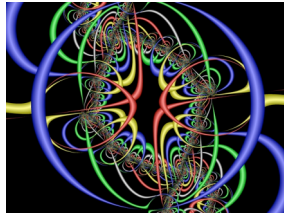# Advanced collection configuration

Course material prepared by

Greenstone Digital Library Project
University of Waikato, New Zealand

# Agenda

- ❖ Formatting
- ❖ Formatting extended metadata
- ❖ Changing metadata sets
- ❖ Collection configuration file
- ❖ GLI modes
- ❖ PDF documents
- ❖ PPT documents
- ❖ Exploding metadata databases
- ❖ Cross collection searching
- ❖ Translate text
- ❖ Full-text tagging
- ❖ Creating a CD-ROM

# Format Statements: Components

- ❖ HTML
- ❖ [Text] — *displays document text*
- ❖ [Title],[Howto] ... — *displays metadata*
- ❖ [link] … [/link] — *links to document*
- ❖ [srclink] … [/srclink] — *links to original file*
- ❖ [icon], [srcicon] — *page/book/bookshelf source icons*
- ❖ *If* and *Or* statements — *conditional processing*

# Format Statements

`{If}{test, if true, if false}`

- ❖ Test can be:

  | [metadata] | exists |
  | [metadata] eq 'value' | equals |
  | [metadata] ne 'value' | not equals |

- ❖ Examples

  `{If}{[ex.Source],(<i>[ex.Source]</i>)}`

  `{If}{[numleafdocs],[ex.Title],[dc.Creator]}`

  `{If}{[ex.FileFormat] eq PDF, [srclink] PDF document [/srclink]}`

# Format Statements

`{Or}{[metadata1],[metadata2],…`

- ❖ Chooses the first metadata that exists

- ❖ Last item can be plain text

- ❖ Examples

  `{Or}{[dls.Title],[dc.Title],[ex.Title],Untitled}`

  `{Or}{[ex.thumbicon],[ex.srcicon]}`

# HTML page

```
<html>
  <head>
    <title>The page title</title>
  </head>
  <body formatting-attributes>
    Page content goes here
  </body>
</html>
```

## HTML elements

`<b>Bold</b>`, `<i>italics</i>`, `<u>underline</u>`

`<br/> a line break`

`<p>a paragraph</p>`

`<table><tr><td>cell content</td></tr></table>`

`<a href="link address">link text</a>`
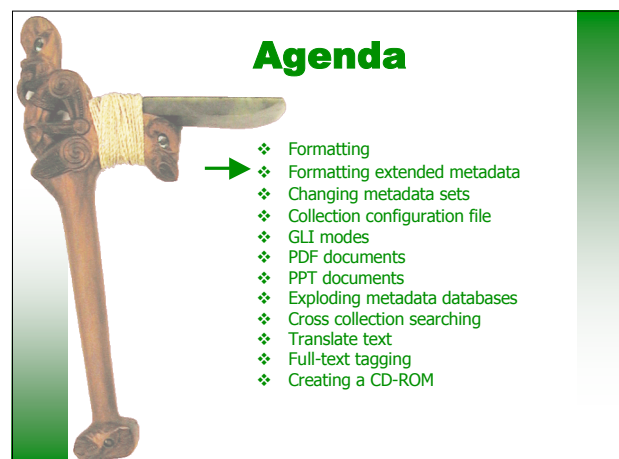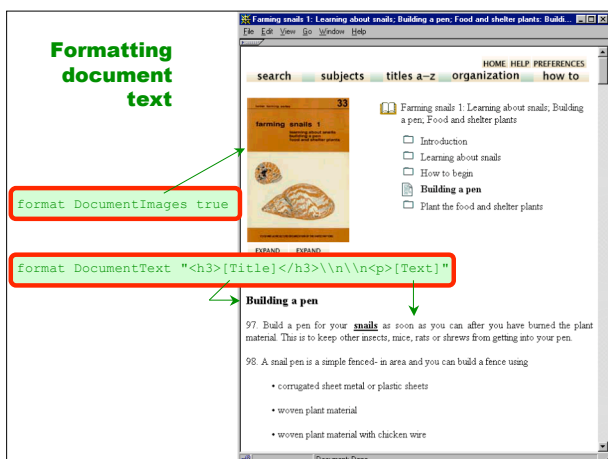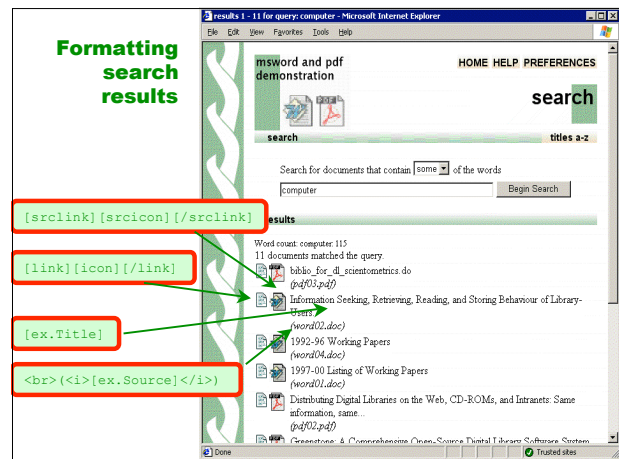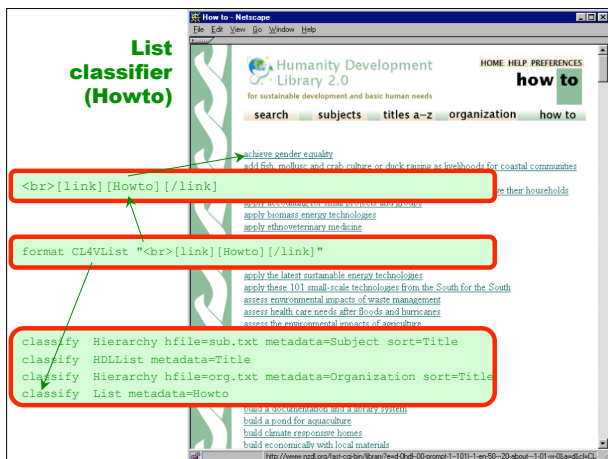
`<img src="path to image"/>`

## Format statements

❖ Defaults provided to give reasonable results for most cases
❖ Interpreted at run-time – no need to rebuild
❖ More information: FAQ

http://www.greenstone.org/cgi-bin/library?a=p&p=faqcustomize#customizeformat

### List classifier (Howto)



```
<br>[link][Howto][/link]

format CL4VList "<br>[link][Howto][/link]"

classify  Hierarchy hfile=sub.txt metadata=Subject sort=Title
classify  HDLLList metadata=Title
classify  Hierarchy hfile=org.txt metadata=Organization sort=Title
classify  List metadata=Howto
```

### Formatting search results



```
[srclink][srcicon][/srclink]

[link][icon][/link]

[ex.Title]

<br>(<i>[ex.Source]</i>)
```

### Formatting document text



```
format DocumentImages true

format DocumentText "<h3>[Title]</h3>\\n\\n<p>[Text]"
```

## Agenda

❖ Formatting
➔ ❖ Formatting extended metadata
❖ Changing metadata sets
❖ Collection configuration file
❖ GLI modes
❖ PDF documents
❖ PPT documents
❖ Exploding metadata databases
❖ Cross collection searching
❖ Translate text
❖ Full-text tagging
❖ Creating a CD-ROM

## Extended metadata

❖ Basic metadata:
    [Title], [Source]                                    (extracted)
    [dc.Subject], [dls.Organization]          (manually assigned)
❖ Extended metadata:
    [parent:Title]                              (Title of parent section)
    [parent(Top):dc.Title]      (dc.Title of document, i.e. top section)
    [sibling:Subject]                    (all Subjects of current section)
    [child:Author]                        (Author of first child section)
    [child(All):Author]                      (Author of all children)
❖ Formatting:
    Sibling(All' and '):Subject
                          (between quotes can specify a separator)

## Format statements: Extended metadata

Snail Farming

        Subject: Agriculture
        Subject: Farming
        Subject: Cuisine
   1. Introduction
     1.1 Snails are good to eat
     1.2 What is snail farming?
   2. Getting started
      2.1 How to prepare the pens

## Format statements: Extended metadata

Snail Farming
    Subject: Agriculture
    Subject: Farming
    Subject: Cuisine
1. Introduction
   1.1 Snails are good to eat
   1.2 What is snail farming?
2. Getting started
   2.1 How to prepare the pens

| | |
|---|---|
| Subject | Agriculture |
| sibling:Subject | Agriculture, Farming, Cuisine |
| sibling(last):Subject | Cuisine |
| Sibling(All' and '):Subject | Agriculture and Farming and Cuisine |
| child:Title | Introduction, Getting Started |
| child(2):Title | Getting started |

## Format statements: Extended metadata

Snail Farming
    Subject: Agriculture
    Subject: Farming
    Subject: Cuisine
1. Introduction
   1.1 Snails are good to eat
   1.2 What is snail farming?
2. Getting started
   2.1 How to prepare the pens

| | |
|---|---|
| parent:Title | Introduction |
| parent(Top):Title | Snail Farming |
| parent(All):Title | Snail Farming, Introduction |
| Parent(Top):Subject | Agriculture |
| Parent(Top):sibling:Subject | Agriculture, Farming, Cuisine |
| Parent(Top):sibling(All' and '):Title | Agriculture and Farming and Cuisine |

## Format statements: Combining If and Or

❖ {If}{test, if true, if false}
❖ {Or}{[metadata1],[metadata2]…}
❖ {Or} can have a conditional as the last element:
❖ {Or}{[Creator],[Editor],
        {If}{[FileFormat] eq "PDF",
        xxx,anonymous}}

## Format statements: Combining If and Or

❖ {If} can have another conditional at 'true' or 'false' position

❖ {If}{[numleafdocs],[Title],
   [dc.Title]{If}{[Date],: [Date]}
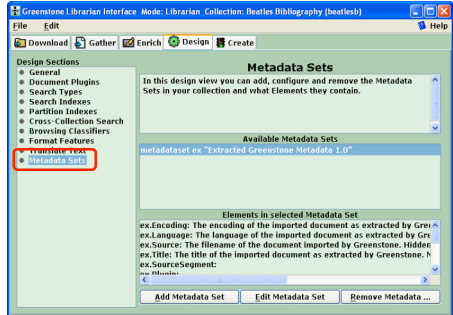   {If}{[Subject], ([Subject], unclassified)}}

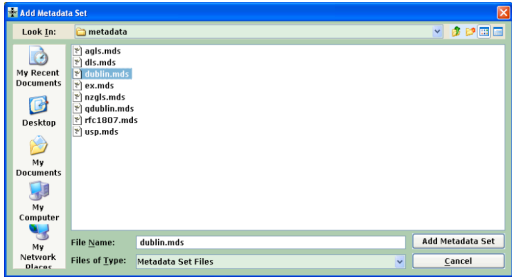Snail farming: 26 Jun 1998 (Small Animal Farming)

## Agenda

- ❖ Formatting
- ❖ Formatting extended metadata
- ➜ ❖ Changing metadata sets
- ❖ Collection configuration file
- ❖ GLI modes
- ❖ PDF documents
- ❖ PPT documents
- ❖ Exploding metadata databases
- ❖ Cross collection searching
- ❖ Translate text
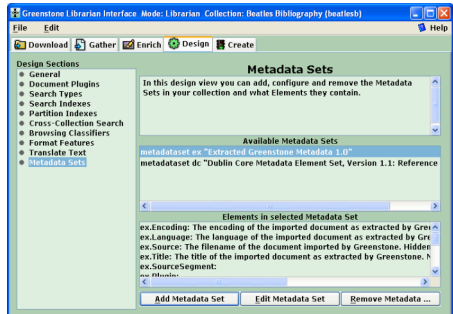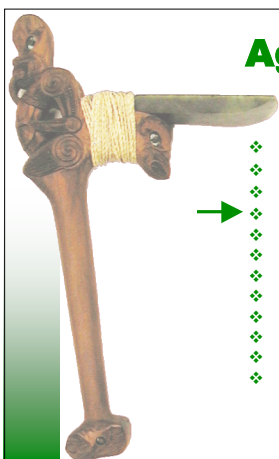- ❖ Full-text tagging
- ❖ Creating a CD-ROM

## Metadata Sets



## Add Metadata Set



## Metadata Sets



## Agenda

- ❖ Formatting
- ❖ Formatting extended metadata
- ❖ Changing metadata sets
- ➜ ❖ Collection configuration file
- ❖ GLI modes
- ❖ PDF documents
- ❖ PPT documents
- ❖ Exploding metadata databases
- ❖ Cross collection searching
- ❖ Translate text
- ❖ Full-text tagging
- ❖ Creating a CD-ROM

## Collection configuration file

- ❖name, icon, etc
- ❖description
- ❖email of creator
- ❖search indexes
- ❖plugins
- ❖classifiers

*how to format*

- ❖documents
- ❖query results
- ❖classifiers

## Documentation and Help

❖ *User's Guide* (user.pdf)

   Includes substantial sections on the GLI
   (Sections 3.1 and 3.2, 36 pp)

❖ Tooltips

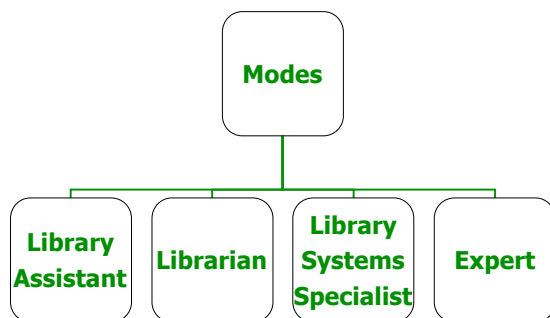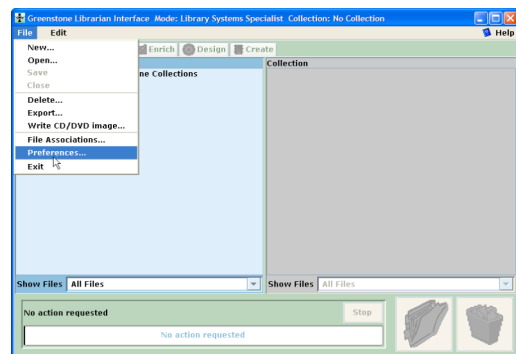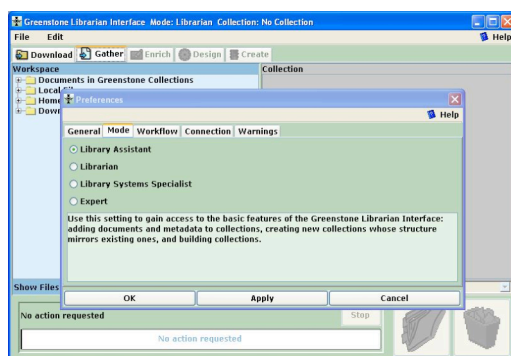❖ On-line help

## Agenda

❖ Formatting
❖ Formatting extended metadata
❖ Changing metadata sets
❖ Collection configuration file
➤ ❖ GLI modes
❖ PDF documents
❖ PPT documents
❖ Exploding metadata databases
❖ Cross collection searching
❖ Translate text
❖ Full-text tagging
❖ Creating a CD-ROM

## GLI Modes

•File->Preferences->Mode

Modes

Library Assistant

Librarian

Library Systems Specialist

Expert

## GLI Mode Setting



## GLI Mode Setting



## GLI Modes

❖ **Library Assistant**
  – Access to basic features: creating new collections;
    adding documents and metadata; building collections
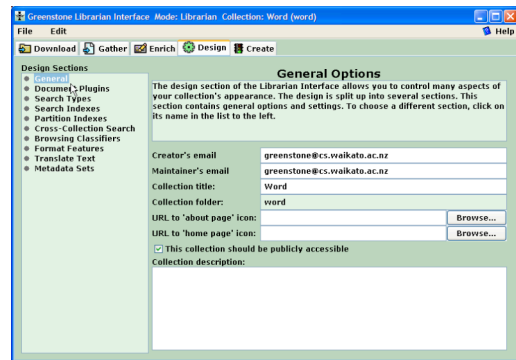  – No design functions are available

❖ **Librarian**
  – Basic features + Design

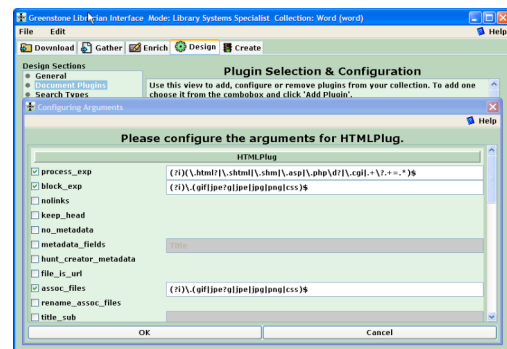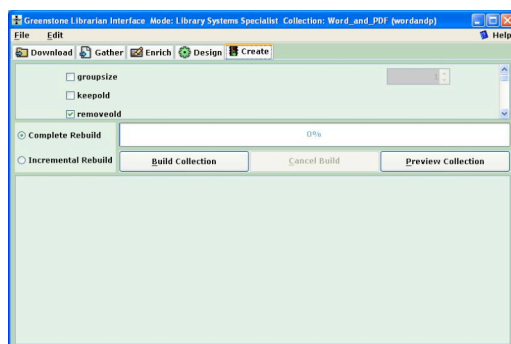## Library Assistant Mode



## Librarian Mode



## GLI Modes

❖ **Library Systems Specialist**
– Full use of GLI

– Formulate regular expressions to make use of formatting features
❖For example: HTML block expression
```
q^(?i)\.(gif|jpe?g|jpe|png|css)$^
```

– Partial options for import and building

## Library Systems Specialist Mode
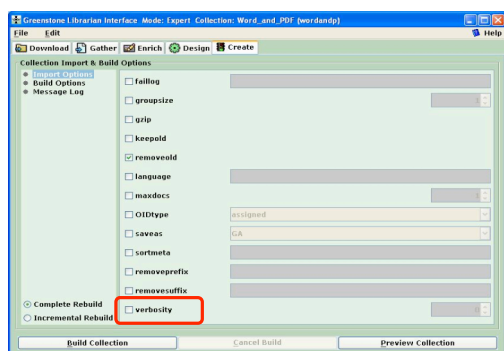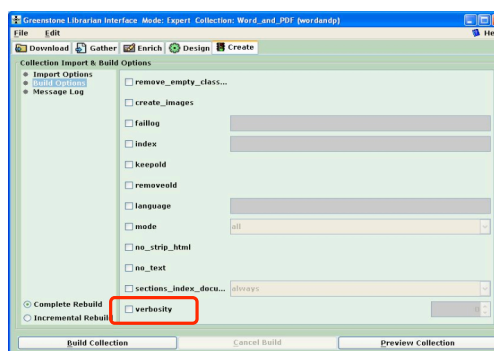


## Library Systems Specialist Mode



## GLI Modes

❖ **Expert**
– All features are enabled
– Recommended for experienced users
– Perform troubleshooting tasks
– Options to control the import and build processes
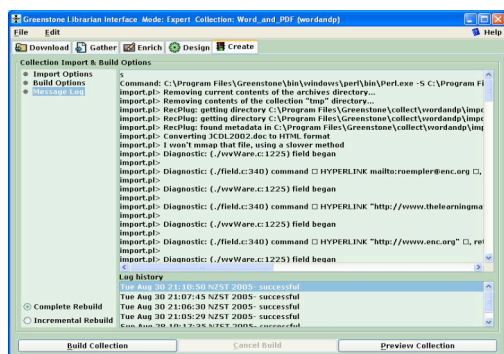– Shows the output from the processes
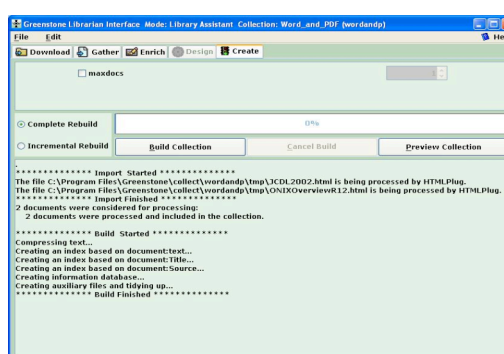
## Expert Mode—Import Options



## Expert Mode—Build Options



## Message Log—Expert Mode
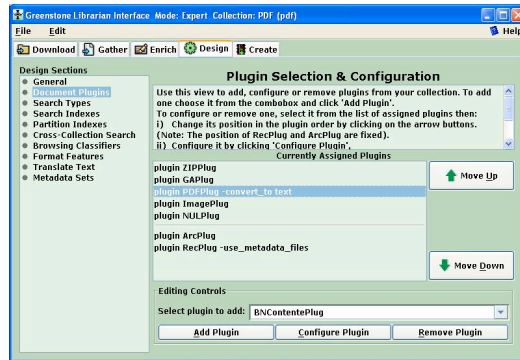


## Building Message—Others



## Agenda



- ❖ Formatting
- ❖ Formatting extended metadata
- ❖ Changing metadata sets
- ❖ Collection configuration file
- ❖ GLI modes
- → ❖ PDF documents
- ❖ PPT documents
- ❖ Exploding metadata databases
- ❖ Cross collection searching
- ❖ Translate text
- ❖ Full-text tagging
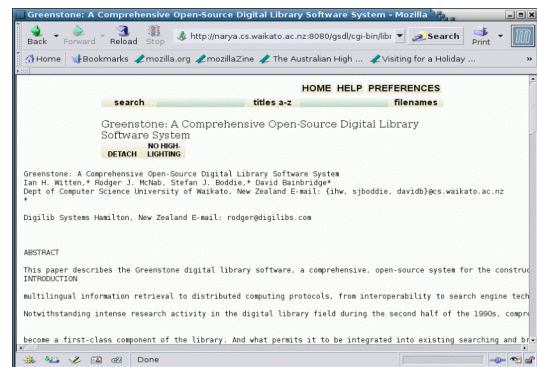- ❖ Creating a CD-ROM

## PDF Document

- ❖ PDF conversions in Greenstone
  1. Text only for Unix   system
  2. HTML
     - ❖ use_sections option
     - ❖ complex option
  3. Image
     - ❖ ImageMagick needs to be installed
     - ❖ For advanced conversions, GhostScript must be installed
     - ❖ Use of convert utility
     - ❖ Convert_to

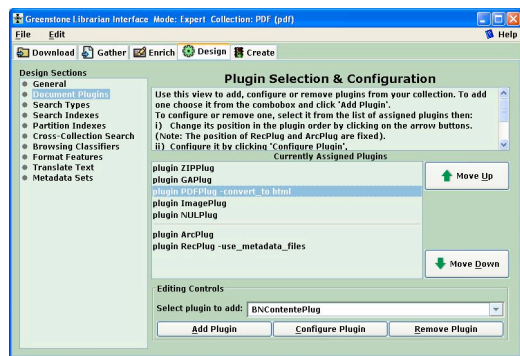       pagedimg_jpg
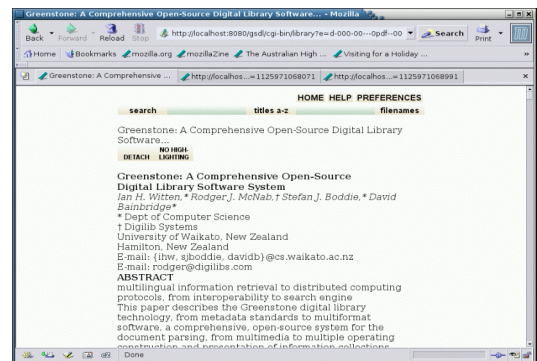       pagedimg_gif
       pagedimg_png
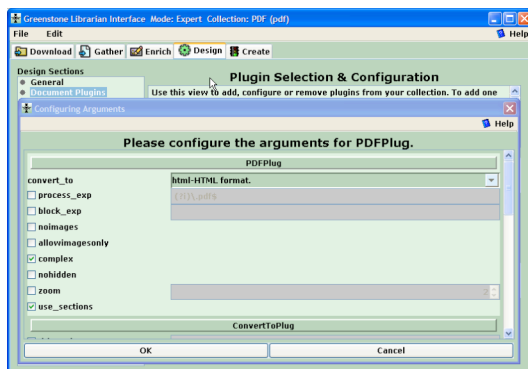
## PDF -> Text

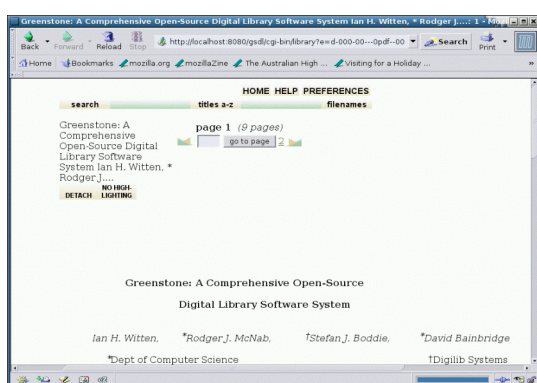## PDF: Text Document Display

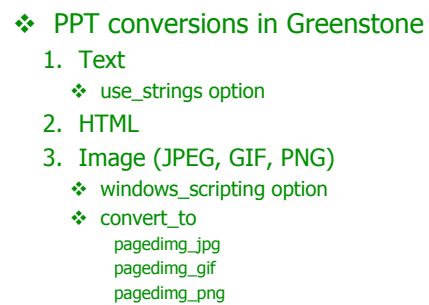## PDF -> HTML

## PDF: HTML Document Display 1

## PDF: use_sections
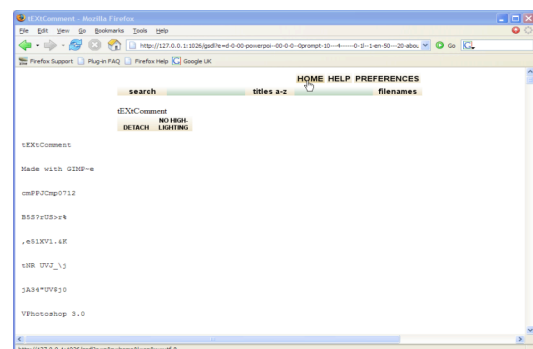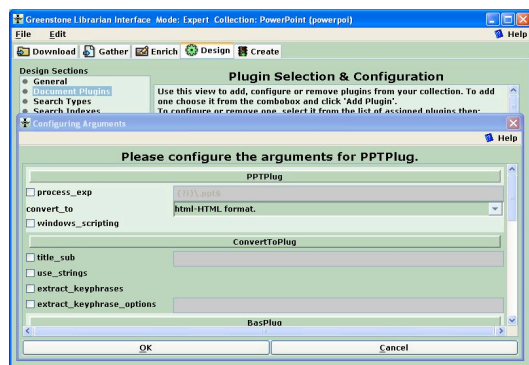
## PDF: HTML Document Display 2

## PDF -> Image



## PDF: Image Document Display



## Agenda



- ❖ Formatting
- ❖ Formatting extended metadata
- ❖ Changing metadata sets
- ❖ Collection configuration file
- ❖ GLI modes
- ❖ PDF documents
- ➤ ❖ PPT documents
- ❖ Exploding metadata databases
- ❖ Cross collection searching
- ❖ Translate text
- ❖ Full-text tagging
- ❖ Creating a CD-ROM

## PowerPoint Document

- ❖ PPT conversions in Greenstone
  1. Text
     - ❖ use_strings option
  2. HTML
  3. Image (JPEG, GIF, PNG)
     - ❖ windows_scripting option
     - ❖ convert_to
       - pagedimg_jpg
       - pagedimg_gif
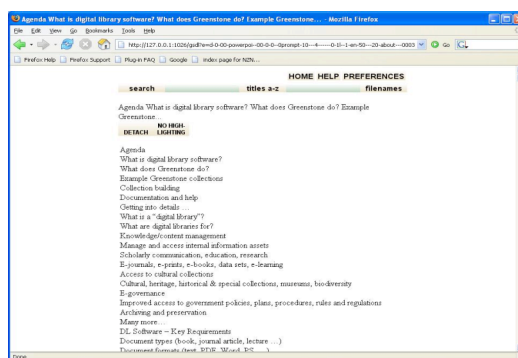       - pagedimg_png

## PPT -> Text


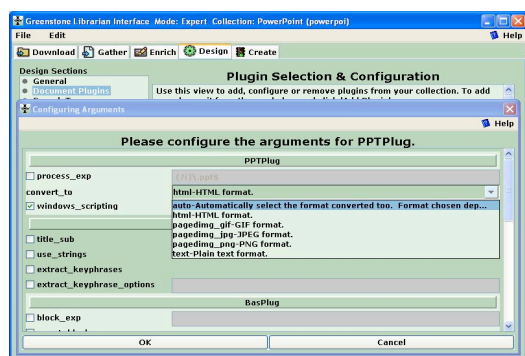
## PPT: Text Document Display
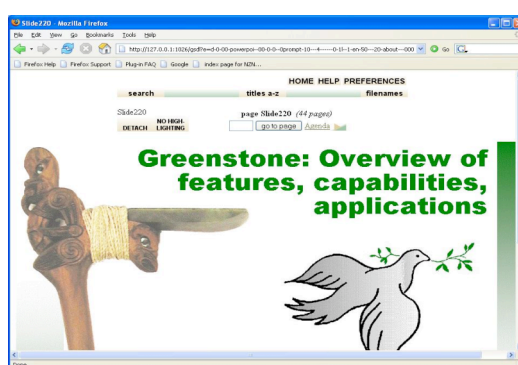
# PPT -> HTML



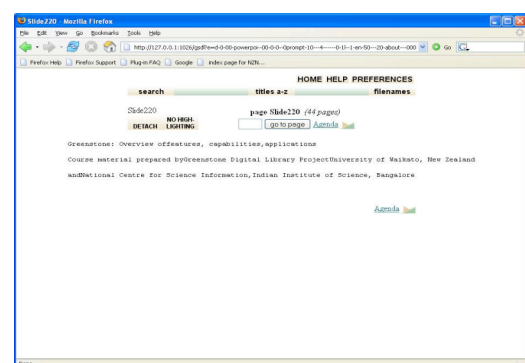# PPT: HTML Document Display



# PPT -> Image



# PPT Image: Image View



# PPT Image: Text View



# Agenda

- ❖ Formatting
- ❖ Formatting extended metadata
- ❖ Changing metadata sets
- ❖ Collection configuration file
- ❖ GLI modes
- ❖ PDF documents
- ❖ PPT documents
- ❖ Exploding metadata databases
- ❖ Cross collection searching
- ❖ Translate text
- ❖ Full-text tagging
- ❖ Creating a CD-ROM

## Exploding metadata databases

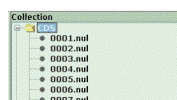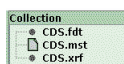❖ The GLI does not show metadata extracted from bibliographic files: CDS/ISIS, MARC etc.



## Why not?

❖ The GLI treats these files just like any other: it shows extracted and assigned metadata, but not the file's contents
  – Double-click the file to open it in the normal editor for this file type eg. WinISIS for CDS/ISIS files
❖ In rare cases you might want to convert these files into Greenstone format, allowing the metadata records to be visible and editable from the GLI
❖ This is irreversible: there's no going back!
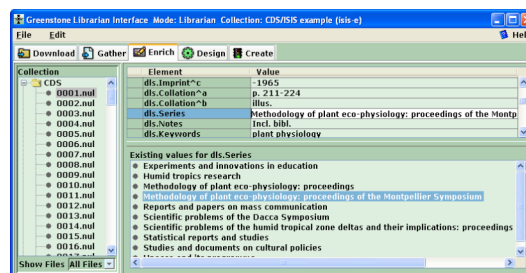  – Keep a backup of the original file

## Exploding Metadata Databases

❖ Files that can be exploded are shown with a green icon in the GLI



❖ Right-click on the file and choose "Explode Metadata Database"

❖ You will be prompted to map the metadata into a metadata set

❖ When complete, the file will be replaced with a folder containing a file for each record



## Exploding Metadata Databases

❖ The metadata can now be viewed and edited:



❖ Note: the GLI is *not* a real database system
  – Only feasible with a small number of records

## Exploding Metadata Databases

❖ Change classifiers, index specifications and format statements to use the namespaced metadata elements

❖ When importing the collection, the files will now be processed by NULPlug rather than the plugin for the original file (eg. ISISPlug)

## Agenda

❖ Formatting
❖ Formatting extended metadata
❖ Changing metadata sets
❖ Collection configuration file
❖ GLI modes
❖ PDF documents
❖ PPT documents
❖ Exploding metadata databases
❖ Cross collection searching
❖ Translate text
❖ Full-text tagging
❖ Creating a CD-ROM

## Cross Collection Searching

- ❖ Select a list of collections

- ❖ Collections need the same indexes

- ❖ User can select which collections to search on Preferences page

- ❖ Format statements applied from original collection
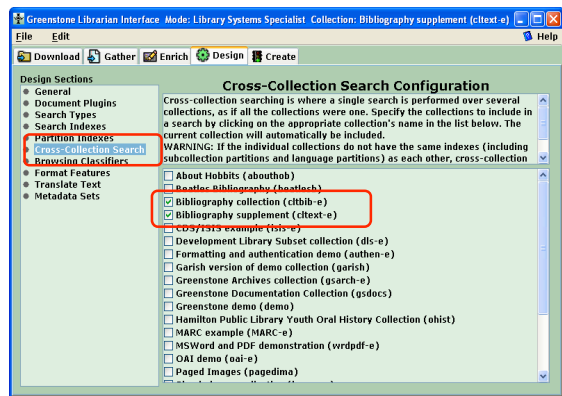
## Bibliography collections

bibliography collection

With about 4,000 bibliography entries, this collection incorporates a form-based search interface that allows fielded searching. It is fairly complex.
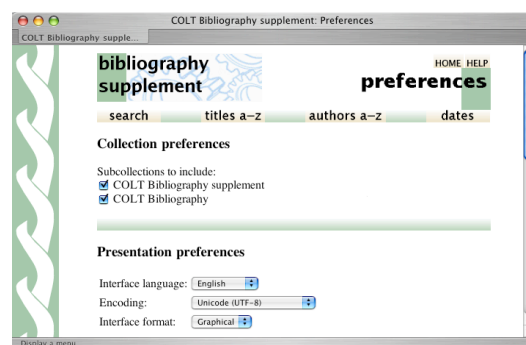
bibliography supplement

This tiny collection of 10 bibliography entries illustrates the "supercollection" facility which searches several collections together, seamlessly. It operates together with the Bibliography collection, and its configuration file is almost the same.
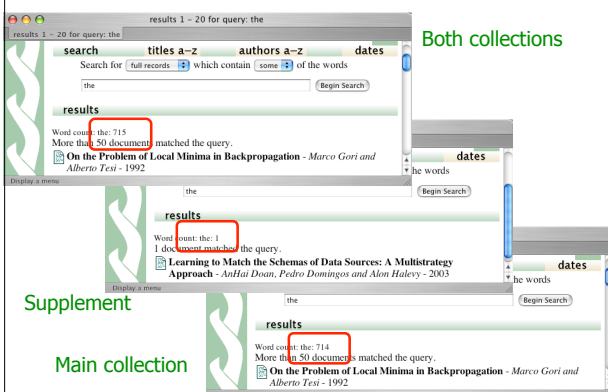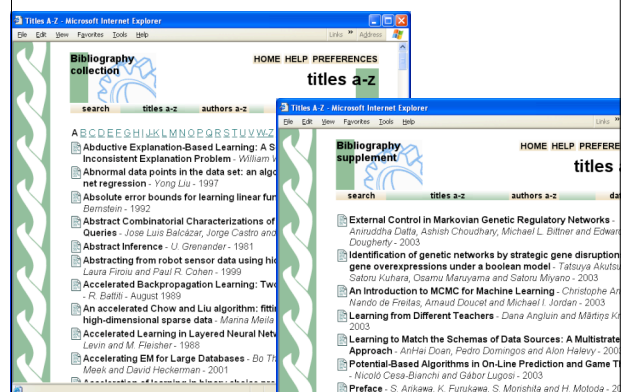
## Cross collection search
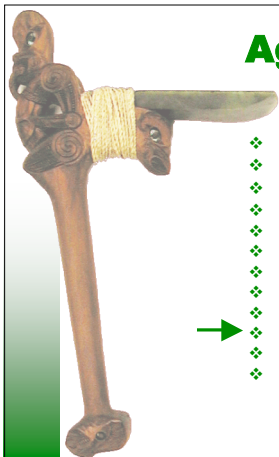


## Bibliography supplement



## Cross-collection searching ...



## ... but not browsing

## Agenda

- ❖ Formatting
- ❖ Formatting extended metadata
- ❖ Changing metadata sets
- ❖ Collection configuration file
- ❖ GLI modes
- ❖ PDF documents
- ❖ PPT documents
- ❖ Exploding metadata databases
- ❖ Cross collection searching
- → ❖ Translate text
- ❖ Full-text tagging
- ❖ Creating a CD-ROM
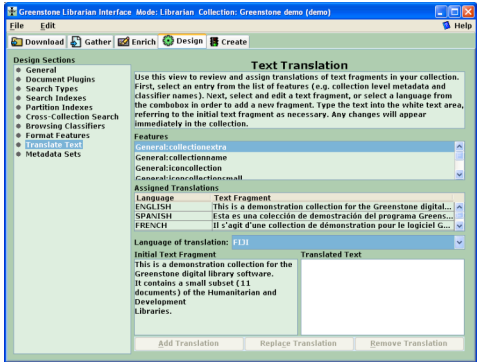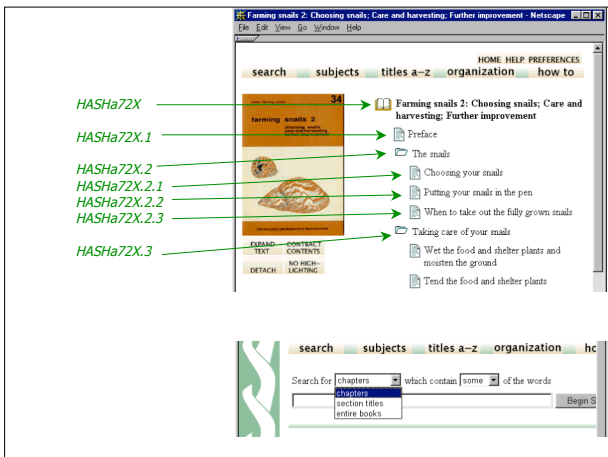
## Translate Text



## Agenda

- ❖ Formatting
- ❖ Formatting extended metadata
- ❖ Changing metadata sets
- ❖ Collection configuration file
- ❖ GLI modes
- ❖ PDF documents
- ❖ PPT documents
- ❖ Exploding metadata databases
- ❖ Cross collection searching
- ❖ Translate text
- → ❖ Full-text tagging
- ❖ Creating a CD-ROM

## Full Text Tagging

- ❖ While creating large digital collections:
  - – the collection must be organized
  - – the larger the collection the greater the need for organization
  - – the larger the documents the greater the need for sections/subsections

- ❖ Greenstone lets you tag the full text of documents

- ❖ Then you can read them hierarchically …

- ❖ … and search them by section



## Full Text Tagging…

To show the hierarchical structure, tag the source files like this:

```
<!--
<Section>
<Description>
<Metadata name="Title">Realizing human rights for
     poor people: Strategies for achieving the
     international development targets</Metadata>
</Description>
-->
     (text of section goes here)
<!-
</Section>
-->
```

## Full Text Tagging...

❖ Section tags define a hierarchical structure
❖ Sections can be nested within other sections
❖ All sections must be nested within a single enclosing section that encompasses the entire document
❖ In the collection configuration file, put

`HTMLPlug -description_tags`

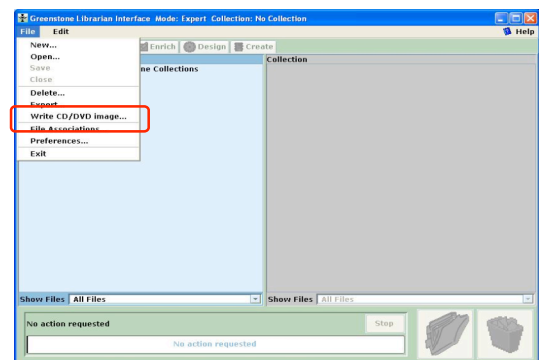❖ Mainly for HTML, but can be used in Word and PDF documents.

## Agenda

❖ Formatting
❖ Formatting extended metadata
❖ Changing metadata sets
❖ Collection configuration file
❖ GLI modes
❖ PDF documents
❖ PPT documents
❖ Exploding metadata databases
❖ Cross collection searching
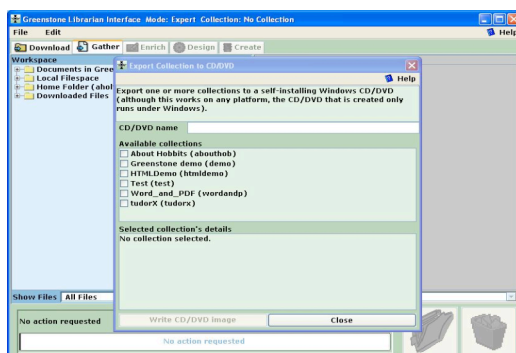❖ Translate text
❖ Full-text tagging
➜ ❖ Creating a CD-ROM

## Creating a CD-ROM

❖ Export any Greenstone collection as a CD-ROM

❖ Self-installing

❖ Windows only (sorry!)

❖ Full Installation of Greenstone

❖ In the Librarian
    File -> Write CD/DVD image

❖ C:\Program Files\Greenstone\tmp\exported_xxx

## Exporting to CD-ROM



## Exporting to CD-ROM



## Note

❖ CD-ROM's created this way have not been tested extensively under different Windows configurations

❖ But they should work on all Windows platforms ...

❖ ... except 3.1/3.11 (is this a problem?)

Creating a CD-ROM:
demo