# Millennium Web Logs

Web Wanderers slowing down your system?

GAELIC Show 'n Tell Workshop
held at the National Library, Pretoria
15 July 2009

Anette Lessing
Library Services, University of Pretoria

# Agenda

- Why analyze logs?

- Logs from Millennium server

  – Setup and retrieving logs

  – Tools to analyze logs

- Robots, Crawlers and Spiders

  – Control by Systems Administrator

  – Robot.txt file on server

# Why analyze logs?

- Support web development
  - Pages most used
  - What scopes are patrons searching
  - What browsers are used by clients
  - Where traffic is coming from

- Control over who is accessing your WebPAC
  - Robots (Web wanderers, Crawlers or Spiders used by search engines to index web content)

# Set-up and retrieving Millennium logs

- Web server logs available as from Release 2007
  - Must run Millennium Rel. 2007 behind Apache WebServer
    (Innovative does not provide a log analyzer)

**Setup to access logs**

➢ Login Manager

➢ Web Master

➢ Activate Web Server Logs

# Retrieving Millennium logs



In Web Master Mode -  Web Server Logs

# Retrieving Millennium logs



Download files and save to your PC

# What the logs look like



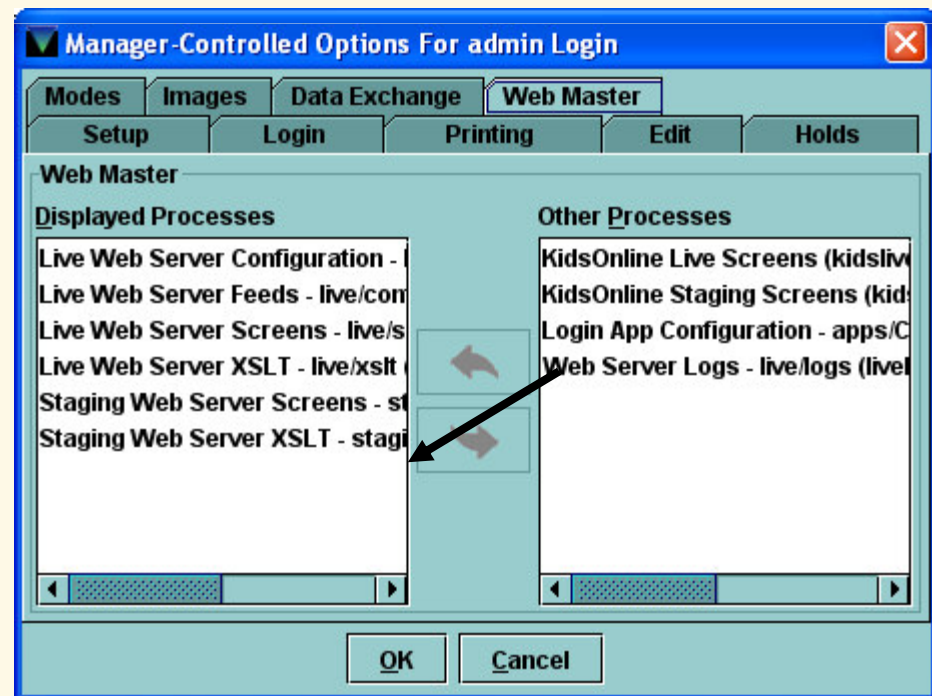Window title: access_log.2009-05-01 - Quick View Plus

Menu: File  Edit  View  Document  Window  Help

Log content excerpt:

80 127.0.0.1 - - [01/May/2009
80 190.82.23.16 - - [01/May/2
www.sciencedirect.com.innop
"http://innopac.up.ac.za/valid
www.sciencedirect.com.innop
es-ES; rv:1.9.0.10) Gecko/20
443 41.15.155.224 - - [01/Ma
HTTP/1.1" 200 4755 "https://
(compatible; MSIE 7.0; Wind
.NET CLR 1.1.4322)" 21739
80 190.82.23.16 - - [01/May/2
www.sciencedirect.com.innop
"http://innopac.up.ac.za/valid
www.sciencedirect.com.innop
es-ES; rv:1.9.0.10) Gecko/20
443 41.15.155.224 - - [01/Ma
/search~S9/?searchtype=X&
&SORT=D&extended=0&SUBMIT=Search HTTP/1.1" 200 7814
"https://innopac.up.ac.za/patroninfo~S9/1128047/getpsearches" "Mozilla/4.0 (compatible; MSIE 7.0;
Windows NT 6.0; GTB6; SLCC1; .NET CLR 2.0.50727; .NET CLR 3.0.04506; .NET CLR 1.1.4322)"
151532
443 41.15.155.224 - - [01/May/2009:02:01:47 +0200] "GET /forward/http://0-
amazon.com.innopac.up.ac.za/images/P/0419201408.01%20.00TLZZZZ HTTP/1.1" 404 124
"https://innopac.up.ac.za/search~S9/?searchtype=X&searcharg=eurocode+2&searchscope=9&sortdr
opdown=-&SORT=D&extended=0&SUBMIT=Search" "Mozilla/4.0 (compatible; MSIE 7.0; Windows
NT 6.0; GTB6; SLCC1; .NET CLR 2.0.50727; .NET CLR 3.0.04506; .NET CLR 1.1.4322)" 12627
443 41.15.155.224 - - [01/May/2009:02:01:47 +0200] "GET /forward/http://0-
amazon.com.innopac.up.ac.za/images/P/0580258238.01%20.00TLZZZZ HTTP/1.1" 404 124
"https://innopac.up.ac.za/search~S9/?searchtype=X&searcharg=eurocode+2&searchscope=9&sortdr
opdown=-&SORT=D&extended=0&SUBMIT=Search" "Mozilla/4.0 (compatible; MSIE 7.0; Windows
NT 6.0; GTB6; SLCC1; .NET CLR 2.0.50727; .NET CLR 3.0.04506; .NET CLR 1.1.4322)" 12806

Status bar: Unknown File Type | The contents of access_log.2009-05-01

Overlay text box:

**Use a log analyzer to interpret eg.**

**Webalizer (only analyze one log at a time)**

**or**

**Web Expert  Lite**

"*Web Log Expert Lite is the most straightforward, and will accept wildcards to analyze all the logs in a directory  o you can have it by month....*" )
**From IUG Listserv**

# Or go the Google Analytics way

- Reports than can be retrieved
  - Traffic sources
  - Visitors
  - Map overlay
  - Site usage
  - What browsers were used
  - And more (can set various filters)

Visits
1  6,299

# Google Analytics set-up

# Google Analytics Reports
## (Trial log for restricted period)

# Google Analytics Reports
(Trial log for restricted period)

# Google Analytics Reports

(Trial log for restricted period)

# Robots

- Known as Web Wanderers, Crawlers or Spiders

- Used by search engines to <u>index web content</u>

- Crawlers try to <u>follow every link</u> embedded in catalog pages

- Used by spammers to <u>scan for email addresses</u>

- Increase load on OPAC and <u>slow system</u> down

# List non-local access attempts allowed

Limit NETWORK access

**Are there many postings for the *SAME*
IP address in very short time span?**

| | DATE | TIME | REMOTE IP ADDRESS | SERVICE NAME |
|---|---|---|---|---|
| 0001 > | 03-07-09 | 08:49:03 | 38.99.13.119 | http |
| 0002 > | 03-07-09 | 07:26:50 | 38.99.44.101 | http |
| 0003 > | 03-07-09 | 07:26:50 | 38.99.44.104 | http |
| 0004 > | 03-07-09 | 07:25:50 | 38.99.44.101 | http |
| 0005 > | 03-07-09 | 07:25:50 | 38.99.44.104 | http |
| 0006 > | 03-07-09 | 07:24:50 | 38.99.44.101 | http |
| 0007 > | 03-07-09 | 07:24:50 | 38.99.44.104 | http |
| 0008 > | 03-07-09 | 07:23:03 | 72.30.87.98 | http |
| 0009 > | 03-07-09 | 07:21:57 | 38.99.44.102 | http |
| 0010 > | 03-07-09 | 07:21:12 | 38.99.44.104 | http |
| 0011 > | 03-07-09 | 07:21:09 | 38.99.44.101 | http |
| 0012 > | 03-07-09 | 07:19:57 | 38.99.44.102 | http (2) |
| 0013 > | 03-07-09 | 07:16:51 | 38.99.44.105 | http |
| 0014 > | 03-07-09 | 07:16:46 | 38.99.13.117 | http |
| 0015 > | 03-07-09 | 07:16:13 | 38.99.44.102 | http |
| 0016 > | 03-07-09 | 07:15:52 | 38.99.44.103 | http |

```
F > FORWARD          S > SORT                 Q > QUIT
J > JUMP             N > Display host NAME     P > PRINT
C > CLEAR log file   T > TOTAL by date
```

# Look up the host info

# Restrict Access on Innopac

Limit NETWORK access

# Robots.txt file on server

- http://<your server/>robots.txt
- Allow legitimate search engines to index main page of catalog (set by III)

# Specific segment for Google Scholar

```
# For the WebBridge Google Scholar Extension. Allows googlebot_IA to crawl
# /screens
User-agent: Googlebot-IA
Disallow: /acquire
Disallow: /airpac
Disallow: /airwkst
Disallow: /articles
Disallow: /availlim
Disallow: /bookill
Disallow: /bookit
Disallow: /circhistlim
Disallow: /circpix
Disallow: /cisti_order
Disallow: /clearhist
Disallow: /documents
Disallow: /donate
Disallow: /extlang
Disallow: /feeds
Disallow: /ftlist
Disallow: /goto
Disallow: /iii
Disallow: /ill
Disallow: /illframe
```

## Allow crawling of both / and /screens

**From: Mark Welge, Robots, Crawlers and spiders…Oh My! Automated Searches and your WebPAC, IUG Anaheim, 17-20 May 2009**

Innovative *interfaces*

IUG

# References

- http://csdirect.iii.com/documentation/weblogs.shtml   (CSDirect search: web logs)
- (Alan Dyck, WebPAC product manager Web Access Log Analysis IUG Anaheim, 17-20 May, 2009
- Mark Welge, Robots, Crawlers and spiders…Oh My! Automated Searches and your WebPAC, IUG Anaheim, 17-20 May 2009.

**Thank you**



From: Mark Welge, Robots, Crawlers and spiders…Oh My! Automated Searches and your WebPAC, IUG Anaheim, 17-20 May 2009