# W3Perl

# A free logfile analyzer

# Features

- Works on Unix / Windows / Mac
  - based on Perl scripts
- Web / FTP / Squid / Email servers
  - Others log format can be added easily
- Web Admin Interface
  - Allow to monitor your reports remotely
- PDF / Email / CSV /RSS reports
  - Get reports in your mailbox each morning
- Logfile parsing or page tagging
  - Read server log or log created from a javascript code

- Real time stats
  - View last entries
- Session tracking
  - Follow each user's path
- Detailed daily / weekly … reports
  - Full reports available on a different scale
- AJAX tools
  - Allow to sort table, to search for data …
- Can be run on a remote host
  - Log files can be retrieved, avoiding the need to compute on the server.

# Monitoring technologies

## Page tagging

- Cons :
  - Third party trust
  - No traffic data
  - Rely on javascript
- Pros :
  - No scripts to run or install
  - No caching problem
  - Can add extra data as screensize
  - Can log on javascript events
  - Doesn't require logfile

## Logfile analysis

- Cons :
  - CPU / Memory hungry
  - Need to install a package
- Pros :
  - No change in your website
  - Archive your logfile to retrieve data later
  - Spider activities
  - Fine tuning

W3Perl can be used with both methods

# Logfile

- Choose carefully your logfile format
  - Combined to get browsers/OS/referer stats
  - Reverse DNS to get country stats (w3perl can do it also)
- Hits vs Accesses
- Rotation / Compression
- Recommendation :
  - Apache : move to combined ECLF
  - IIS : use W3C or add extra fields to IIS format
- Create your own logfile if not available
  - Install the Page Tagging package first
- CLF logfile format :

%host %null %login %date %hourshift %method %page %protocol %status %requetesize

# Logfile format

- Web : CLF, ECLF, W3C, IIS …
- Squid : Native, CLF
- FTP : ProFTP, xferlog
- Mail : Exim / Postfix / Sendmail

Your can define your own format using predefined fields

=> %host, %date, %page, %protocol, %status, %requetesize …

# Data mining

- Main stats
  - Page, Hosts, Directories, File, Trafic, Country, City, Scripts ….



- Time stats
  - Real time, hours, days, weeks, months and years



- Additional stats
  - Referer, Browsers, OS, Errors, Sessions, RSS, URL mapping …

# Using W3Perl

- Installation
  - Check requirements (Perl and Fly/Flydraw)
  - Choose Windows Binaries / Linux Packages or Tarball
- Configuration
  - Customize your stats : one configuration per report
  - Build your config from a Web administration interface
- Running the scripts
  - Master script : cron-w3perl.pl
    - -a to initialize
    - -e to update
    - -d <days> to compute the last days
  - Remotely via the Web administration interface or
  - Daily via a crontab
  - Load configuration via the –c flag
- Watch the results !

# Installation

- Unix
  - RPM / Debian package
  - Tarball
    - Scripts -> /cgi-bin/w3perl/
    - Resources -> /htdocs/w3perl/
- Windows binaries (include predefined configuration files)
  - Apache
  - IIS
  - Abyss
  - No server
- Logfile access is not mandatory
  - A javascript tag to be inserted in your web pages can create logfiles for you
- Plug-in
  - PDF, Email, GeoIP ...

Avoid installing the package on a running server

# Installation - Unix

- **RPM / Debian package**
  - Fly or Flydraw package needed
- **Tarball**
  - Install the fly or flydraw package
  - Extract w3perl package in your www document root
  - Edit the install.pl script
  - Change paths according to your system :
    - Perl path
    - CGI path
    - W3Perl package path
  - Run the install.pl
  - Use the web administration interface to build a config file
- **Plug-in**
  - Install optional third-party software : PDF, Email, GeoIP ...

# Installation - Windows

- ActivePerl
  - Install Perl first

- Binaries
  - IIS / Apache / Abyss server with default values
  - No server

- Configuration file
  - Apache / IIS / Abyss are ready to go
  - No server => give log files information
  - Customize values using a web interface

- Plug-in
  - Install optional third-party software : PDF, Email, GeoIP ...

# Configuration

- Create a configuration file
  - Path, threshold, display, log format …
- If a web server is running :
  - Use the web administration interface
- If not :
  - You can use the online tool on w3perl website
  - Or edit manually the config.pl file
- Predefined config files :
  - Mandriva / Ubuntu / Debian
  - IIS / Apache / Abyss
  - FTP / Squid / Postfix / Exim
  - No logfile

# Web Administration Interface

- Manage configuration files
  - Create / Modify / Clone / Delete
- Run remotely your stats
  - Show status with a progress bar
- View stats reports
  - All your reports on one page
- Update package

# Build a configuration file

- Customize display
  - Threshold, Graphs, URL / Hostname mapping, language …
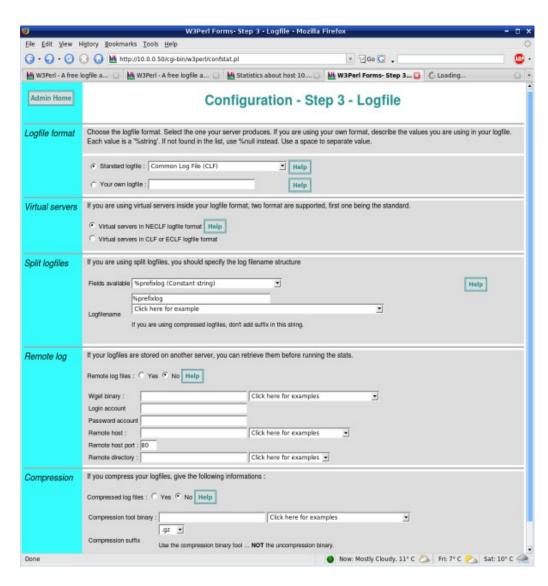
- Logfiles
  - Paths, format, compression, split, filename

- Filtering rules
  - Hosts, URL, Countries, Directories, Robots to exclude
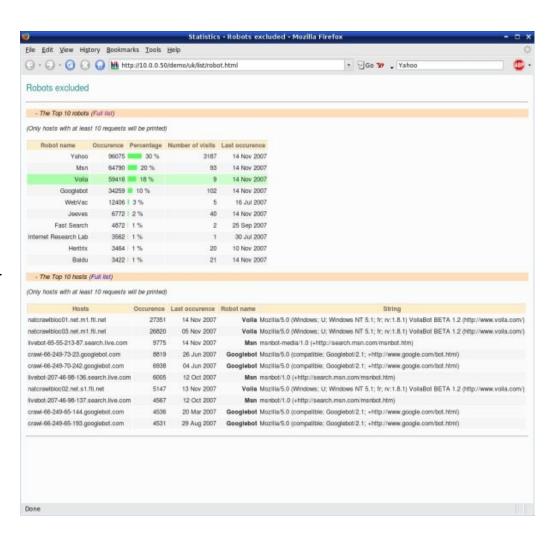
- Processing options
  - Reverse DNS, GeoIPCity, Domain stats, Scripts to run

- Automatic update

- Online tool available



Avoid to build configuration file manually

# Filtering

- Robots detection
  - To reject spider visits
- Referer spam
- Countries
  - You can match your country
- URL
  - Selecting a specific web area
  - User's report
- Directories
  - Avoid scanning private area
- Server's status code
- Hosts
  - Exclude list of hosts

# Running

- From a web admin
  - Remotely
  - Get report in one click
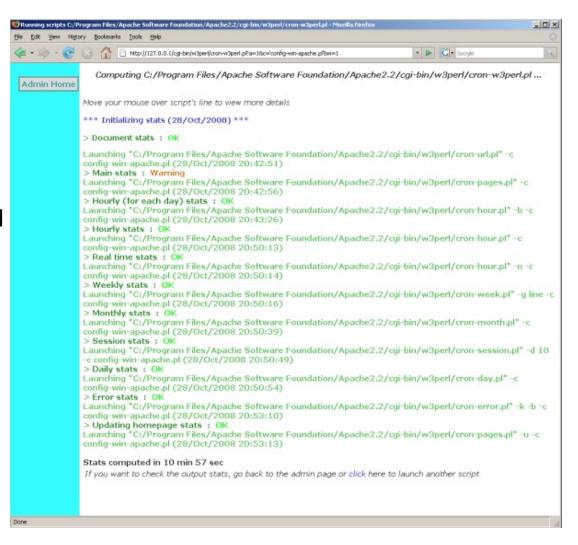  - Click on 'Init' to start
- From command lines
  - Use -h for help
  - Master script : cron-w3perl
  - Flag –c to load a config
  - Flag –a to init
  - Flag –e to update
- Windows menu
  - Init / Update
  - Install / Uninstall
- Options
  - Date selection
  - Graph type
  - … and many more



Don't mix web admin and command line run
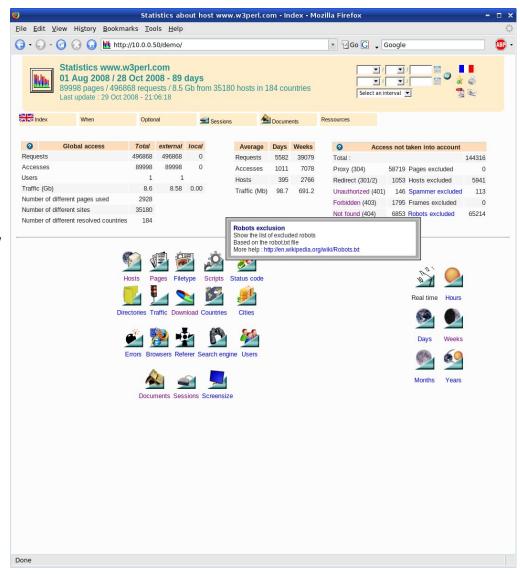
# Reports

- Top menu
  - Summary
  - Tools : date selection, pdf exportation, language ...
- Pure CSS menu
  - Navigate through reports
  - Main/When/Optional/Sessions/Documents/Resources
- Global summary
  - Hosts / Pages / Files / Traffic / Users / Exclude
- Icons field
  - Select reports
  - Updated when new reports are available

# Real time

- Events monitoring
  - Updated every minute
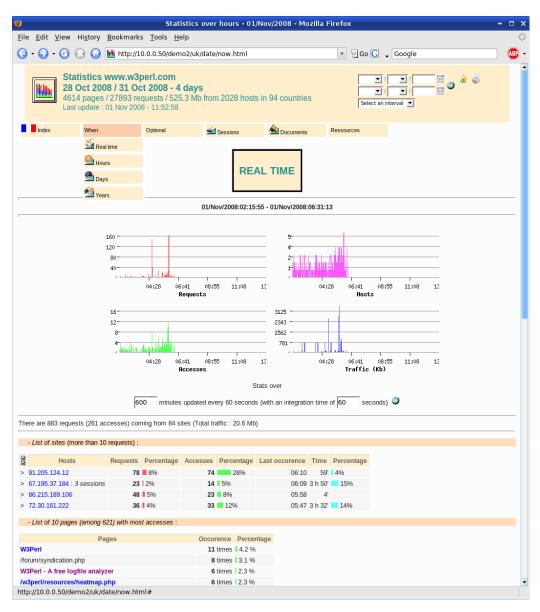  - Hits, Accesses, Traffic, Hosts…
- Graphics
  - Activity over the last hours
  - Length and integration can be changed in real time
- Show latest :
  - hosts / pages / scripts / files / traffic
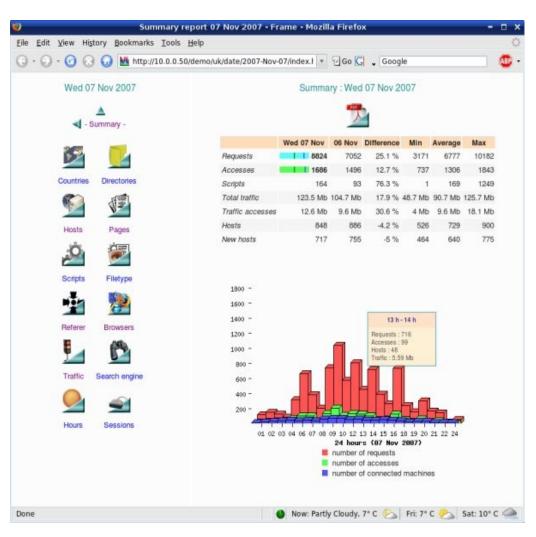- AJAX
  - Click on tables to sort them
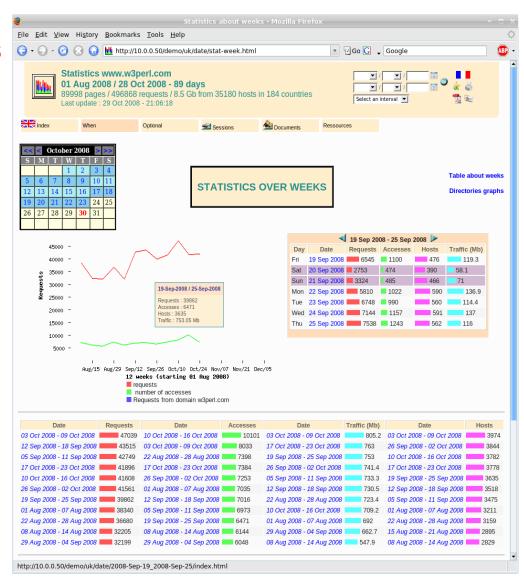  - Click on hosts to view page read

# For each day

- Summary
  - Hits, Accesses, Traffic, Hosts…

- Graphic
  - Activity over 24 hours
  - Show peaks

- Complete reports
  - Countries / Directories / Hosts / Pages / Filetype / Traffic / Session / Referer / Browsers / Scripts / …

- Can be send via email
  - Html report
  - PDF as an attachment

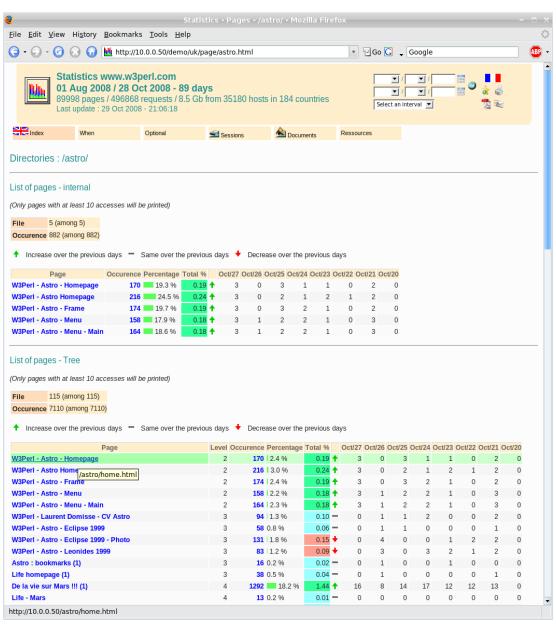# Weekly / Monthly / Yearly

- Same template as daily reports
  - Hits, Accesses, Traffic, Hosts…
- Graphics
  - Hits / Accesses versus time
  - Hosts versus time
- Full list
  - Not restricted to the top ten
- External / Domain split
- Interactive graphics
  - Link to others time stats
  - Display popup
- Can be send via email
  - Html report
  - PDF as an attachment

# Pages stats
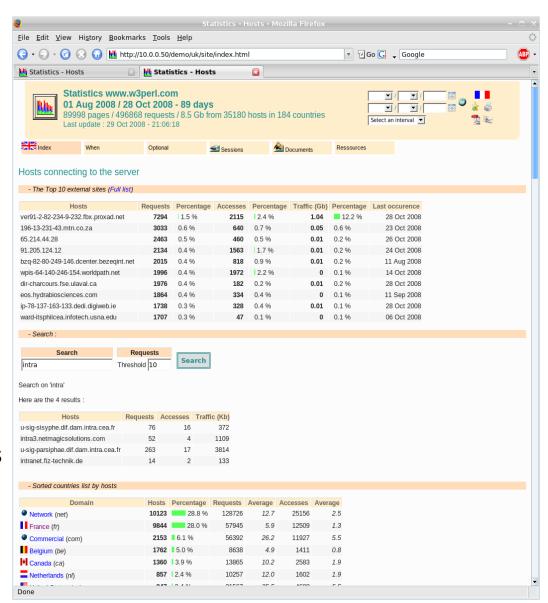
- **Based on file extension**
  - Html, php …
- **Most popular pages**
  - Full list available
- **Detailed pages**
  - View hosts/date/occurrence
- **Search**
  - To find easily a target
  - To compare a group
  - Match URL or page's title
- **Link to hosts**
  - View hosts for each page
- **Over the last few days**
  - Graphic activity
  - Increase / Decrease
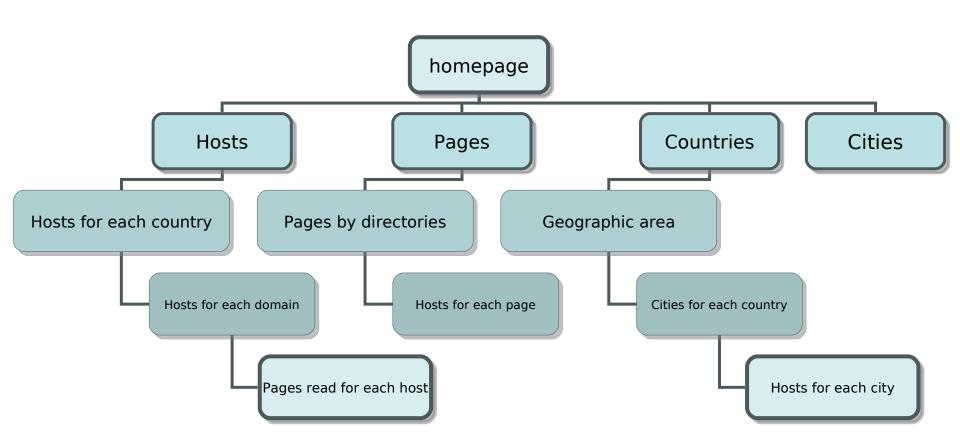- **Sorted by directory**

# Hosts stats

- Can exclude robots
  - Based on a file
- Most popular hosts
  - Full list available
- For each host
  - List of pages red and when
  - Hits / Accesses / Traffic
- Search
  - To find easily a target
  - To compare a group
- Explore by country domain
  - Country => Domain => Hosts => Host's pages

# Navigation



```
                    homepage
        ┌──────────────┼──────────────┬──────────────┐
      Hosts          Pages        Countries        Cities
        │              │              │
 Hosts for each    Pages by      Geographic area
   country         directories
        │              │              │
 Hosts for each   Hosts for      Cities for each
   domain         each page        country
        │                             │
 Pages read for                  Hosts for each
   each host                        city
```
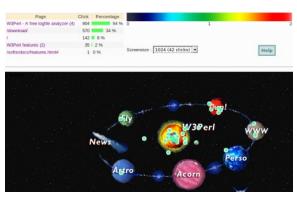
# Others

- **Using plug-in**
  - GeoIP to get country stats from IP logfile
  - GeoLiteCity to get cities stats
  - When adding a javascript :
    - Heat map
    - Screen resolution / Color depth / Java support
- **Directories**
  - Hits / Accesses / Traffic for each path level
- **Filetype**
  - Check traffic by file extension (find easily unsolicited files)
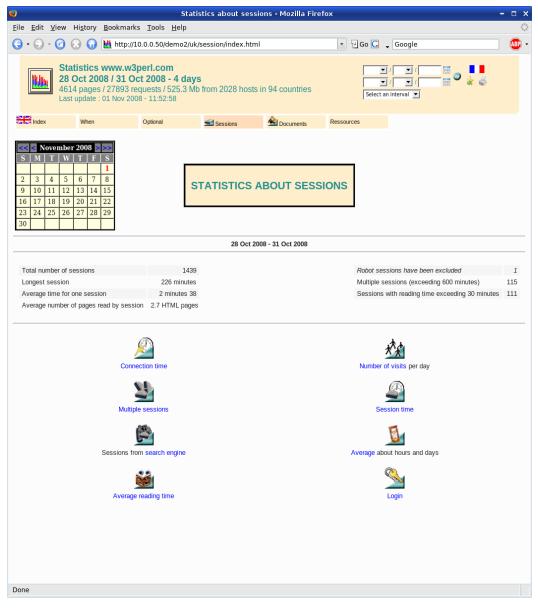
- **RSS**
  - Who is reading your flux
- **Status code**
  - Check client / server errors
- **Search engine**
  - Which keyword has been used to reach your website ?
- **Referer**
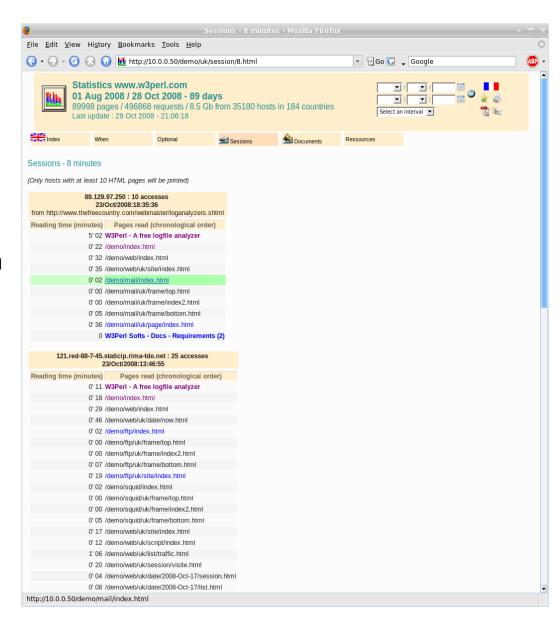  - Where does people come from ?

# Session problem

- Internet is an endless protocol
  - No delimiter
- Host = Visits ?
  - Same IP shared by different people
    - Proxy
    - NAT
    - Computer room
  - Same person with different IP
    - DHCP lease
    - Use different computer (home/work)
- Setting a Timeout
  - Usually 20 minutes without requests
- Session stats
  - Track user's path

# Follow user's path

- **For each session**
  - Host / Date / Referer / Accesses
  - View pages / Time spent
  - First / Ending page
- **Fidelity**
  - How many session for each host
  - How long they stay
- **Session length**
  - Histogram
- **Page reading time**
  - Most popular page (avoid page jump)
- **Session hourly and daily histogram**

# Squid / FTP / Mail

Add some extra reports

- Squid
  - Native and CLF format
  - Basic proxy stats (TCP status code, Traffic, elapsed time)
  - Destination sorted by countries
  - Users stats if authentified users found
- FTP
  - ProFTP, xferlog format
  - Users stats (files downloaded, traffic, from, date)
  - External / Local Top ten users
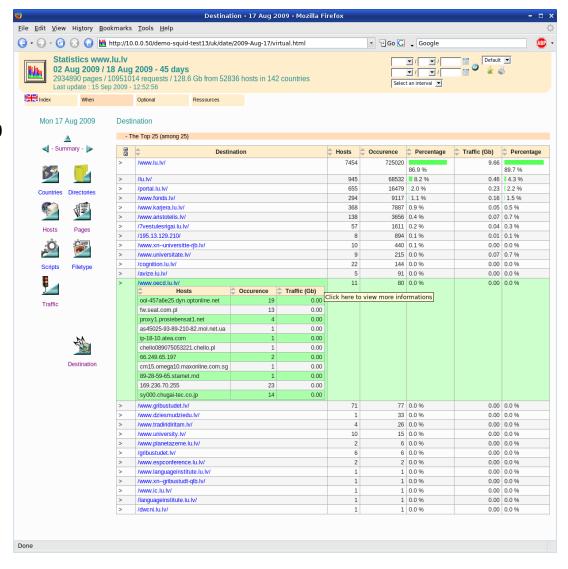  - Incoming or Outgoing reports
- Mail
  - Can read Exim / Sendmail / Postfix
  - Emails list per users / domain / countries
  - Local / External server

# Intranet

- Can map users to IP
  - Fill the resolv_users.csv
- Local domain
  - Regular expression to map domain IP
  - Parameter to exclude external hosts
- Works with secure server

Squid report showing :

- Daily destination : which websites have been seen

- List of hosts for each destination : occurrence and traffic

# Futur developements

- More logfile format support
  - PureFTP, Domino, Firewall …
- Web service
- More AJAX tools
  - Dynamic graphs
- More administration tools
  - Referer spammer management
- More tools for data-mining
- More mail reports
  - spam detection, bounced emails, delay
- Improve Windows installer
- Increase processing speed

# Others packages

## Logfile analysis

- AWStats
  - Widely used
  - No user's path track
- Analog
  - Fastest in the world
  - Lack of features
- Webalizer
  - Only small updates from 2002
  - Fast
- AWFFull
  - Webalizer fork
- Visitors
  - Fast
  - Command line based

## Page tagging

- Piwik
  - Require a database
  - Javascript code to include
- BBClone
  - PHP tag to include
- CrawTrack
  - Crawler stats
  - Require a database
- Google Analytics
  - No access to data
  - Popular

# That's all

- W3Perl homepage
  - http://www.w3perl.com
- New release every 3 months
- Demo : http://www.w3perl.com/demo/
  - Watch by yourself
- Feedbacks welcome
- Help available
  - FAQ
  - Forum
- Mailing list
  - Only new release update
- Comments ?
  - Send them to domisse@w3perl.com