

WebVoyáge with a Wrapper



Michael Doran, Systems Librarian
doran@uta.edu

Ex Libris Users of North America (ELUNA) Meeting
Long Beach, CA - Session 48.2 - July 31, 2008

Once upon a time...



"If you will just configure me right,
I will turn into a handsome OPAC."



Michael Doran, Systems Librarian

doran@uta.edu

What is a handsome OPAC?

- Aesthetically handsome
- Functionally handsome
- An OPAC is "handsome" if it
 - is simple to use
 - is intuitive to use
 - makes it easy to find stuff
- "Only librarians like to search, everybody else likes to find."
– Roy Tennant



Michael Doran, Systems Librarian

doran@uta.edu

Simple searches



Michael Doran, Systems Librarian

doran@uta.edu

WebVoyáge simple search

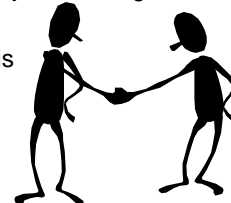


Michael Doran, Systems Librarian

doran@uta.edu

The secret handshakes

- last name, first name for author searches
- no initial articles for title searches
- Library of Congress subject headings
- Boolean operators
- what an index browse is



Michael Doran, Systems Librarian

doran@uta.edu

WebVoyage simple search (after)



- keyword anywhere search
- words within quotes are treated as a phrase
- other words are automatically Boolean ANDed*
- relevancy ranked results*

Hmmm... that's **Google**-like
in functionality



Michael Doran, Systems Librarian

doran@uta.edu

Code



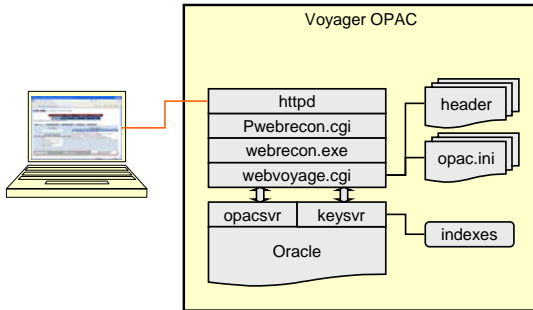
```
<HTML>
<HEAD>
<TITLE>WebVoyage</TITLE>
<META data>
</META data>
</HEAD>
<BODY>
<P>
Yada, yada, yada...
</P>
</P>
<FORM ACTION="/webrecon.cgi">
<INPUT TYPE="text">
<INPUT TYPE="submit">
</FORM>
<P>
More yada, yada, yada...
</P>
</BODY>
</HTML>
```



Michael Doran, Systems Librarian

doran@uta.edu

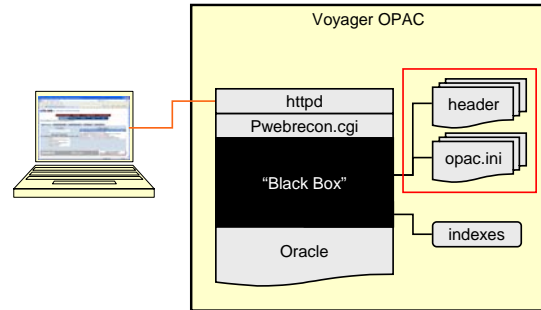
WebVoyage server-side back end



Michael Doran, Systems Librarian

doran@uta.edu

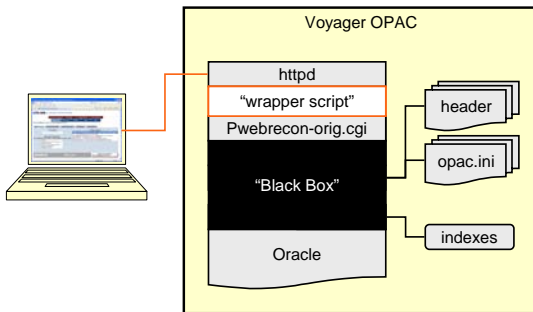
WebVoyage is a "black box"



Michael Doran, Systems Librarian

doran@uta.edu

They call it a wrapper



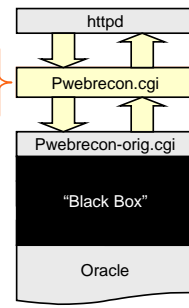
Michael Doran, Systems Librarian

doran@uta.edu

Basic wrapper script

```
#!/usr/bin/perl
$data_stream = `./Pwebrecon-orig.cgi`;
print $data_stream;
exit;
```

```
<HTML>
<HEAD>
<TITLE>WebVoyage</TITLE>
<META data>
</META data>
</HEAD>
<BODY>
<P>
Yada, yada, yada...
</P>
</P>
<FORM ACTION="/webrecon.cgi">
<INPUT TYPE="text">
<INPUT TYPE="submit">
</FORM>
<P>
More yada, yada, yada...
</P>
</BODY>
</HTML>
```



Michael Doran, Systems Librarian

doran@uta.edu

Do your thing to that datastream

- aka "screen scraping" 

"A technique in which a computer program extracts data from the display output of another program. The key element that distinguishes screen scraping from regular parsing is that the output being scraped was intended for final display to a human user, rather than as input to another program, and is therefore usually neither documented nor structured for convenient parsing." [from Wikipedia]

- text wrangling 

- add text
- delete text
- rearrange text



Michael Doran, Systems Librarian

doran@uta.edu

Example – adding text

- Voyager's "header.htm" file
 - is inserted after the <body> tag
 - okay for display tags, but not for others
- Wrapper script can insert elements within the <head> tag
 - metadata
 - JavaScript
 - CSS



Michael Doran, Systems Librarian

doran@uta.edu

Example – adding text

```
#!/usr/bin/perl
$data_stream = `./Pwebrecon-orig.cgi`;
$meta_code =
    qw(<link rel="stylesheet" type="text/css" href="/css/my.css">);
$data_stream -- s#</HEAD>#$meta_code</HEAD>#i;
print $data_stream;
exit;
```

```
<HTML>
<HEAD>
<meta http-equiv="Content-Type" Content="text/html;charset=UTF-8">
<TITLE>Library Catalog - University of Texas at Arlington</TITLE>
<link rel="stylesheet" type="text/css" href="/css/my.css">
</HEAD>
<BODY onLoad="document.querybox.Search_Arg.focus()"...
```

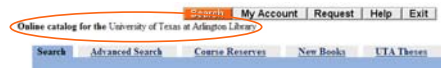


Michael Doran, Systems Librarian

doran@uta.edu

Example – removing text

```
#!/usr/bin/perl
$data_stream = `./Pwebrecon-orig.cgi`;
$data_stream -- s#<TD><STRONG>.*?</STRONG>University of Texas at
Arlington Library</TD>#i;
print $data_stream;
exit;
```

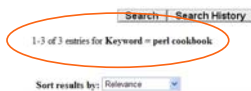


Michael Doran, Systems Librarian

doran@uta.edu

Example – rearranging text

```
#!/usr/bin/perl
$data_stream = `./Pwebrecon-orig.cgi`;
$data_stream -- s#Search request: (.*)</TD>.*Results: (.*)
entries.#<br />/$2 entries for <b>$1</b>#s;
print $data_stream;
exit;
```

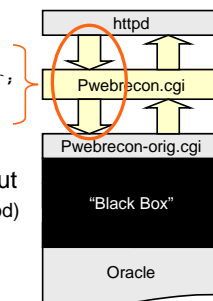


Michael Doran, Systems Librarian

doran@uta.edu

Wrapper script redux

```
#!/usr/bin/perl
$data_stream = `./Pwebrecon-orig.cgi`;
print $data_stream;
exit;
```



- Read and parse form input
 - QUERY_STRING (get method)
 - STDIN (post method)



Michael Doran, Systems Librarian

doran@uta.edu

Truncation adaptation



Michael Doran, Systems Librarian

doran@uta.edu

Incoming data



```
#!/usr/bin/perl
$data_stream = `./Pwebrecon-orig.cgi`;
print $data_stream;
exit;
```

QUERY STRING

```
Search_Arg=samuel+clem*&Search_Code=GKEY%5E*&
PID=D3hcrVvigATy0bZXTmXMK61ori&
SEQ=20070527145935&CNT=50&HIST=1
```



Michael Doran, Systems Librarian

doran@uta.edu

Incoming data



```
#!/usr/bin/perl
ReadParse();
$data_stream = GetOrigDataStream();
sub GetOrigDataStream {
    $data_stream = `./Pwebrecon-orig.cgi`;
    return $data_stream;
}
print $data_stream;
exit;
```



Michael Doran, Systems Librarian

doran@uta.edu

Example – truncation adaptation

```
#!/usr/bin/perl
ReadParse();
$data_stream = GetOrigDataStream();
sub GetOrigDataStream {
    $search_arg = $formdata{'Search_Arg'};
    $search_arg =~ s/\*/?/;
    if ($ENV{'QUERY_STRING'}) {
        $ENV{'QUERY_STRING'} =~ s/Arg=.*?&Search/Arg=$search_arg&Search/;
    }
    $data_stream = `./Pwebrecon-orig.cgi`;
    return $data_stream;
}
print $data_stream;
exit;
```



Michael Doran, Systems Librarian

doran@uta.edu

Example – truncation adaptation

Search | Search History | **Results List** | My Account | Request | Help | Exit

1-50 of 169 entries for Keyword = samuel clem*

« Previous 1 51 101 151 Next »

Sort results by: [Relevance] [v] Refine Search Results

#	Title	Author	Date
1	Travels of Paul Fildred Wilson, and the comedi: These enterwonderous trave, by Mark Twain (Samuel L. Clemens)	Clemens, Samuel Langhorne, 1835-1910	1894
2	All Sin, a dramatic work, by Mark Twain (Samuel L. Clemens), Edited by Frederick Anderson	Clemens, Samuel Langhorne, 1835-1910	1961
3	Innocent relatives, Mobile, Browne, and Mark Twain in the Holy Land, by Frankie Walker	Walker, Frankie Dickerson, 1900-	1974
4	Documentation, by S. C. Bradford	Bradford, Samuel Clemens, 1878-1948.	1950
5	Favorite works of Mark Twain (Samuel L. Clemens)	Clemens, Samuel Langhorne, 1835-1910	1939



Michael Doran, Systems Librarian

doran@uta.edu

Other input data munging

- fix Voyager 6.x GKEY/TKEY/SKEY keyword "multiple spaces" no hits bug (Support Web incident #131344)
`$search_arg =~ s/ / /g;`
- deal with "right single quotation mark" vs. "apostrophe" in search input issue
- allow for ISBNs with dashes*



* (combined output/input) data munging



Michael Doran, Systems Librarian

doran@uta.edu

Is a wrapper right for you?

- requires some programming expertise
- requires lots (and lots) of testing
 - test platform
 - ideally a Voyager test server
 - separate WebVoyage instance (a la preview server)
 - law of unintended consequences
- extra layer makes WebVoyage more brittle
 - more dependencies, e.g. with opac.ini
- upgrades more complicated



Michael Doran, Systems Librarian

doran@uta.edu

Getting started

- wrappers are programming-language-neutral, however...
- Perl is good
 - designed for text processing
 - robust regular expressions
 - is already on your system
 - example wrappers available
- it's fine to think big...
- ... but start small



Michael Doran, Systems Librarian

doran@uta.edu

Resources

- Michael Doran, University of Texas at Arlington
Presentation: "WebVoyage with a Wrapper"
Source code:
<http://rocky.uta.edu/doran/wrapper/>
- Ere Maijala, National Library of Finland
European EndUser 2006 presentation*:
"Enhancement scripts for WebVoyage OPAC"
(no longer on the web – Contact Ere Maijala)
Source code:
<http://www.nationallibrary.fi/libraries/linnea/pwebrecon2.html>



Michael Doran, Systems Librarian

doran@uta.edu

A small start

- copy original Pwebrecon.cgi
`cp -p Pwebrecon.cgi Pwebrecon-orig.cgi`
- create Pwebrecon.cgi wrapper template


```
#!/usr/bin/perl
$data_stream = `./Pwebrecon-orig.cgi`;
print $data_stream;
exit;
```
- add desired feature
- test

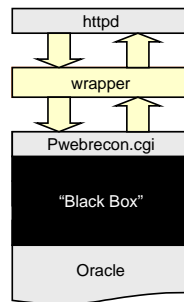


Michael Doran, Systems Librarian

doran@uta.edu

Will wrappers work with Tomcat?

- Is it **possible** to apply a wrapper to the Voyager 7 Tomcat WebVoyage?
- Is it **desirable** to apply a wrapper to the Voyager 7 Tomcat WebVoyage?
 - V7 Tomcat XML files
 - V7 Cascading Style Sheets

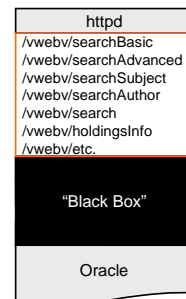


Michael Doran, Systems Librarian

doran@uta.edu

Will wrappers work with Tomcat?

- Is it **possible** to apply a wrapper to the Voyager 7 Tomcat WebVoyage?
 - No "Pwebrecon.cgi" under Tomcat WebVoyage (in fact no CGI period)
 - No one single controlling executable to wrap



Michael Doran, Systems Librarian

doran@uta.edu