



# W3Perl

---

## Statistiques de consultation de sites Web

Laurent Domisse  
(INRP)



# Rôle

---

- **Transparence / Service**
- **Surveillance**
  - Bande passante
  - Impact événementiel
- **Amélioration du site**
  - Accroître la visibilité
  - Connaître son public (langues, références...)
  - Résoudre les problèmes

**TOP SECRET**



# Différentes méthodes

---

- **Compteur**
  - Insertion de balises...
- **Fichier de log**
  - Serveur web...
- **Audit**
  - Application propriétaire...



# Fichier de log

---

- **Serveur**
  - Unité de mesure : Requêtes/Pages
- **Gestion des logs**
  - Rotation
  - Compression
- **Format des logs**



# Format des logs

---

- **De nombreuses variantes**
  - CLF : standard
  - ECLF, W3C, IIS, Lotus, Squid...
- **ECLF**

```
198.45.0.12 - - [24/Mar/2002:15:13:27 -0500] "GET /fr/ HTTP/1.1" 200 45  
"http://www.mydomain.com/" "Mozilla/4.0 (compatible; MSIE 6.0; Windows NT 5.0)"
```

```
host login passwd date hour methode URL protocole status taille referer agent
```



# Extraction

---

- **Stats de base**
  - Page, Site, Répertoires, Fichier, Téléchargement, Trafic, Pays, Scripts
- **Stats temporelles**
  - Heures, Jours, Semaines, Mois
- **Stats optionnelles**
  - Références, navigateur, OS, virus, session, erreur, URL mapping



# W3Perl

---

- **Perl / ActivePerl**
  - Script par stats, librairie
- **Serveurs Unix-Linux / NT**
- **Fichier de configuration**
  - Définition des règles
- **Interface Web d 'administration**
- **Mode incrémental**



# Installation

---

- Perl / Fly
- Windows (9x/NT/XP/2000)
  - IIS / Apache ou sans serveur
- Unix
  - Install.pl
    - Scripts --> /cgi-bin/w3perl
    - Ressources --> /htdocs/w3perl
  - RPM en cours



# Fichier de config.

---

- Un fichier par serveur / stats
- Gestion par interface web
- Etapes
  - Serveur/Fichier de log
  - Affichage
  - Filtres
  - Préférences



# Interface d'administration

---

- **Fichier de configuration**
  - Créer / Modifier / Effacer / Dupliquer
- **Lancement**
  - Contrôle à distance
- **Mise à jour**

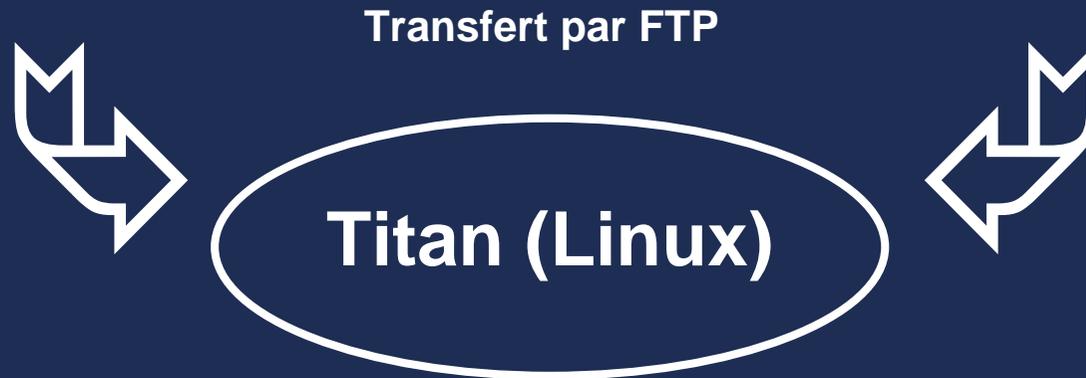


# Configuration INRP

---

Netscape Proxy (NT)  
ECLF

Squid (Linux)  
CLF + refer + agent



Grep sur la date courante  
Filtrage des images  
Concaténation + Compression



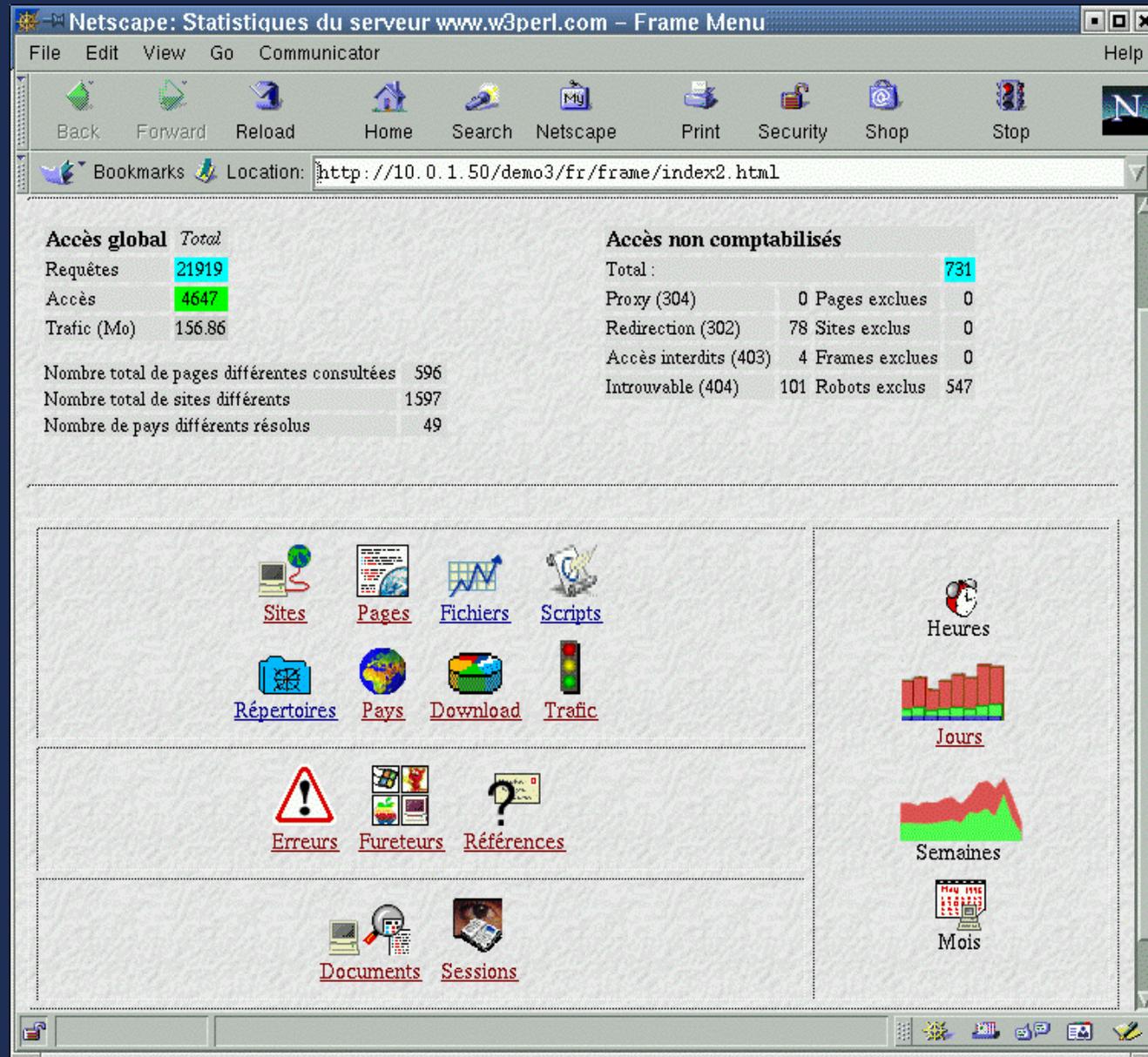
# Lancement

---

- **Initialisation / Incrémental**
- **Script de lancement**
  - Gère les autres scripts
  - Date de lancement
  - Option de lancement
  - Crontab
- **Interface web d'administration**



# Page principale



Bourse aux Outils



28-03-2002

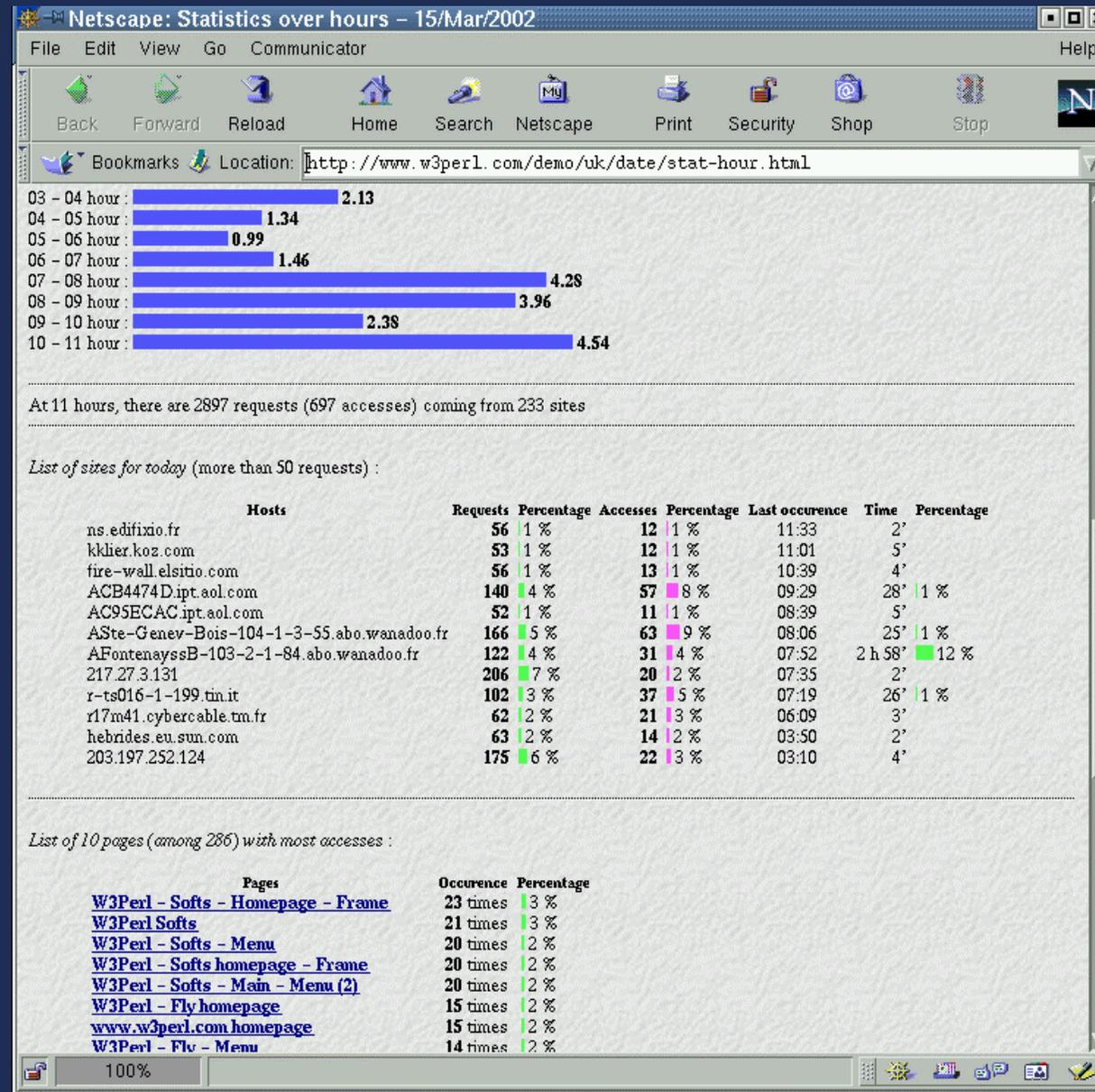
# Proxy

---

- **Cache web**
  - réduire la bande passante
  - sécurité
- **Transaction**
  - 200 : OK
  - 304 : Not modified



# Stats par heures

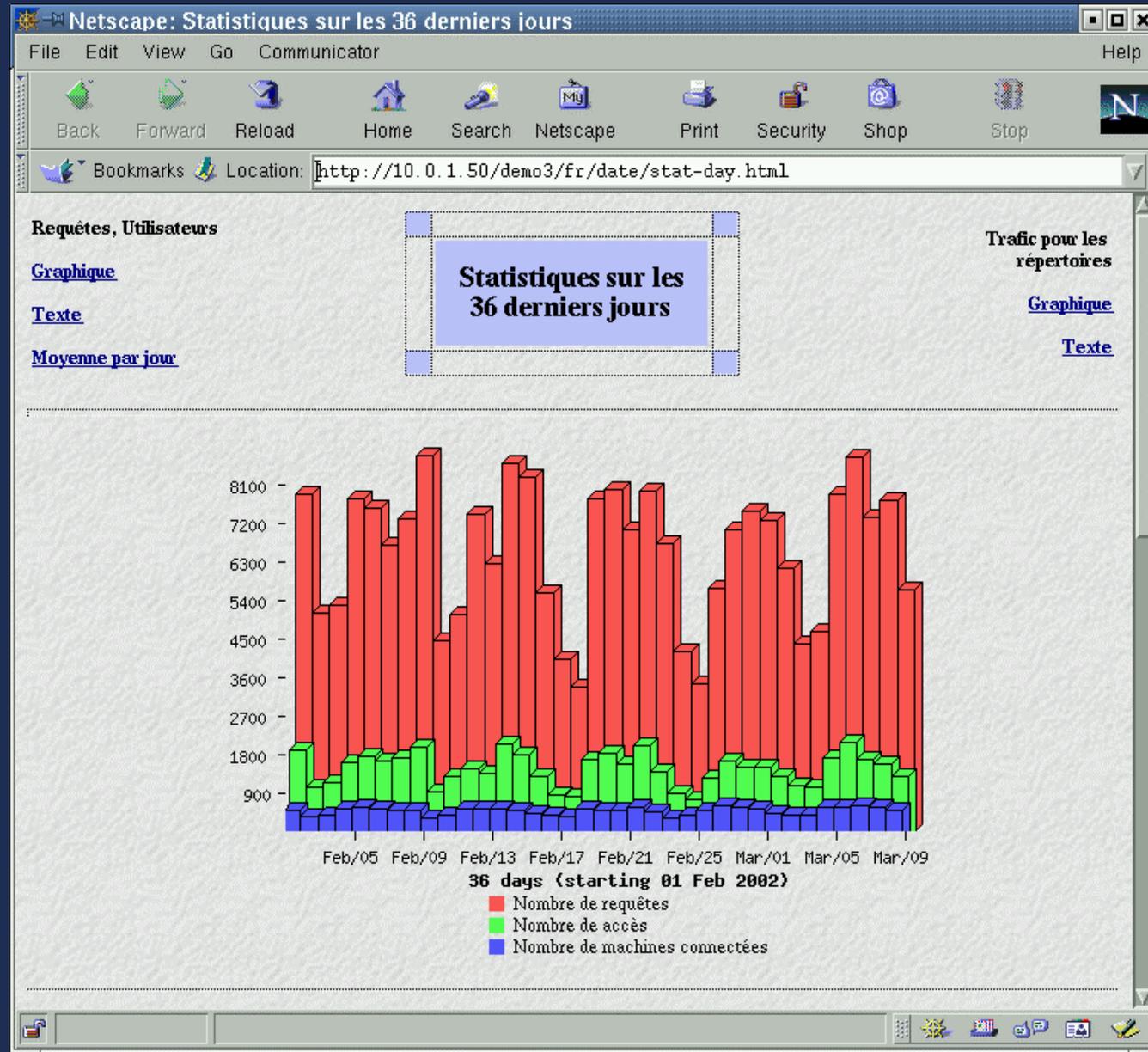


Bourse aux Outils



28-03-2002

# Stats sur les jours

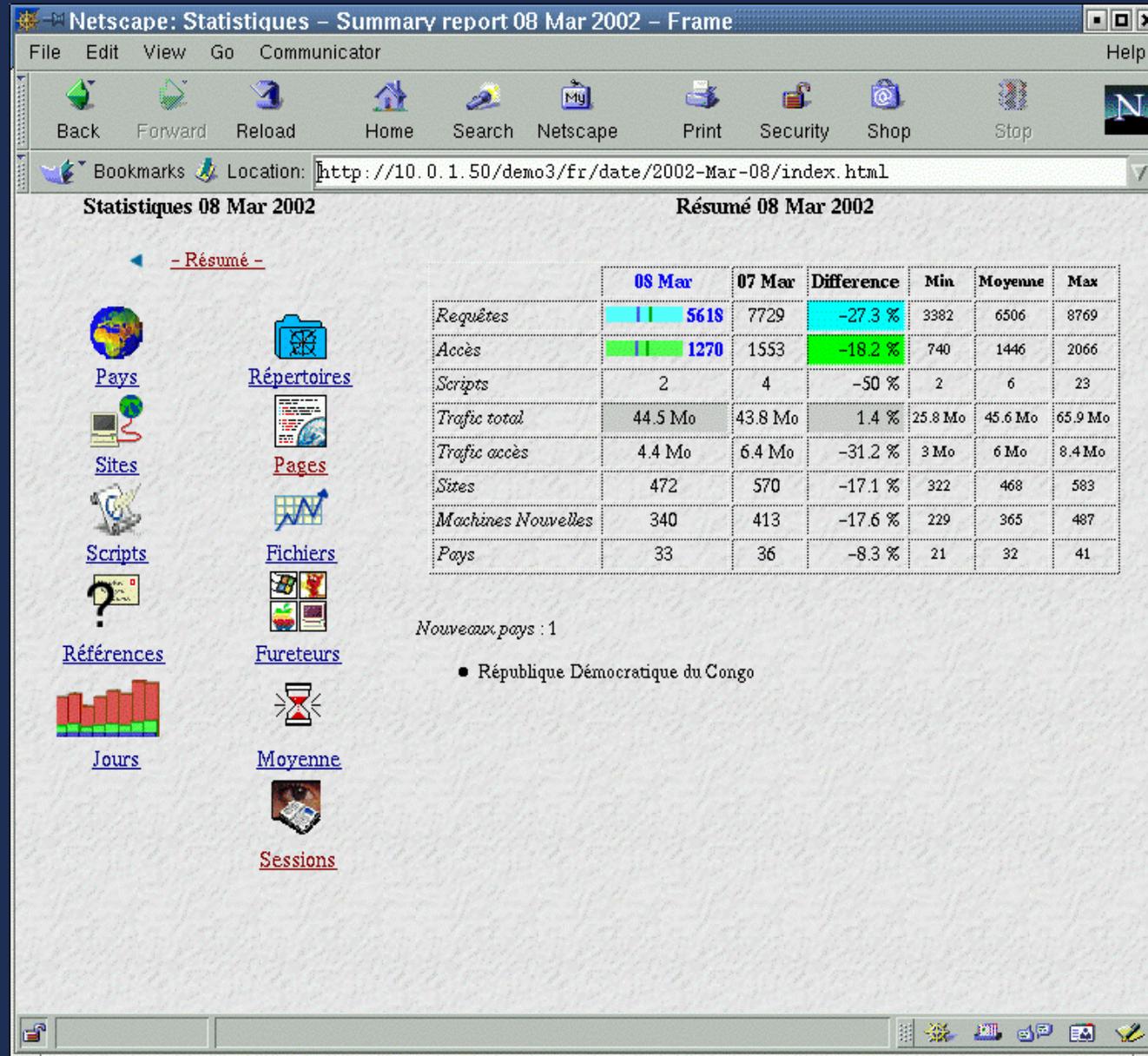


Bourse aux Outils



28-03-2002

# Stats quotidiennes

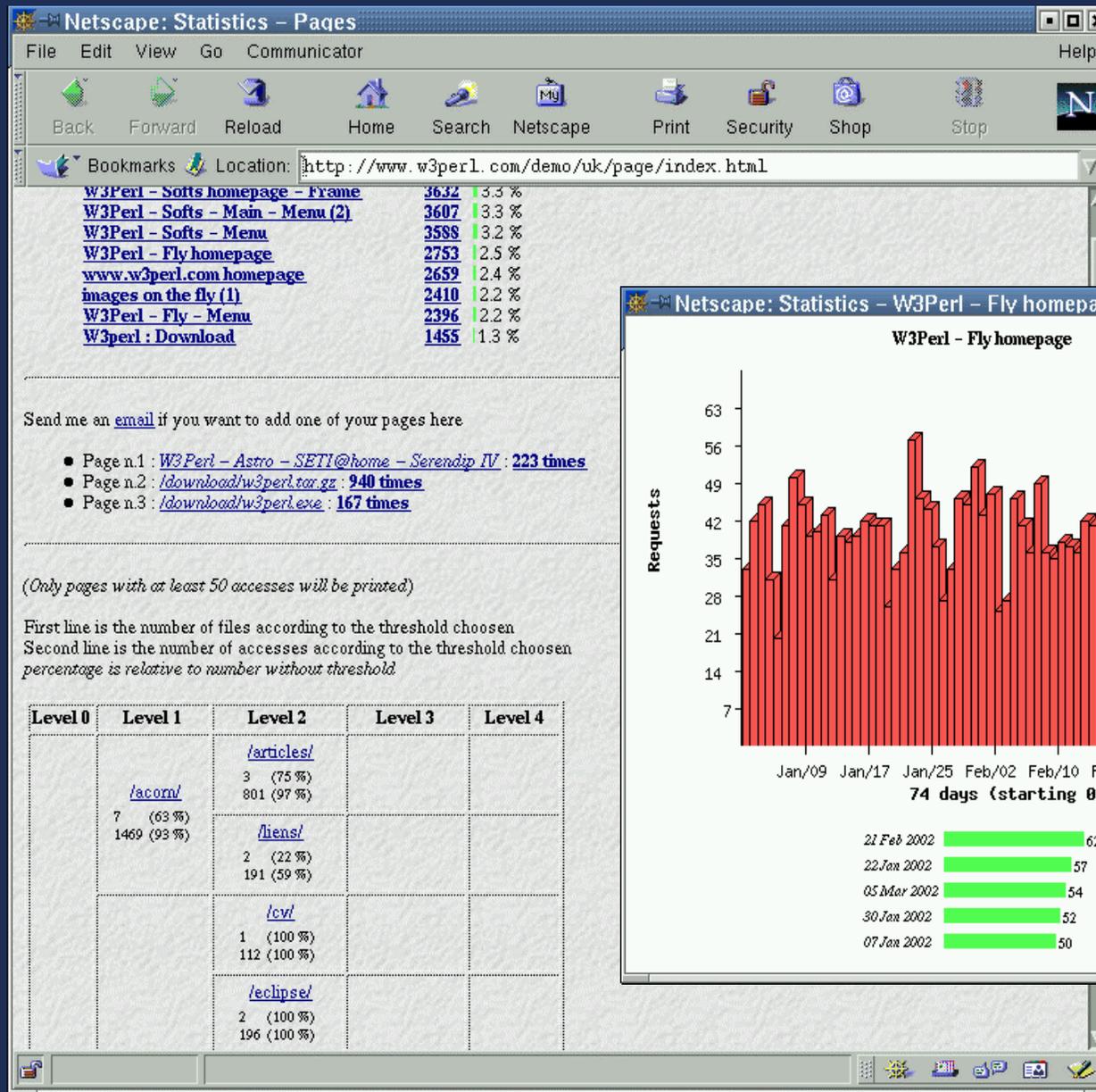


Bourse aux Outils



28-03-2002

# Stats sur les pages



Bourse aux Outils



28-03-2002

# Stats sur les pays

Netscape: Statistiques - pays

File Edit View Go Communicator Help

Back Forward Reload Home Search Netscape Print Security Shop Stop

Location: http://10.0.1.50/demo3/fr/pays/index.html

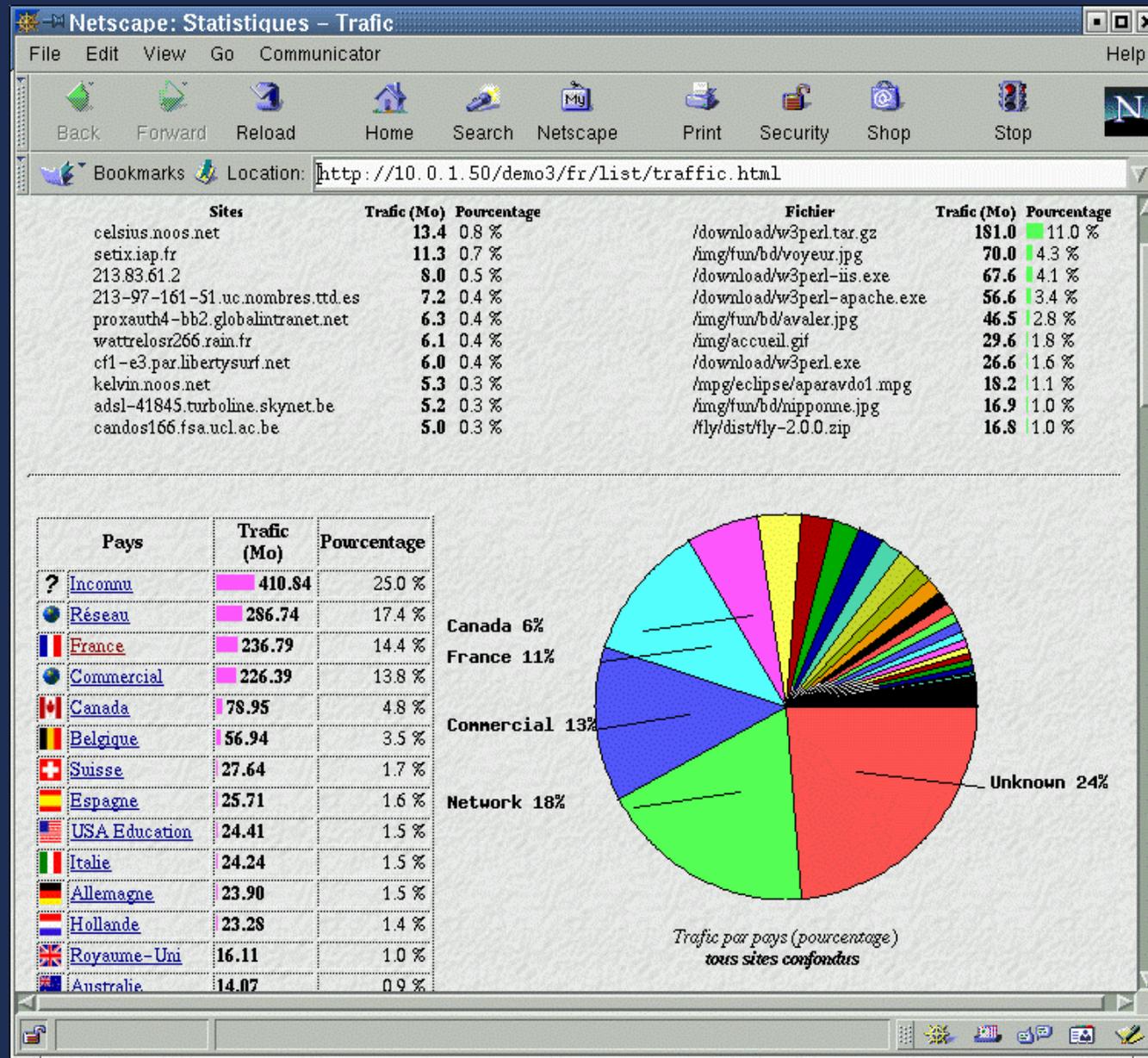
Pays	Requêtes	%	Accès	%	Sites	%	Trafic (Mo)	%
 Réseau	46486	19.8	9954	19.1	2267	17.2	286.7	17.4
 Commercial	32657	13.9	7281	14.0	1956	14.9	226.4	13.8
 France	31225	13.3	7742	14.9	2346	17.8	236.8	14.4
 Canada	8181	3.5	2105	4.0	860	6.5	78.9	4.8
 Belgique	6078	2.6	1809	3.5	464	3.5	56.9	3.5
 Italie	4461	1.9	877	1.7	99	0.8	24.2	1.5
 USA Education	4292	1.8	880	1.7	256	1.9	24.4	1.5
 Hollande	3544	1.5	689	1.3	139	1.1	23.3	1.4
 Suisse	3460	1.5	824	1.6	197	1.5	27.6	1.7
 Allemagne	3243	1.4	768	1.5	220	1.7	23.9	1.5
 Royaume-Uni	3010	1.3	551	1.1	152	1.2	16.1	1.0
 Espagne	2842	1.2	625	1.2	131	1.0	25.7	1.6
 Australie	2038	0.9	390	0.7	103	0.8	14.1	0.9
 Japon	1985	0.8	427	0.8	113	0.9	6.3	0.4
 Brésil	1801	0.8	346	0.7	51	0.4	12.0	0.7
 Pologne	1553	0.7	348	0.7	37	0.3	11.7	0.7
 Danemark	1342	0.6	271	0.5	48	0.4	5.5	0.3
 Suède	1305	0.6	260	0.5	79	0.6	13.0	0.8
 Fédération de Russie	1253	0.5	204	0.4	29	0.2	8.0	0.5
 Argentine	1143	0.5	190	0.4	19	0.1	4.9	0.3
 Finlande	1005	0.4	184	0.4	43	0.3	8.4	0.5
 Autriche	992	0.4	163	0.3	48	0.4	3.0	0.2
 Taiwan	911	0.4	255	0.5	20	0.2	5.5	0.3
 Organisations	847	0.4	194	0.4	59	0.4	4.9	0.3
 Israël	743	0.3	143	0.3	26	0.2	3.3	0.2

Bourse aux Outils



28-03-2002

# Stats téléchargements

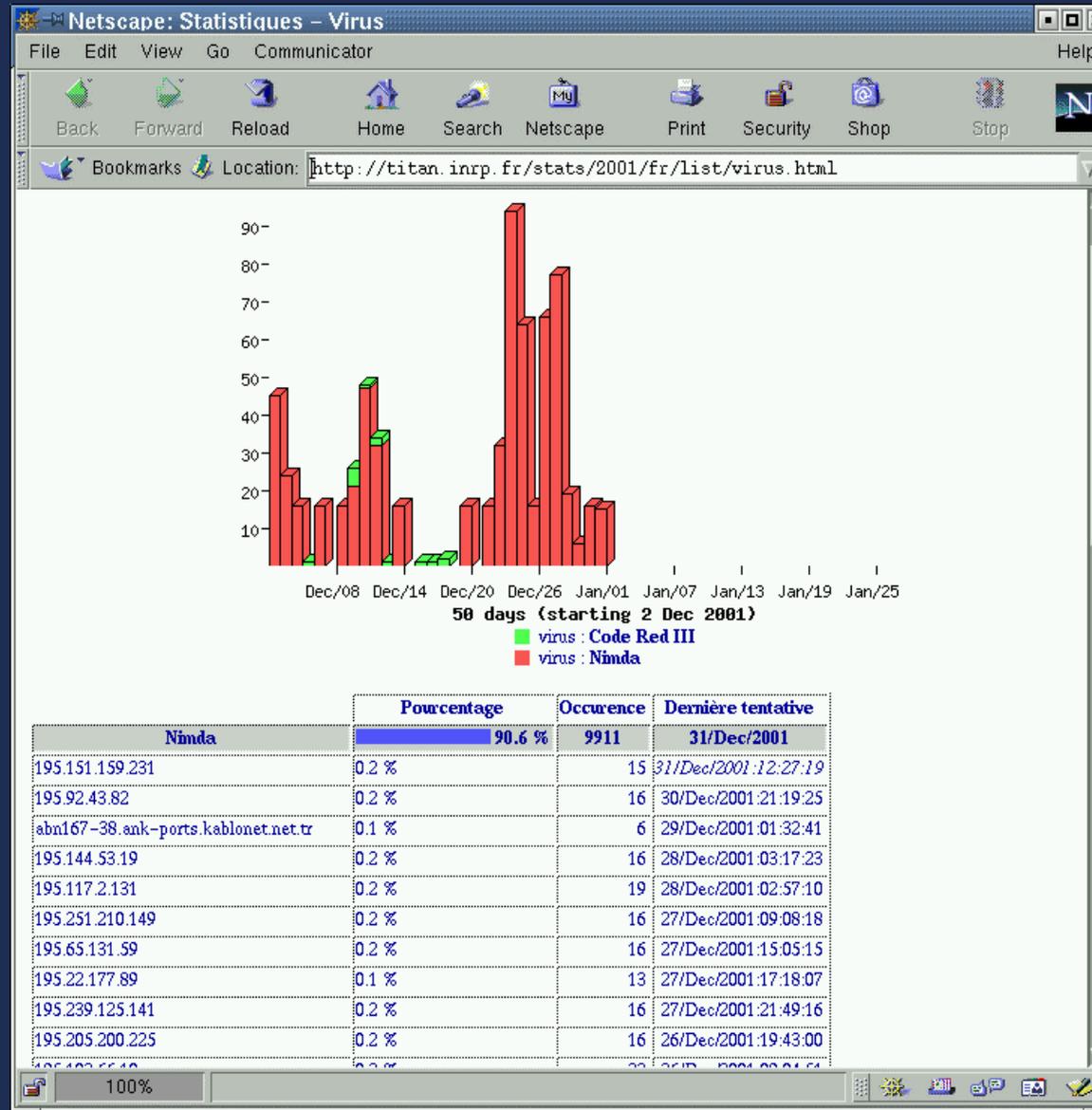


Bourse aux Outils



28-03-2002

# Stats sur les virus

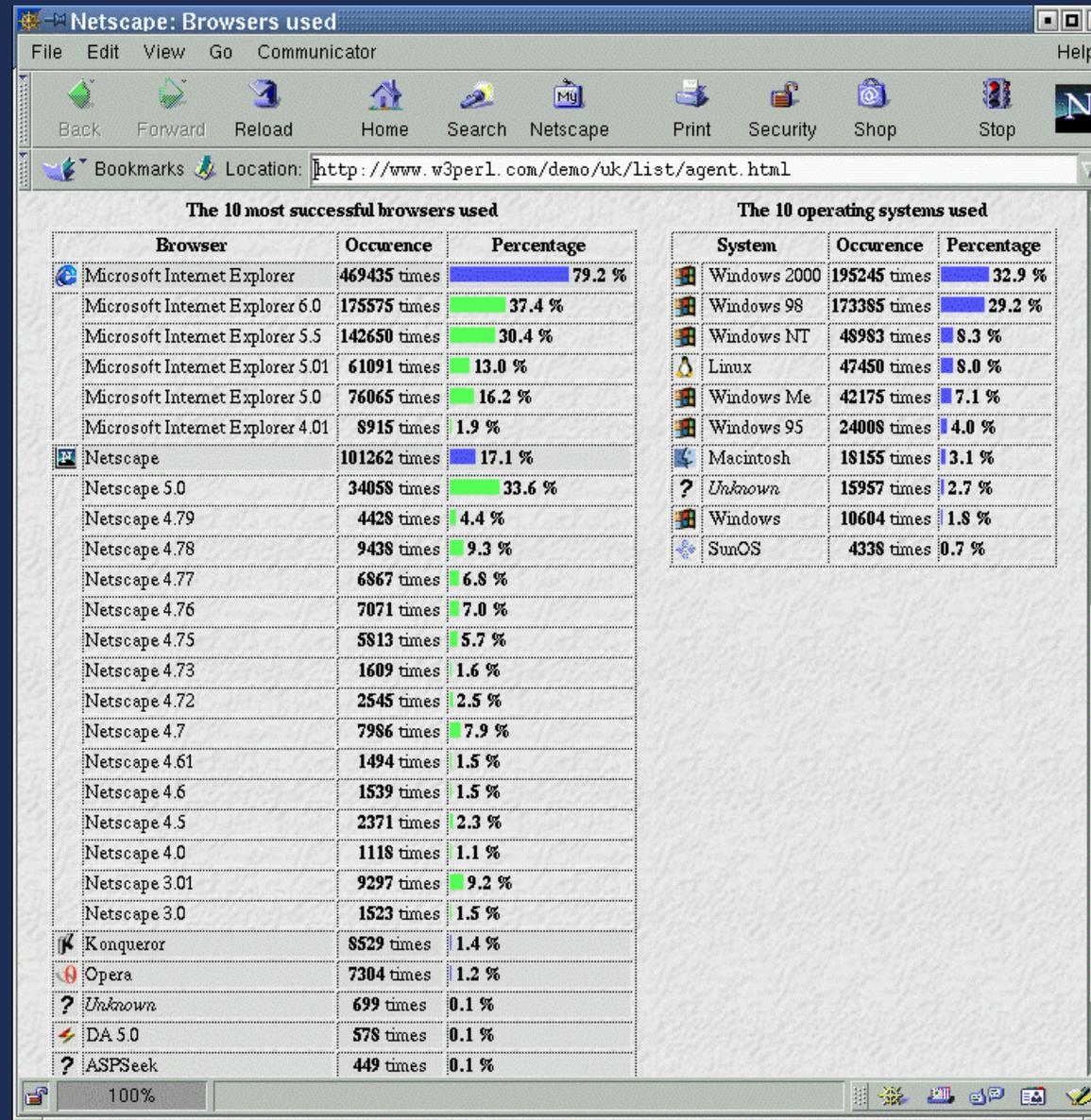


Bourse aux Outils



28-03-2002

# Stats sur les navigateurs



Bourse aux Outils



28-03-2002

# URL mapping

The screenshot shows the Netscape browser window titled "Netscape: Statistiques sur les fichiers HTML". The address bar contains the URL "http://10.0.1.50/demo3/fr/document/index.html". The main content area displays the following statistics:

**Moyenne :**

- *Fichiers HTML* : **5 liens, 6 images** et ont un poids de **3 Kilo-octets**
- *Images* ont un poids de **15.5 Kilo-octets**
- *Fichiers HTML* ont un poids de **97.0 Kilo-octets**
- *Temps de chargement* pages HTML : **19 secondes** (avec un modem 28800)

Below the statistics, there is a grid of links to various analysis tools:

 <a href="#">Liens</a>	 <a href="#">Images</a>	 <a href="#">Pages les plus récentes</a>
 <a href="#">Poids des pages HTML</a>	 <a href="#">Répertoires</a>	 <a href="#">Arborescence</a>
 <a href="#">Titre identique</a>	 <a href="#">Liens absolus</a>	 <a href="#">Liens symboliques</a>
 <a href="#">WIDTH, HEIGHT</a>	 <a href="#">ALT</a>	 <a href="#">Liens erronés</a>
 <a href="#">Pages inutilis</a>	 <a href="#">Pages sans index</a>	

Bourse aux Outils



28-03-2002

# Session : problème

---

- **Pas de support de début / fin**
- **Visite = Machine?**
  - Même machine (IP)
    - Proxy, Provider, Foyer
  - Même personne
    - Provider, DHCP, multi-postes



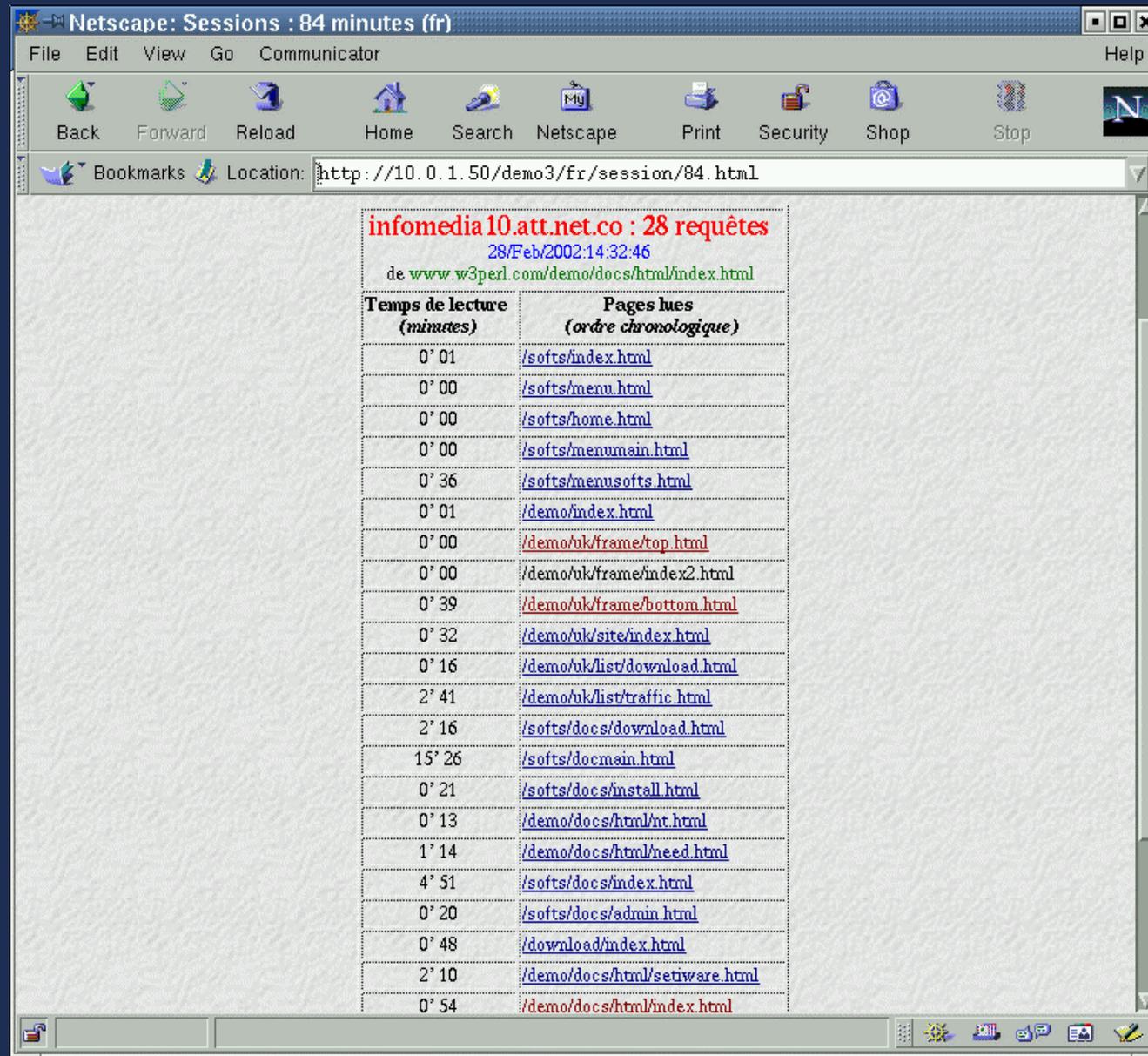
# Session : solution

---

- **Période d'inactivité / Timeout**
- **Filtrage**
  - Basé sur les types de fichiers
- **Résultats**
  - Exploration précise de chaque visite
  - Robots, Login, Fidélité
  - Durée de connexion, par pages



# Exemple de visite



The screenshot shows a Netscape browser window with the title "Netscape: Sessions : 84 minutes (fr)". The address bar contains the URL "http://10.0.1.50/demo3/fr/session/84.html". The main content area displays a table of page visit statistics for the domain "infomedia 10.att.net.co".

infomedia 10.att.net.co : 28 requêtes  
28/Feb/2002:14:32:46  
de www.w3perl.com/demo/docs/html/index.html

Temps de lecture (minutes)	Pages lues (ordre chronologique)
0' 01	<a href="#">/softs/index.html</a>
0' 00	<a href="#">/softs/menu.html</a>
0' 00	<a href="#">/softs/home.html</a>
0' 00	<a href="#">/softs/menomain.html</a>
0' 36	<a href="#">/softs/menusofts.html</a>
0' 01	<a href="#">/demo/index.html</a>
0' 00	<a href="#">/demo/uk/frame/top.html</a>
0' 00	<a href="#">/demo/uk/frame/index2.html</a>
0' 39	<a href="#">/demo/uk/frame/bottom.html</a>
0' 32	<a href="#">/demo/uk/site/index.html</a>
0' 16	<a href="#">/demo/uk/list/download.html</a>
2' 41	<a href="#">/demo/uk/list/traffic.html</a>
2' 16	<a href="#">/softs/docs/download.html</a>
15' 26	<a href="#">/softs/docmain.html</a>
0' 21	<a href="#">/softs/docs/install.html</a>
0' 13	<a href="#">/demo/docs/html/nt.html</a>
1' 14	<a href="#">/demo/docs/html/need.html</a>
4' 51	<a href="#">/softs/docs/index.html</a>
0' 20	<a href="#">/softs/docs/admin.html</a>
0' 48	<a href="#">/download/index.html</a>
2' 10	<a href="#">/demo/docs/html/setiware.html</a>
0' 54	<a href="#">/demo/docs/html/index.html</a>

Bourse aux Outils



28-03-2002

# Futur

---

- **Moyen terme**
  - RPM
  - Squid natif, FTP, IRC....
  - Détection des pirates
- **Long terme**
  - MySQL
  - PDF



# Autres outils gratuits

---

	(+)	(-)
<b>AWStats 3.2</b>	Temps réel	CGI
<b>Analog 5.21</b>	Vitesse	Affichage
<b>Webalizer 2.0</b>	Robuste	Session
<b>W3Perl 2.86</b>	Complet	Vitesse



# Documentation

---

Vous trouverez sur le site de l'application toute la documentation online (y compris cette présentation avec les commentaires associés)

<http://www.w3perl.com/softs/>

Bourse aux Outils



28-03-2002

