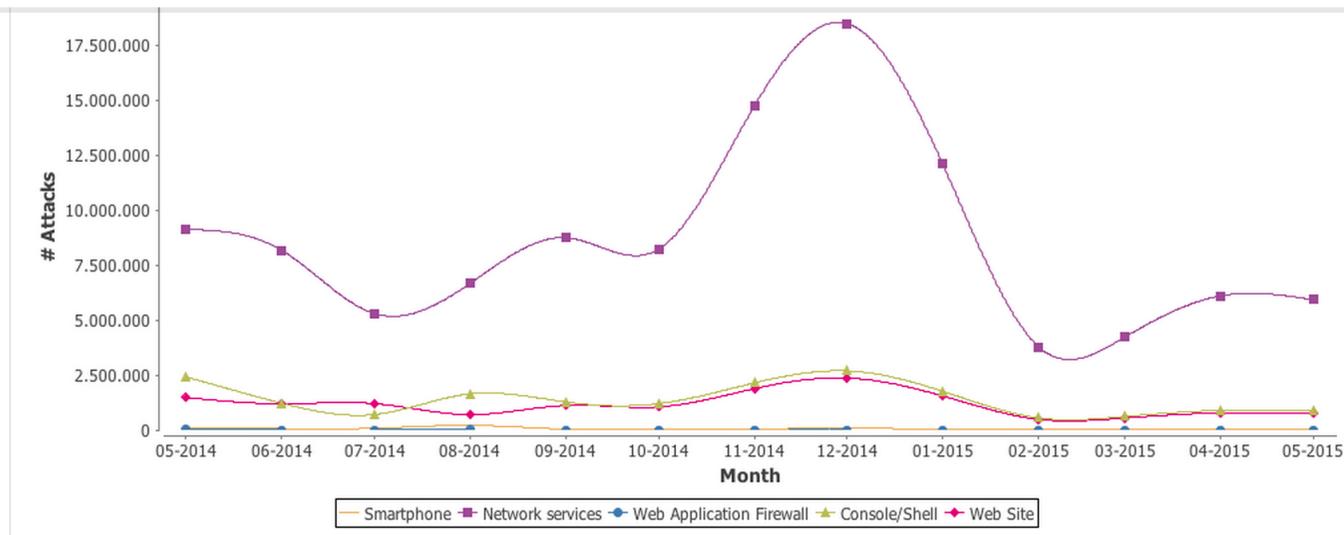
STATISTICS

INFO

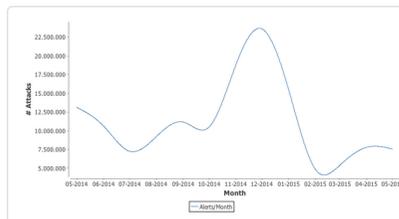
DOWNLOAD

IMPRINT

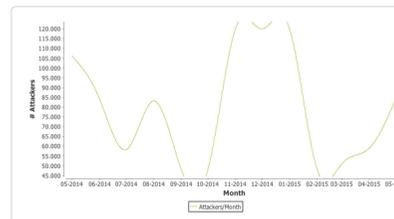
DTAG ▾



Overall sum of attacks per Month



Overall sum of attackers per Month



Passwords



ONG – Maffias – Etats



« 1984 » est-il derrière ou devant?

L'année 1984 est en arrière...

Les dispositifs de prise d'information sont en avance

Faut-il craindre la généralisation de la surveillance?

Ou la perte de souveraineté numérique au profit de nouveaux acteurs?

