

## Adding metadata—and using it



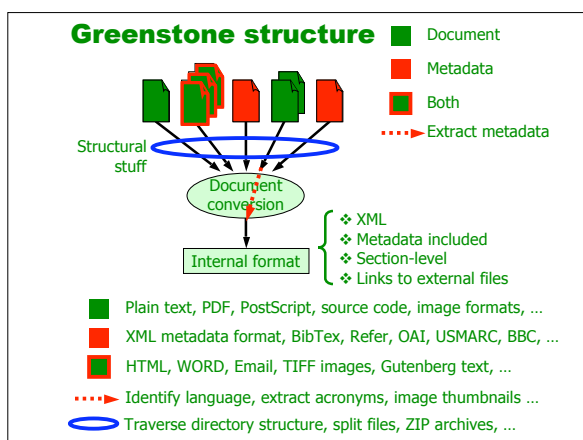
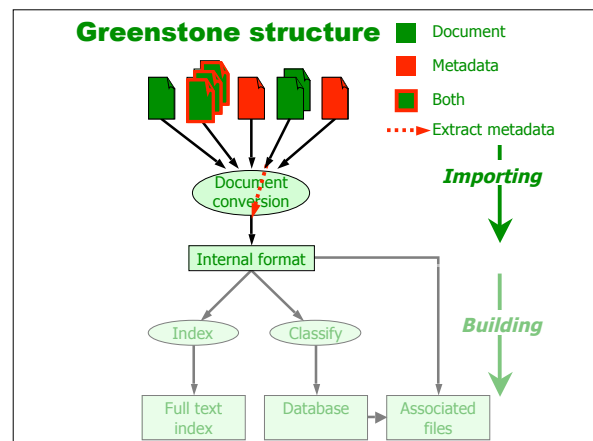
Course material prepared by  
Greenstone Digital Library Project  
University of Waikato, New Zealand

## Agenda

- ❖ Overview of collection building
- ❖ Metadata sets
- ❖ Adding metadata in GLI (+demo)
- ❖ GLI tricks
- ❖ Review: Searching and browsing
- ❖ General Options
- ❖ Plugins
- ❖ Searching – indexes
- ❖ Browsing – classifiers
- ❖ Simple formatting
- ❖ Form search
- ❖ Partitioning indexes
- ❖ PHIND phrase index
- ❖ CDS/ISIS
- ❖ GEMS: metadata set editing

## Collection building process

- ❖ Input: Set of documents + metadata
  - Various formats, e.g. Word, PDF, HTML, images ...
- ❖ Import: Conversion to common internal (XML) format
  - Uses third party tools
- ❖ Build: Create indexes, browsing structures, metadata database
- ❖ Output: Greenstone collection



## Agenda

- ❖ Overview of collection building
- ❖ Metadata sets
- ❖ Adding metadata in GLI (+demo)
- ❖ GLI tricks
- ❖ Review: Searching and browsing
- ❖ General Options
- ❖ Plugins
- ❖ Searching – indexes
- ❖ Browsing – classifiers
- ❖ Simple formatting
- ❖ Form search
- ❖ Partitioning indexes
- ❖ PHIND phrase index
- ❖ CDS/ISIS
- ❖ GEMS: metadata set editing

## Metadata Sets

- ❖ Standard metadata sets
  - Librarian added
  - Good quality, but high cost
  - Many types
    - ❖ Dublin Core (dc)
    - ❖ Development Library Subset (dls)
    - ❖ New Zealand Government Locator Service (nzgls)
- ❖ Greenstone extracted metadata (ex)
  - Automatically extracted
  - Non-editable
  - May be poor quality, but cheap

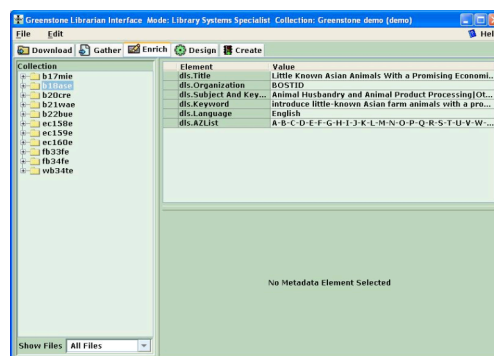
## Dublin Core metadata

Metadata	Tag	Definition
<b>Title</b>	dc.Title	Name given to the resource
<b>Creator</b>	dc.Creator	Entity primarily responsible for making the content of the resource
<b>Subject and keywords</b>	dc.Subject	Topic of the content of the resource
<b>Description</b>	dc.Description	Account of the content of the resource
<b>Publisher</b>	dc.Publisher	Entity responsible for making the resource available
<b>Contributor</b>	dc.Contributor	Entity responsible for making contributions to the content of the resource
<b>Date</b>	dc.Date	Date associated with an event in the life cycle of the resource
<b>Resource type</b>	dc.Type	Nature or genre of the content of the resource
<b>Format</b>	dc.Format	Physical or digital manifestation of the resource
<b>Resource identifier</b>	dc.Identifier	Unambiguous reference to the resource within a given context: this is the object identifier or OID
<b>Source</b>	dc.Source	Reference to a resource from which the present resource is derived
<b>Language</b>	dc.Language	Language of the intellectual content of the resource
<b>Relation</b>	dc.Relation	Reference to a related resource
<b>Coverage</b>	dc.Coverage	Extent or scope of the content of the resource
<b>Rights management</b>	dc.Rights	Information about rights held in and over the resource

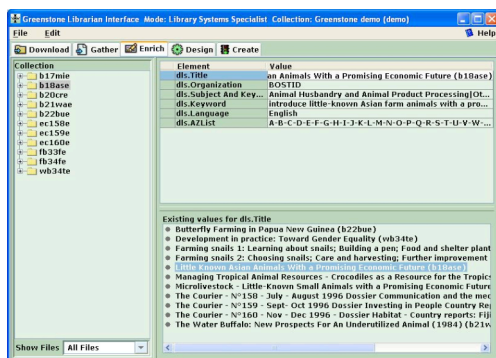
## DLS metadata

Metadata	Definition
dls.Title	Title
dls.Keyword	Keyword (originally "How To")
dls.Organization	Organization that produced the document
dls.Language	Language of the document
dls.Subject and Keywords	Hierarchical subject metadata
dls.AZList	A-Z range into which this document's title falls

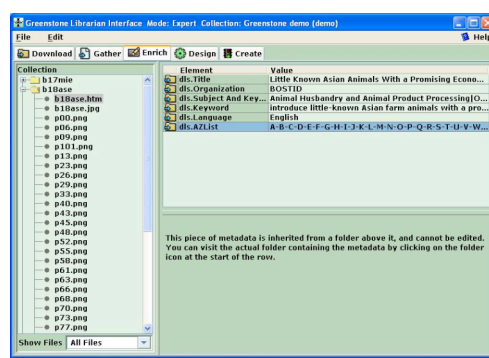
## Example of DLS metadata



## Example of DLS metadata



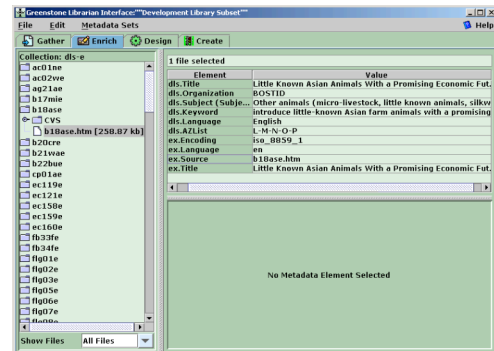
## Example of DLS metadata



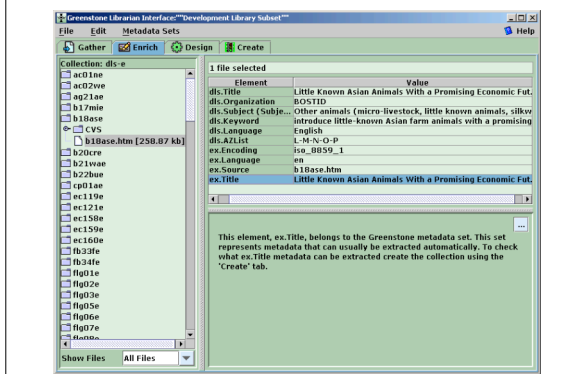
## Extracted metadata

Metadata	Definition
ex.Title	Title extracted from HTML, Word, PDF
ex.Source	Name of source file, e.g. Apache.html
ex.Language	Language of document, e.g. en, sp, fr
ex.Encoding	Encoding of document, e.g. iso_8859_1
...	...
ex.Acronym	Acronyms that appear in the document
...	...

## Example of extracted metadata



## Example of extracted metadata



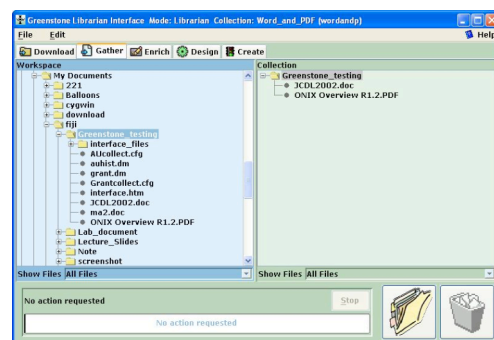
## Agenda

- ❖ Overview of collection building
- ❖ Metadata sets
- ❖ Adding metadata in GLI (+demo)
- ❖ GLI tricks
- ❖ Review: Searching and browsing
- ❖ General Options
- ❖ Plugins
- ❖ Searching – indexes
- ❖ Browsing – classifiers
- ❖ Simple formatting
- ❖ Form search
- ❖ Partitioning indexes
- ❖ PHIND phrase index
- ❖ CDS/ISIS
- ❖ GEMS: metadata set editing

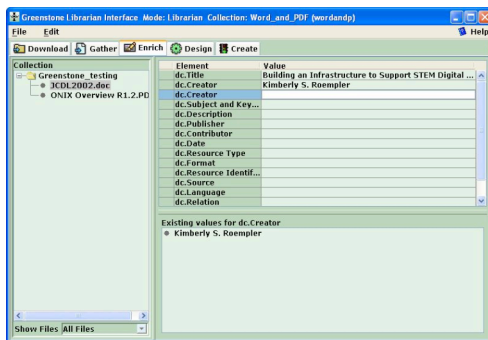
## GLI: Adding metadata

- ❖ Enrich pane
- ❖ Folder level or document level
- ❖ Specified metadata set
- ❖ Can view extracted metadata but not edit

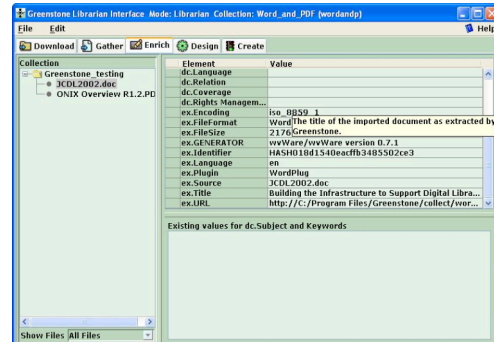
## GLI: adding metadata



## GLI: adding metadata



## Extracted Metadata



## Greenstone Librarian Interface demo: adding metadata

## Hierarchical Metadata

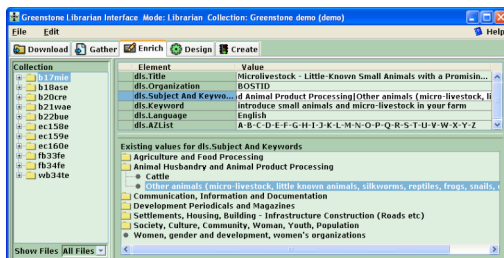
- ❖ Metadata is often hierarchical:



- ❖ How is this done in Greenstone?

## Hierarchical Metadata

- ❖ In the GLI's Enrich pane, use the pipe character (|) to separate the levels



- ❖ Preview the hierarchy in the value tree

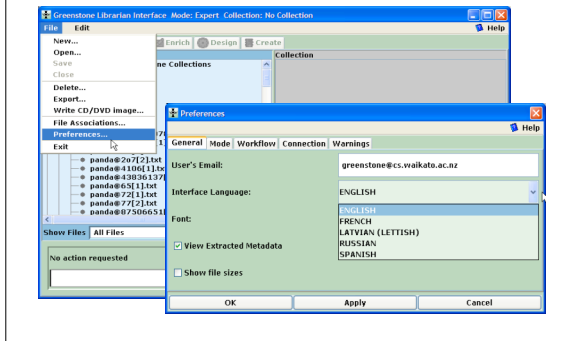
## Agenda

- ❖ Overview of collection building
- ❖ Metadata sets
- ❖ Adding metadata in GLI (+demo)
- ❖ GLI tricks
- ❖ Review: Searching and browsing
- ❖ General Options
- ❖ Plugins
- ❖ Searching – indexes
- ❖ Browsing – classifiers
- ❖ Simple formatting
- ❖ Form search
- ❖ Partitioning indexes
- ❖ PHIND phrase index
- ❖ CDS/ISIS
- ❖ GEMS: metadata set editing



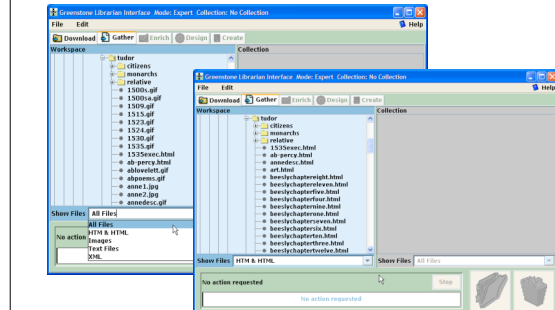
## GLI Tricks

### ❖ General Preferences



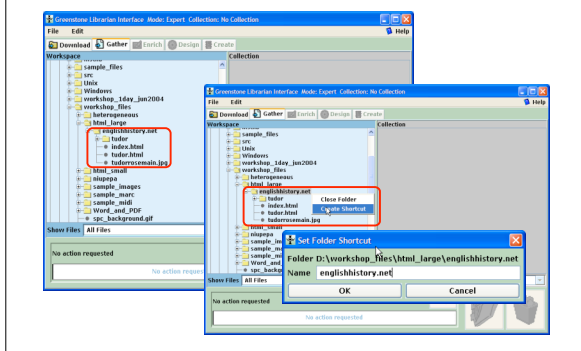
## GLI Tricks

### ❖ File Filtering

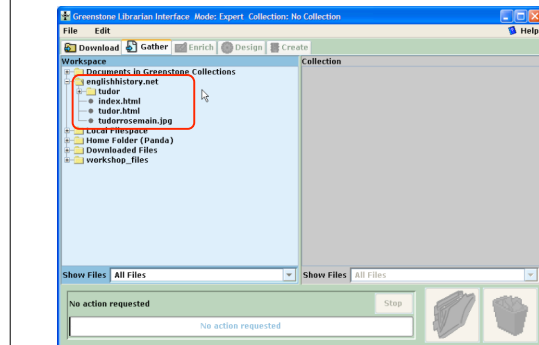


## GLI Tricks

### ❖ Create file folder shortcut

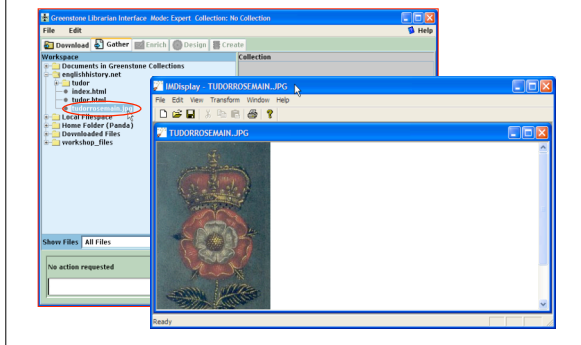


## GLI Tricks



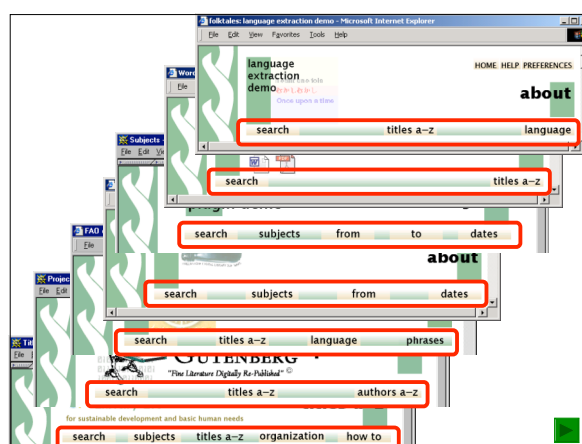
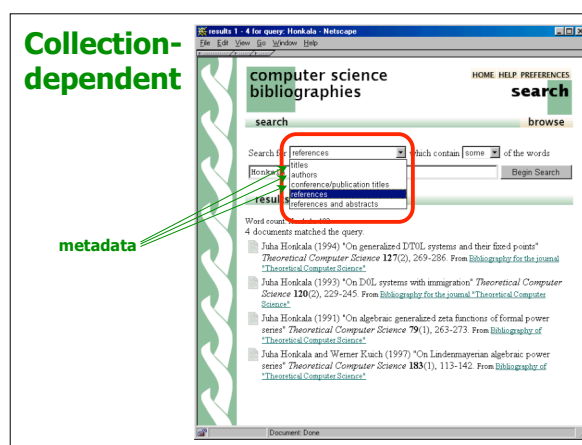
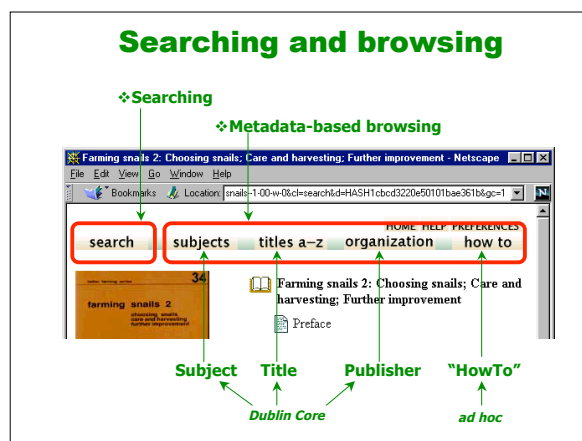
## GLI Tricks

### ❖ Viewing files: Double-clicking



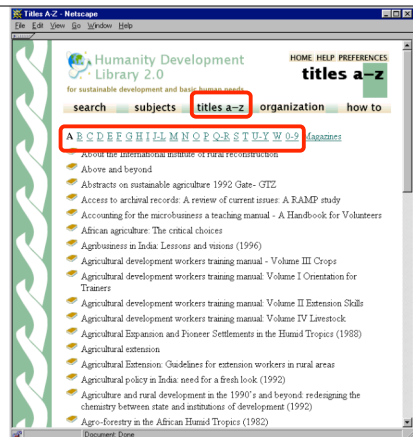
## Agenda

- ❖ Overview of collection building
- ❖ Metadata sets
- ❖ Adding metadata in GLI (+demo)
- ❖ GLI tricks
- ❖ Review: Searching and browsing
- ❖ General Options
- ❖ Plugins
- ❖ Searching – indexes
- ❖ Browsing – classifiers
- ❖ Simple formatting
- ❖ Form search
- ❖ Partitioning indexes
- ❖ PHIND phrase index
- ❖ CDS/ISIS
- ❖ GEMS: metadata set editing

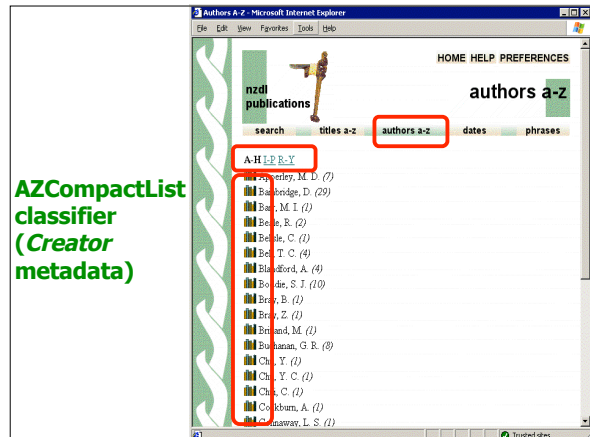


## Browsing using classifiers

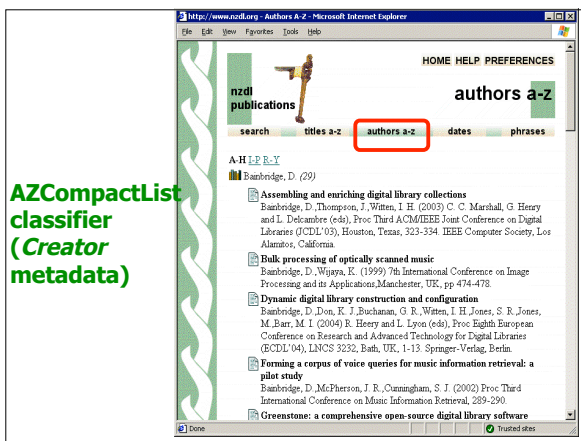
**AZList classifier (Title metadata)**



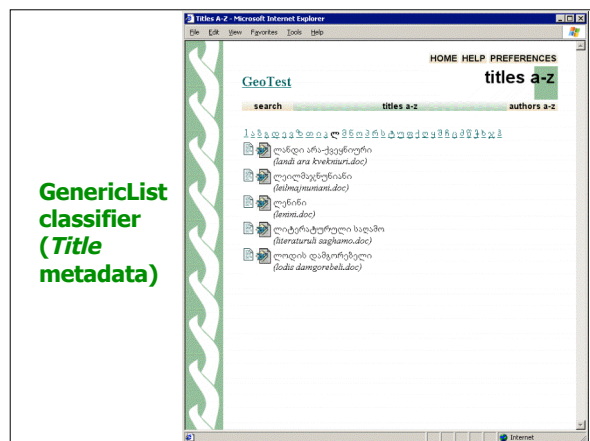
**AZCompactList classifier (Creator metadata)**



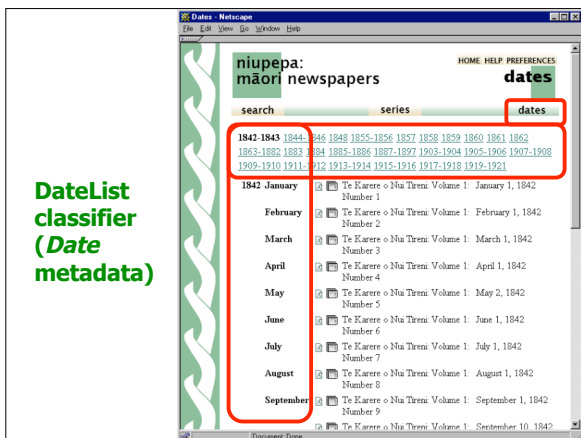
**AZCompactList classifier (Creator metadata)**



**GenericList classifier (Title metadata)**

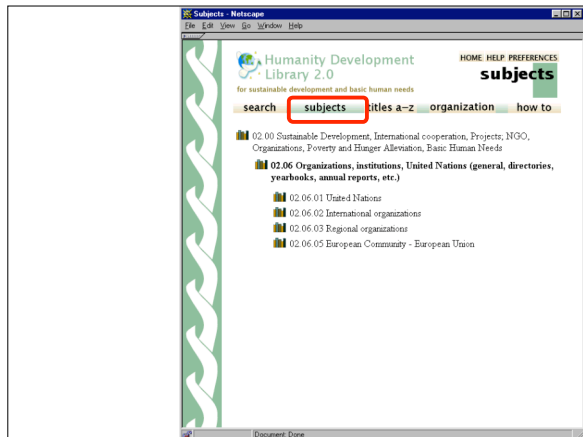


**DateList classifier (Date metadata)**



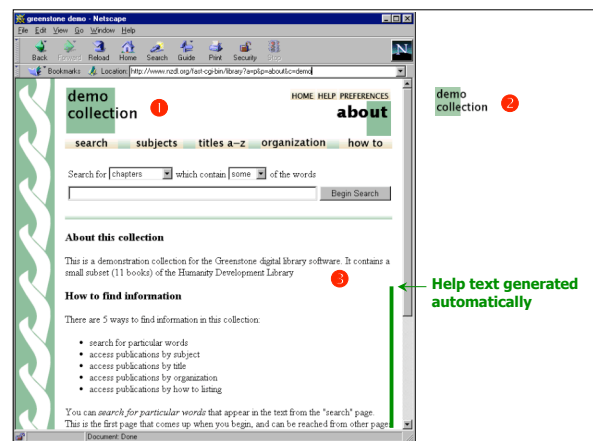
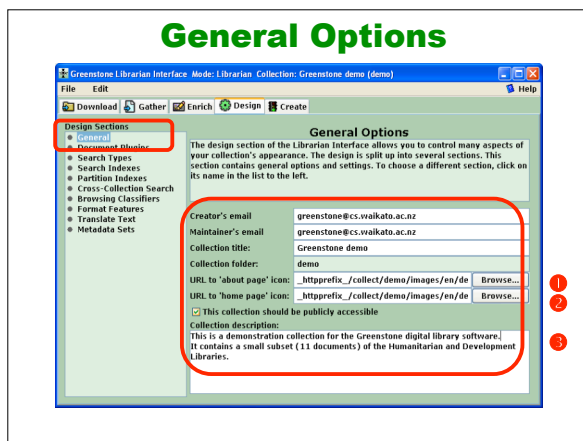
**Hierarchy classifier (Subject metadata)**





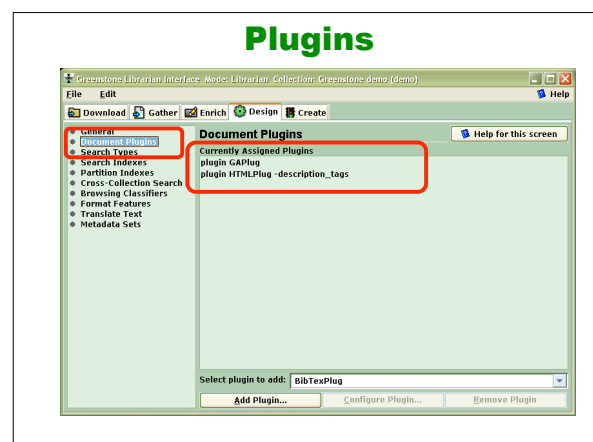

## Agenda

- ❖ Overview of collection building
- ❖ Metadata sets
- ❖ Adding metadata in GLI (+demo)
- ❖ GLI tricks
- ❖ Review: Searching and browsing
- ❖ General Options
- ❖ Plugins
- ❖ Searching – indexes
- ❖ Browsing – classifiers
- ❖ Simple formatting
- ❖ Form search
- ❖ Partitioning indexes
- ❖ PHIND phrase index
- ❖ CDS/ISIS
- ❖ GEMS: metadata set editing

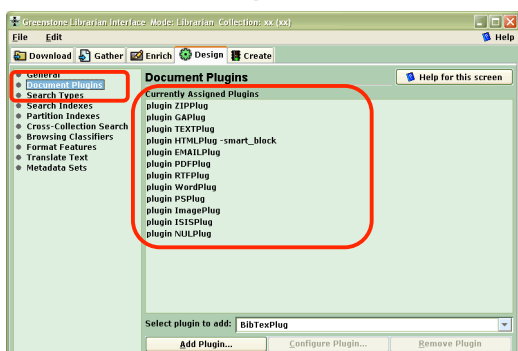



## Agenda

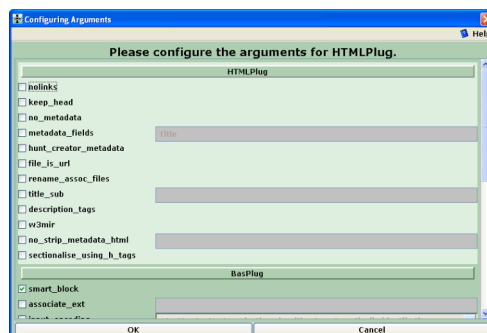
- ❖ Overview of collection building
- ❖ Metadata sets
- ❖ Adding metadata in GLI (+demo)
- ❖ GLI tricks
- ❖ Review: Searching and browsing
- ❖ General Options
- ❖ Plugins
- ❖ Searching – indexes
- ❖ Browsing – classifiers
- ❖ Simple formatting
- ❖ Form search
- ❖ Partitioning indexes
- ❖ PHIND phrase index
- ❖ CDS/ISIS
- ❖ GEMS: metadata set editing



## Plugins



## Plugin Options



## Plugins

Used by collection-building software to accomplish format-specific parsing of source documents

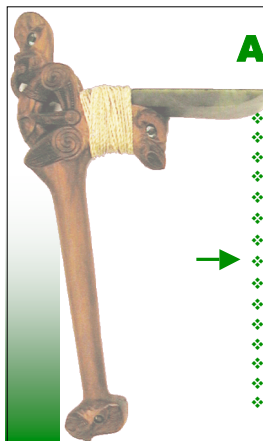
Plugin pipeline: files are passed to each plugin in turn until one is found that can process it

- ❖ GAPug processes doc.xml files generated during *import*
- ❖ ArcPlug processes filelist in *archives.inf*
- ❖ RecPlug recurses through a directory structure

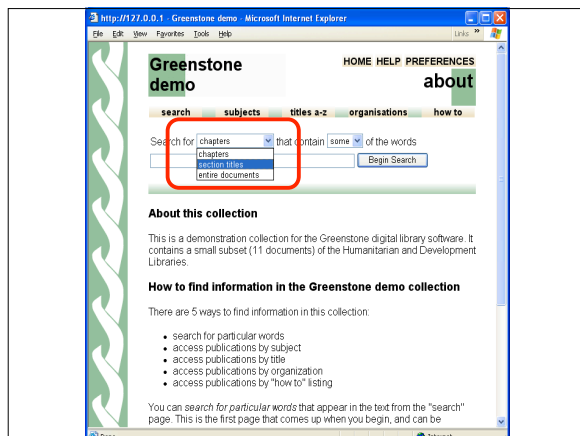
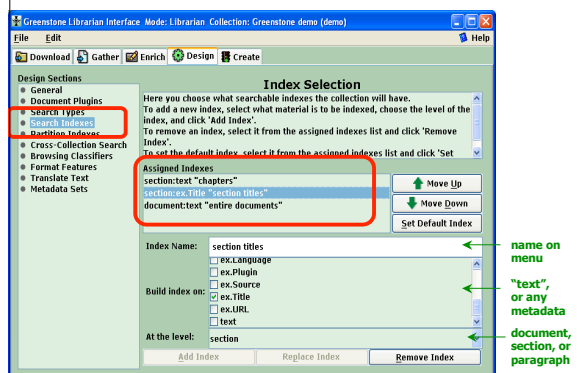
TEXTPlug  
HTMLPlug  
EMAILPlug  
WORDPlug  
RTFPlug  
PDFPlug  
PSPlug  
ImagePlug  
PPTPlug  
ISISPlug  
TCCPlug ...

## Agenda

- ❖ Overview of collection building
- ❖ Metadata sets
- ❖ Adding metadata in GLI (+demo)
- ❖ GLI tricks
- ❖ Review: Searching and browsing
- ❖ General Options
- ❖ Plugins
- ❖ Searching – indexes
- ❖ Browsing – classifiers
- ❖ Simple formatting
- ❖ Form search
- ❖ Partitioning indexes
- ❖ PHIND phrase index
- ❖ CDS/ISIS
- ❖ GEMS: metadata set editing



## Search Indexes

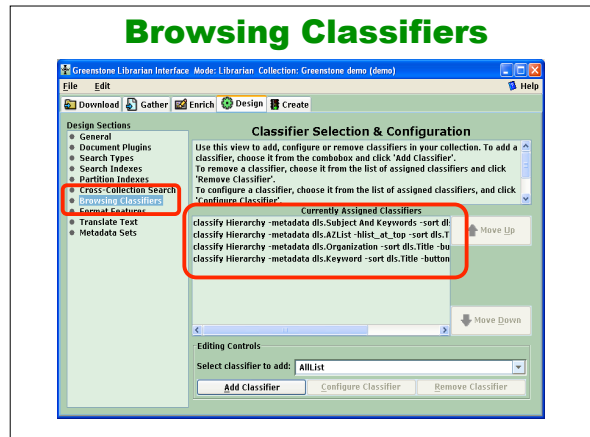
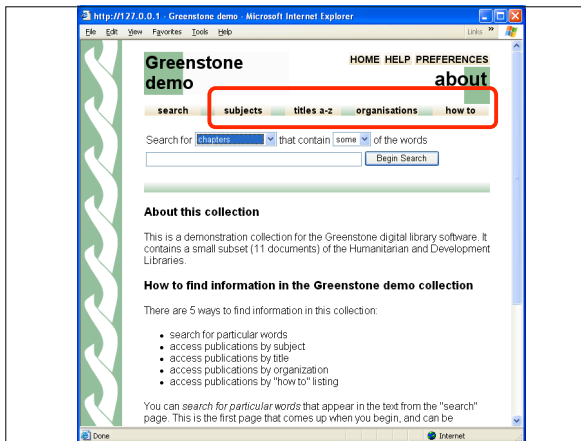




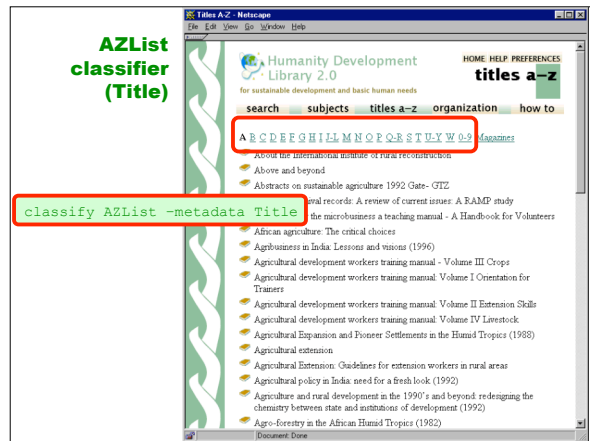
## Agenda

- ❖ Overview of collection building
- ❖ Metadata sets
- ❖ Adding metadata in GLI (+demo)
- ❖ GLI tricks
- ❖ Review: Searching and browsing
- ❖ General Options
- ❖ Plugins
- ❖ Searching – indexes
- ➔ ❖ Browsing – classifiers
- ❖ Simple formatting
- ❖ Form search
- ❖ Partitioning indexes
- ❖ PHIND phrase index
- ❖ CDS/ISIS
- ❖ GEMS: metadata set editing

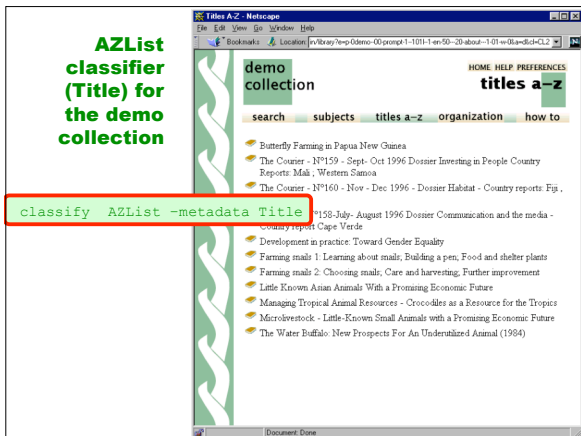
## Browsing Classifiers

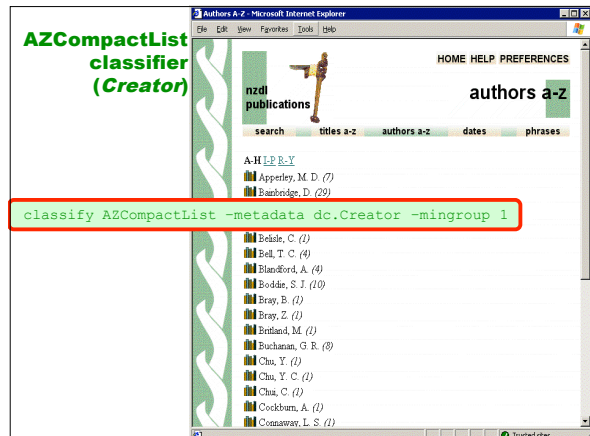
## AZList classifier (Title)



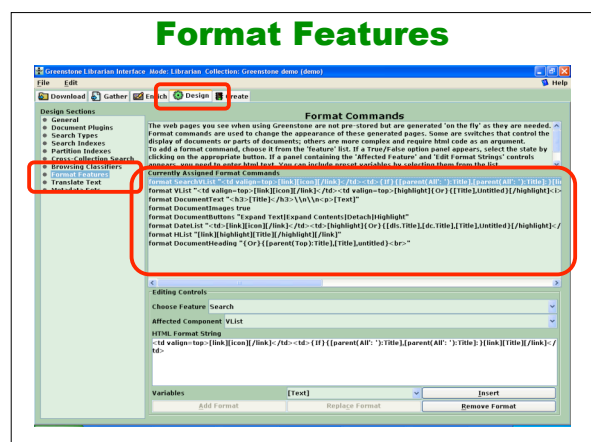
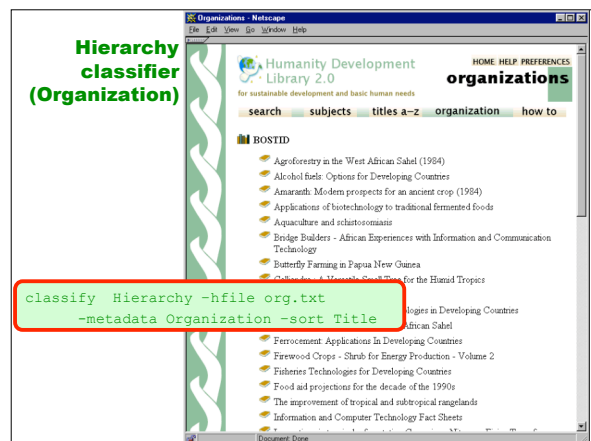
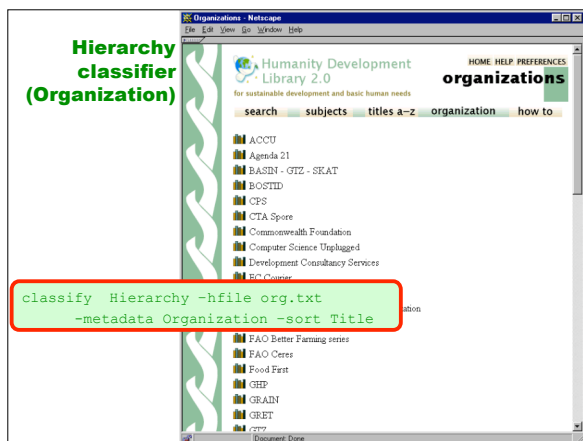
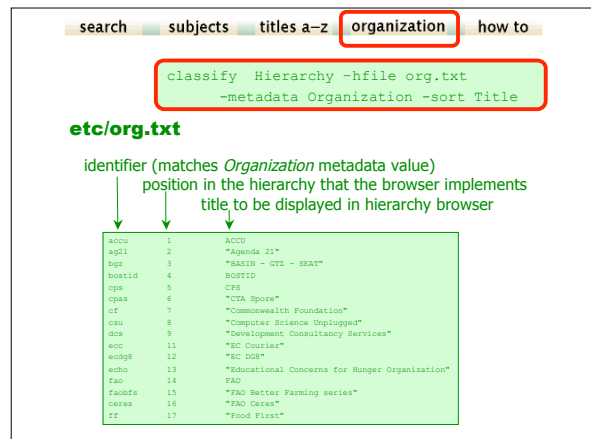
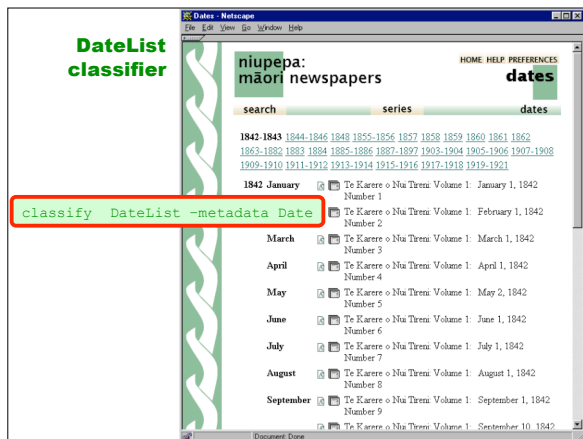
## AZList classifier (Title) for the demo collection



## AZCompactList classifier (Creator)



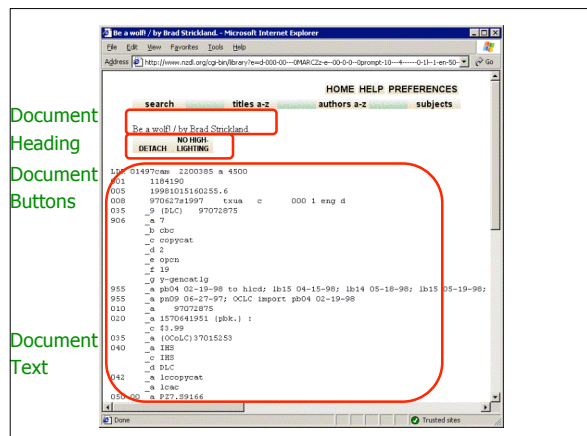
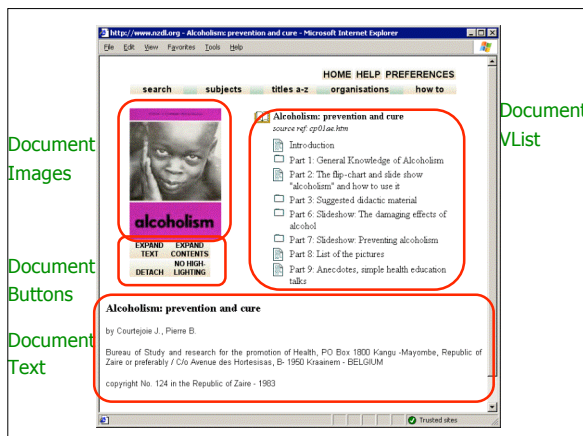
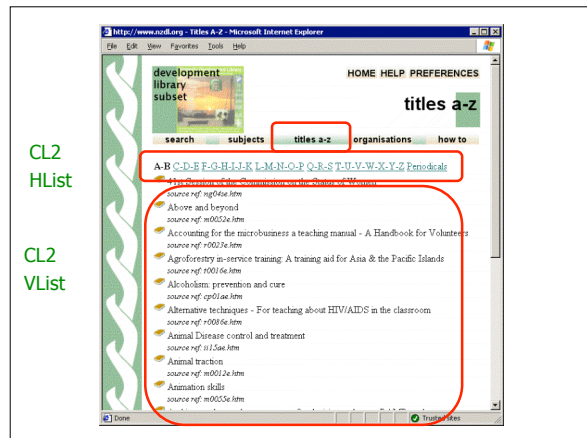
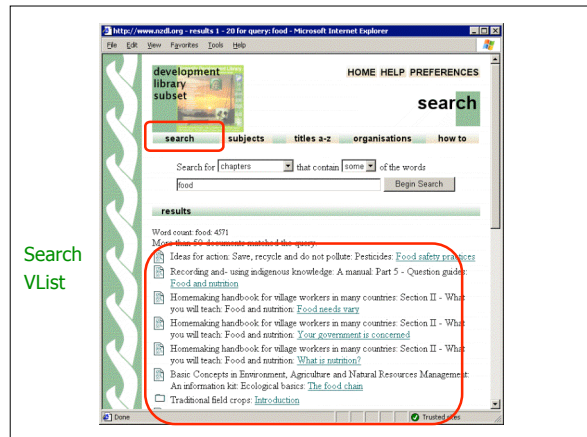






## Format Statements

- ❖ **VList** – vertical lists: search, classifiers, document table of contents
- ❖ **HList** – horizontal lists: classifiers
- ❖ **SearchVList**, **CL1VList** – individual vertical lists
- ❖ **DocumentImages/DocumentHeading/DocumentContents** – document header display
- ❖ **DocumentText** – document display





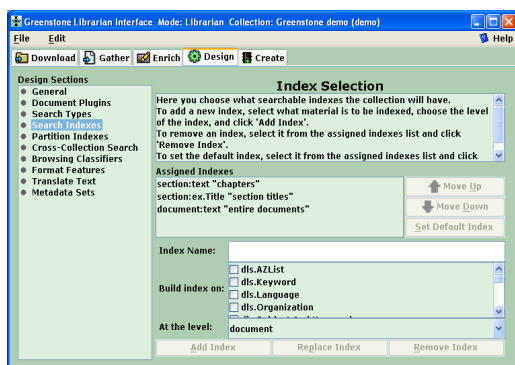
## Agenda

- ❖ Overview of collection building
- ❖ Metadata sets
- ❖ Adding metadata in GLI (+demo)
- ❖ GLI tricks
- ❖ Review: Searching and browsing
- ❖ General Options
- ❖ Plugins
- ❖ Searching – indexes
- ❖ Browsing – classifiers
- ❖ Simple formatting
- ❖ Form search
- ❖ Partitioning indexes
- ❖ PHIND phrase index
- ❖ CDS/ISIS
- ❖ GEMS: metadata set editing

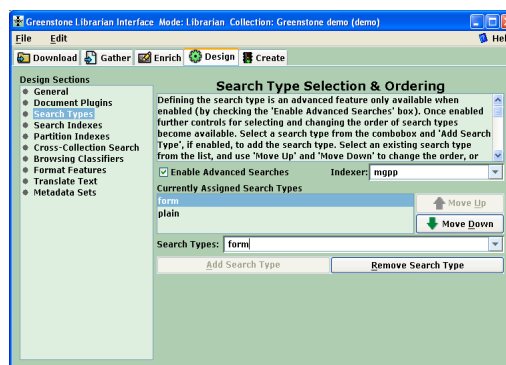
## Search Types

- ❖ “Enable search types” – allows fielded searching
- ❖ MGPP/Lucene
- ❖ Form/plain – can switch on preferences
- ❖ Different index specification

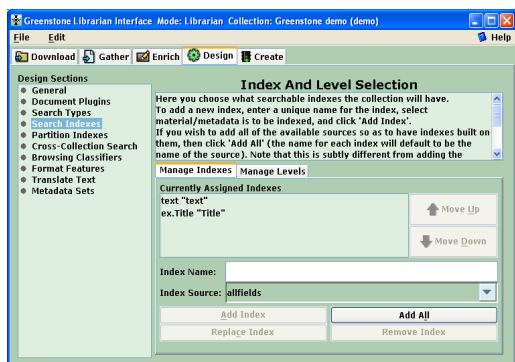
## Default Search Indexes



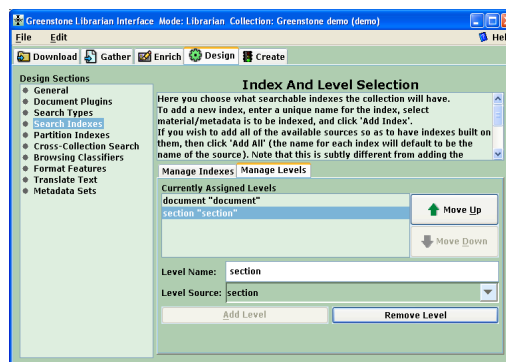
## Search Types



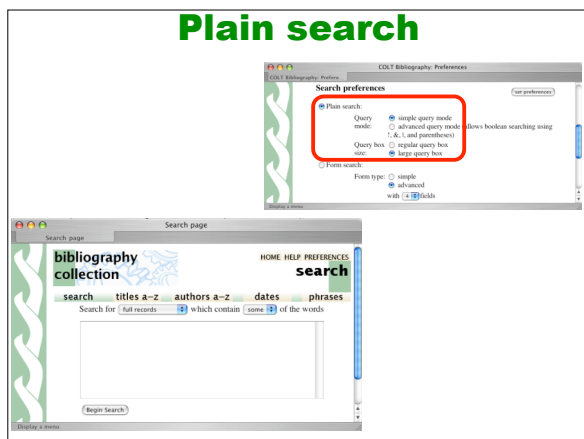
## Form Search Indexes: Indexes



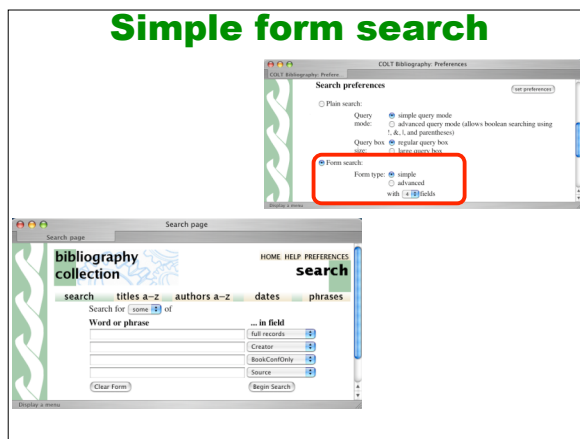
## Form Search Indexes: Levels



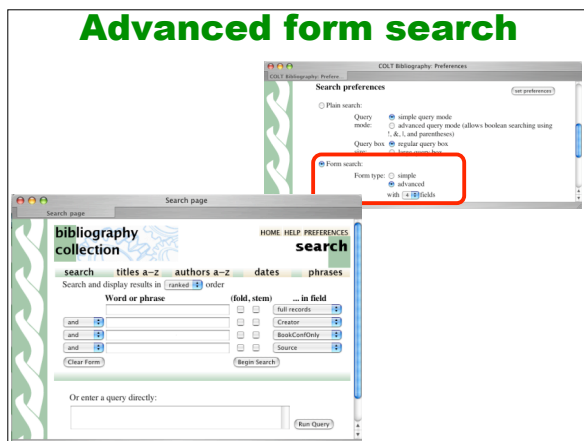
## Plain search



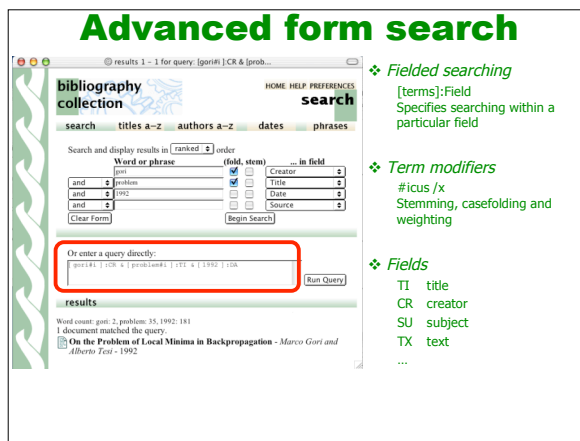
## Simple form search



## Advanced form search

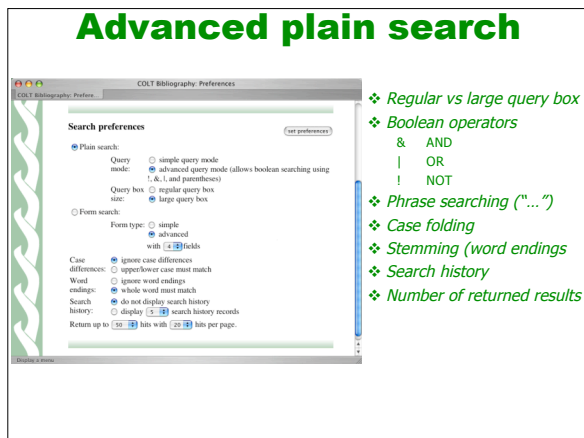


## Advanced form search



- ❖ **Fielded searching**  
[terms]:Field  
Specifies searching within a particular field
- ❖ **Term modifiers**  
#icus /x  
Stemming, casefolding and weighting
- ❖ **Fields**  
TI title  
CR creator  
SU subject  
TX text  
...

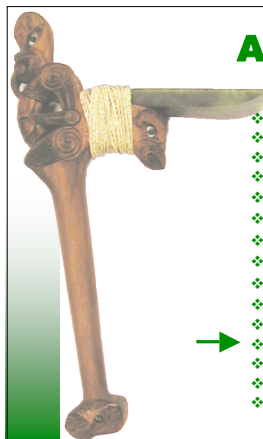
## Advanced plain search



- ❖ **Regular vs large query box**
- ❖ **Boolean operators**  
& AND  
| OR  
! NOT
- ❖ **Phrase searching ("...")**
- ❖ **Case folding**
- ❖ **Stemming (word endings)**
- ❖ **Search history**
- ❖ **Number of returned results**

## Agenda

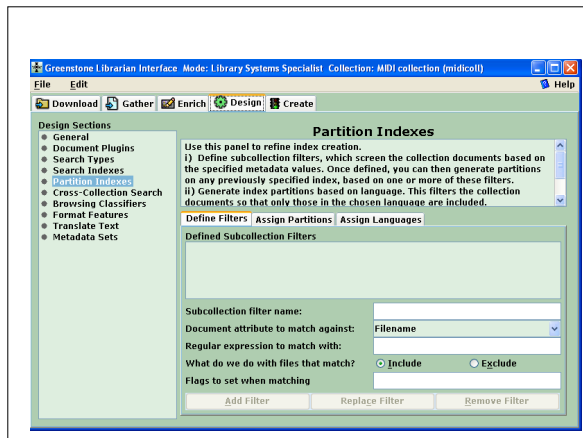
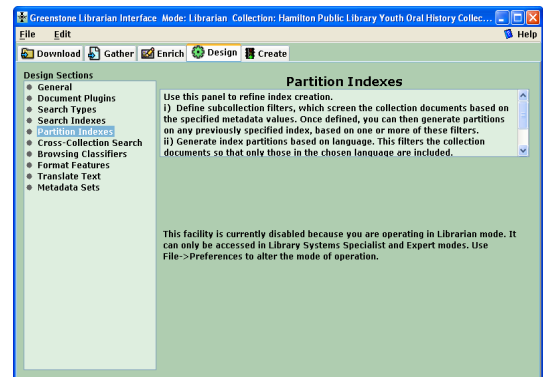
- ❖ Overview of collection building
- ❖ Metadata sets
- ❖ Adding metadata in GLI (+demo)
- ❖ GLI tricks
- ❖ Review: Searching and browsing
- ❖ General Options
- ❖ Plugins
- ❖ Searching – indexes
- ❖ Browsing – classifiers
- ❖ Simple formatting
- ❖ Form search
- ❖ Partitioning indexes
- ❖ PHIND phrase index
- ❖ CDS/ISIS
- ❖ GEMS: metadata set editing



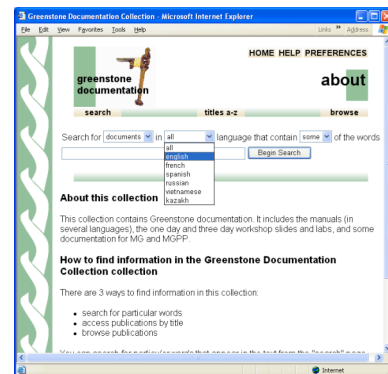
## Partition Indexes

- ❖ Searching within a single collection
- ❖ The indexes are split based on subcollection definitions
  - E.g. language, file types, metadata
- ❖ Browsing classifiers contain all the documents
- ❖ Library Systems Specialist mode

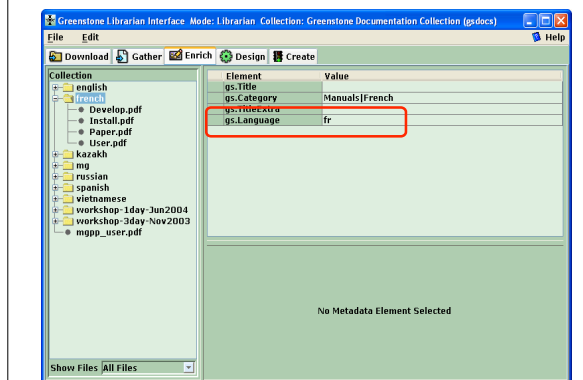
## Partition Indexes



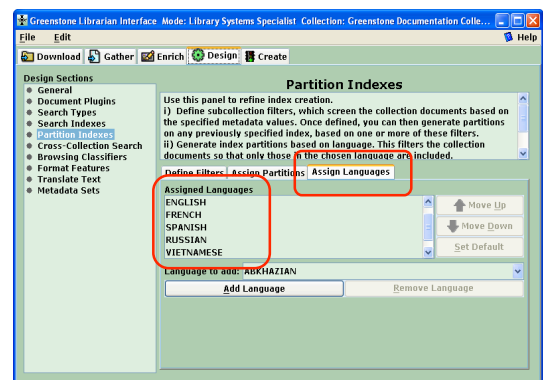
## Language partitions



## Language metadata



## Language partitions



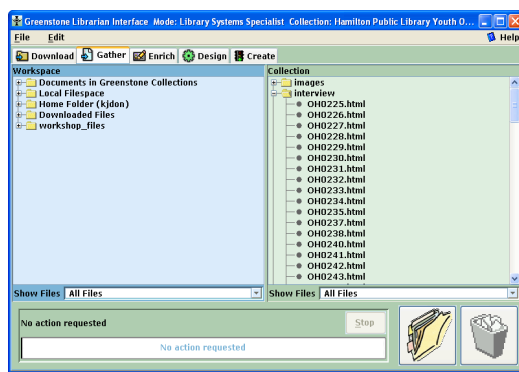
## Document type partitions



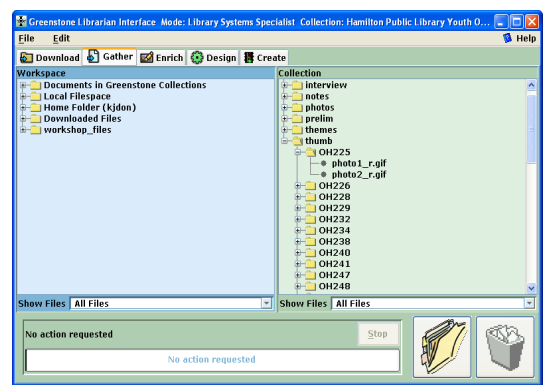
## Document type partitions



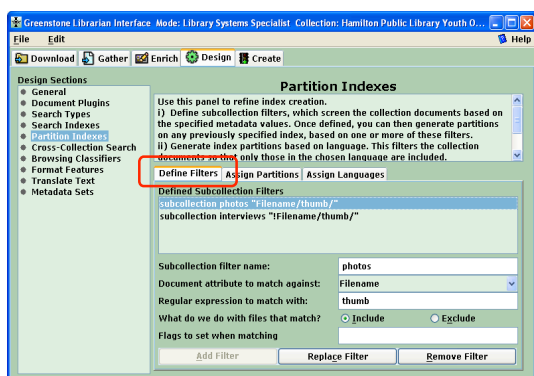
## Collection documents



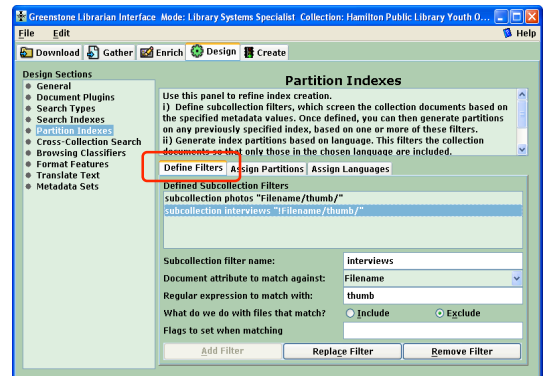
## Collection documents



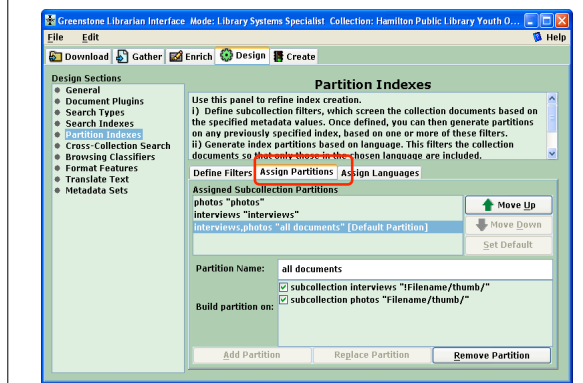
## Subcollection filters



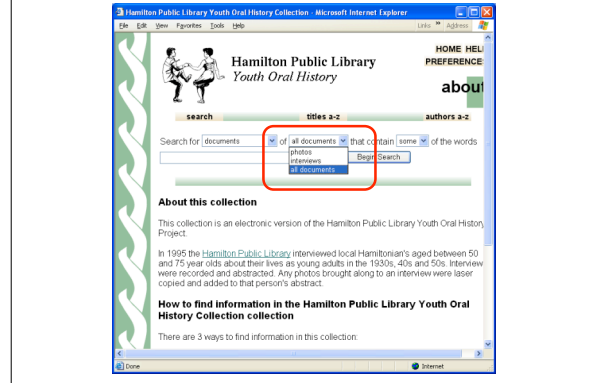
## Subcollection filters



## Assign partitions



## Document type partitions

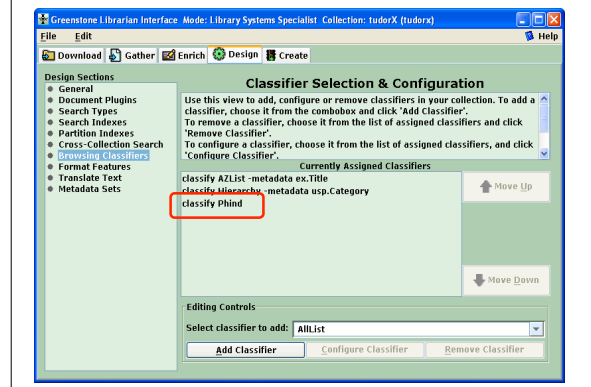


## Agenda

- ❖ Overview of collection building
- ❖ Metadata sets
- ❖ Adding metadata in GLI (+demo)
- ❖ GLI tricks
- ❖ Review: Searching and browsing
- ❖ General Options
- ❖ Plugins
- ❖ Searching – indexes
- ❖ Browsing – classifiers
- ❖ Simple formatting
- ❖ Form search
- ❖ Partitioning indexes
- ❖ PHIND phrase index
- ❖ CDS/ISIS
- ❖ GEMS: metadata set editing



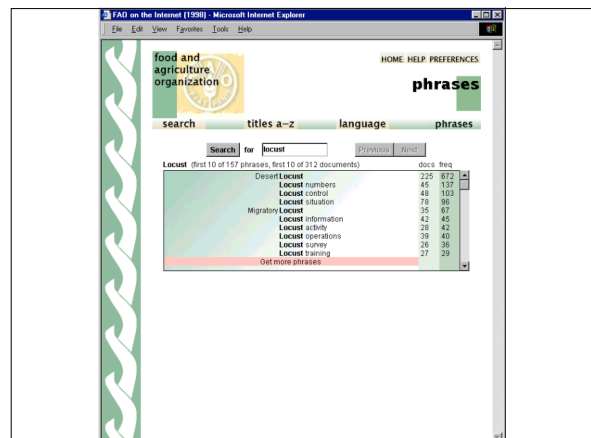
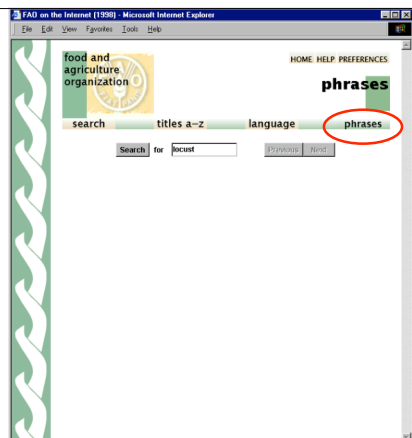
## Phind Phrase Index



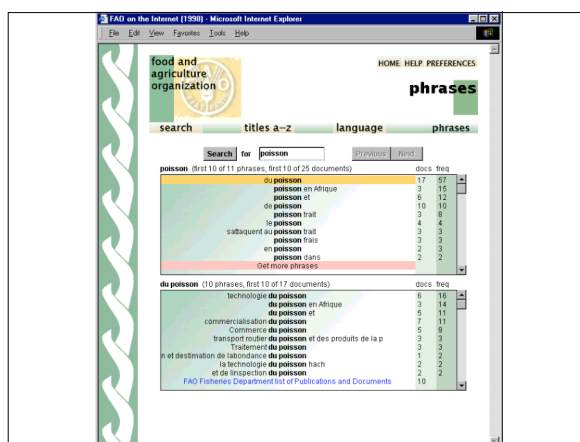
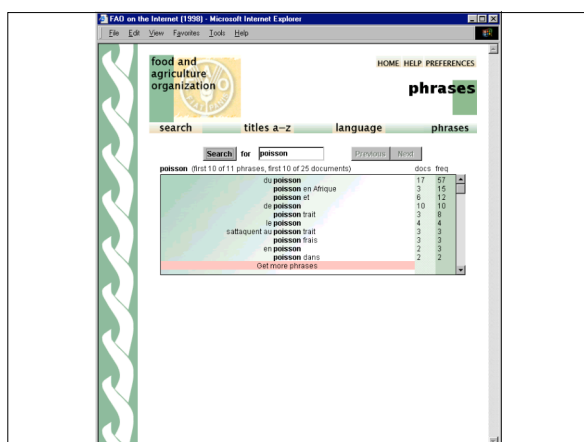
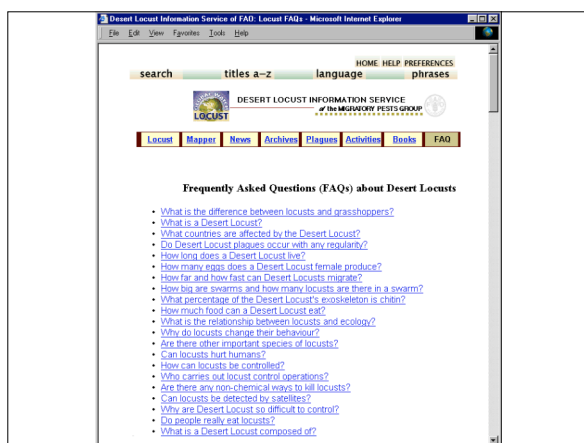
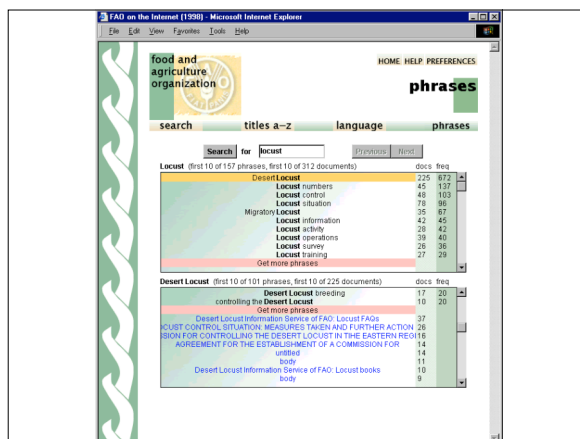
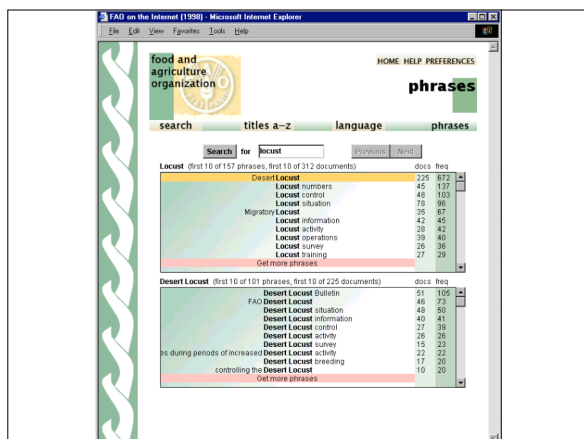
## PHIND

### Hierarchical index of phrases

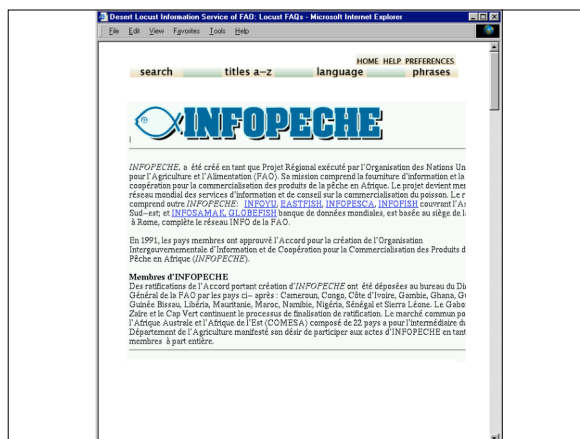
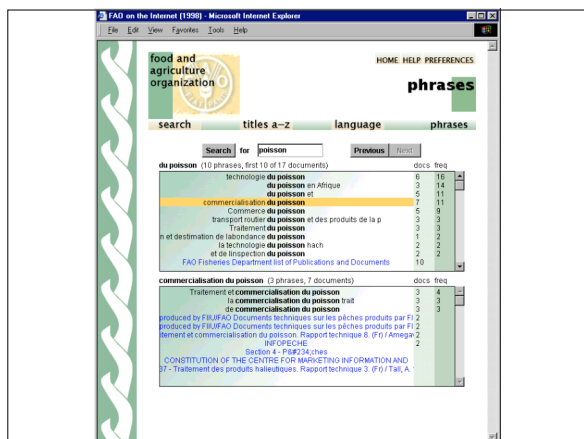
- ❖ What's in this collection?
- ❖ Is it any good?
- ❖ What coverage for topic X?
- ❖ My query returned too much/little, what now?












## Agenda

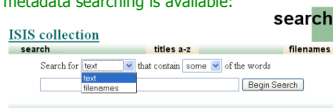
- ❖ Overview of collection building
- ❖ Metadata sets
- ❖ Adding metadata in GLI (+demo)
- ❖ GLI tricks
- ❖ Review: Searching and browsing
- ❖ General Options
- ❖ Plugins
- ❖ Searching – indexes
- ❖ Browsing – classifiers
- ❖ Simple formatting
- ❖ Form search
- ❖ Partitioning indexes
- ❖ PHIND phrase index
- ❖ CDS/ISIS
- ❖ GEMS: metadata set editing

## CDS/ISIS

- ❖ Bibliography collections are typically fairly complex:
  - Form searching
  - Customised query result and browse lists
  - Customised document display
- ❖ Let's work through creating a simple collection using a small CDS/ISIS database describing a set of film slides  
(More information in the "Bibliography collection" and "CDS/ISIS" documented example collections)

## CDS/ISIS

- ❖ Add the CDS/ISIS files to a new collection
- ❖ After building, let's view the collection:
  - No metadata searching is available:



- The titles classifier is completely empty!



## CDS/ISIS

- ❖ More problems:
  - The filenames classifier is useless!

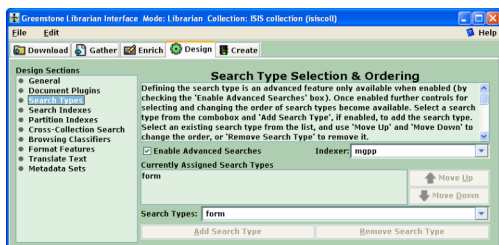


- The document display isn't very pretty:



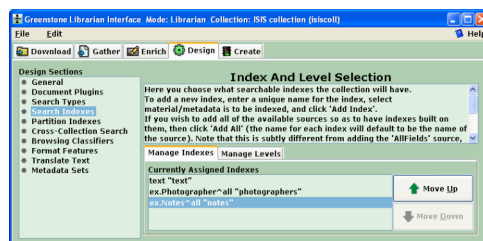
## CDS/ISIS: Metadata searching

- ❖ To enable form searching, go to the "Search Types" area in the GLI's Design pane
  - Tick "Enable Advanced Searches" on
  - Add the "form" search type, and remove "plain"



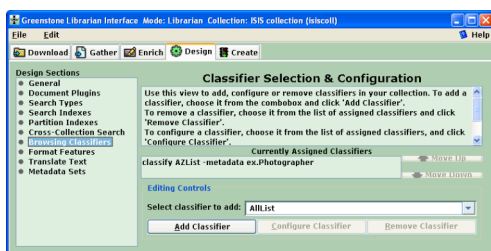
## CDS/ISIS: Metadata Searching

- ❖ Add metadata indexes in the "Search Indexes" part of the GLI's Design pane
  - Add indexes for Photographer and Notes metadata
  - Remove the useless Source and Title indexes



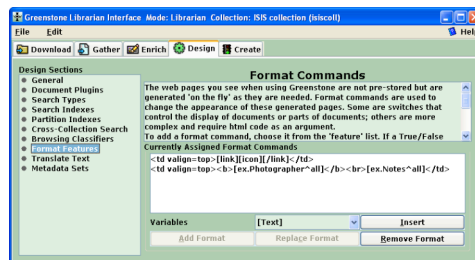
## CDS/ISIS: Better browsing

- ❖ Remove the existing (useless) classifiers for Title and Source metadata, and add a new one for Photographer



## CDS/ISIS: Better browsing

- ❖ Change the VList format statement to display the Photographer and Notes metadata:



## CDS/ISIS: Document display

- ❖ Next, let's change the DocumentText format statement to show the Photographer and Notes metadata:
 

```
<center><table width= pagewidth ><tr><td>Photographer:
</td><td>[ex.Photographer*all]</td></tr><tr><td>Notes:
</td><td>[ex.Notes*all]</td></tr></table></center>
```
- ❖ Then, let's remove those annoying "Detach" and "Highlight" buttons by setting DocumentButtons to empty
- ❖ Lastly, clear DocumentHeading to remove the "untitled" at the top of the document

## CDS/ISIS: Finished!

- ❖ Metadata searching now available:



- ❖ Better browsing facilities:



## CDS/ISIS: Finished!

### ❖ Document display improved:

search Photographer

Photographer: Australia: Bureau of Mineral Resources  
United States Geological Survey

Notes: Slides for presentations on the petroleum potential of the SOPAC region  
(New Ireland (PNG), Solomon Islands, Tonga and Vanuatu), obtained  
from N.F. Eason

### ❖ What could still be improved?

- More metadata indexes, classifiers
- Display all fields in the document display
- Nice images for classifiers
- ...?

## Agenda

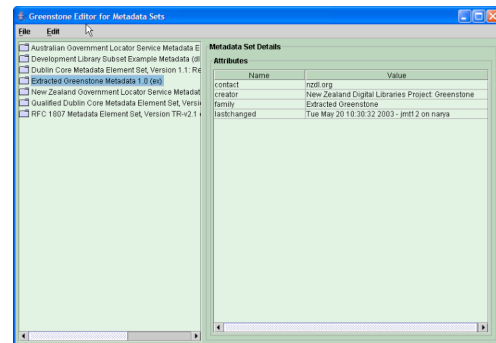
- ❖ Overview of collection building
- ❖ Metadata sets
- ❖ Adding metadata in GLI (+demo)
- ❖ GLI tricks
- ❖ Review: Searching and browsing
- ❖ General Options
- ❖ Plugins
- ❖ Searching – indexes
- ❖ Browsing – classifiers
- ❖ Simple formatting
- ❖ Form search
- ❖ Partitioning indexes
- ❖ PHIND phrase index
- ❖ CDS/ISIS
- ❖ GEMS: metadata set editing

## Greenstone Editor for Metadata Sets - GEMS

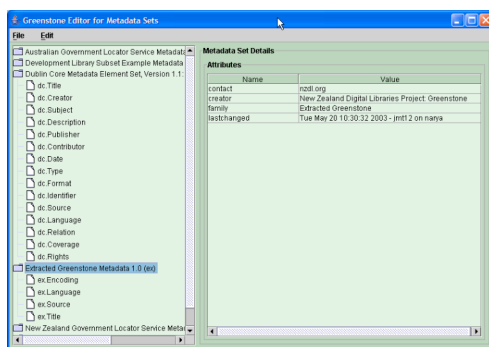
### ❖ What can GEMS do?

- Create a new metadata set
- Add metadata elements to an existing set
- Define and edit metadata attributes
- Inherit from an existing metadata set
- Import an existing metadata set
- Sets can be used in GLI

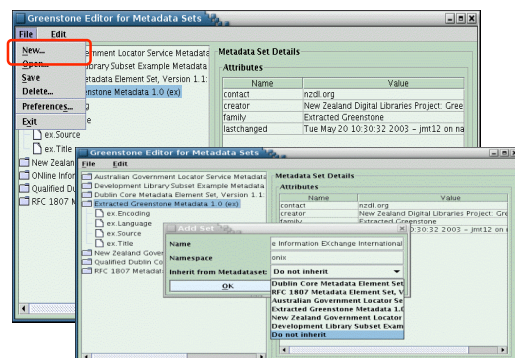
## GEMS Interface



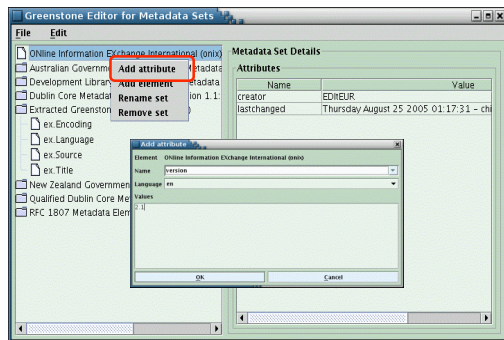
## GEMS Interface



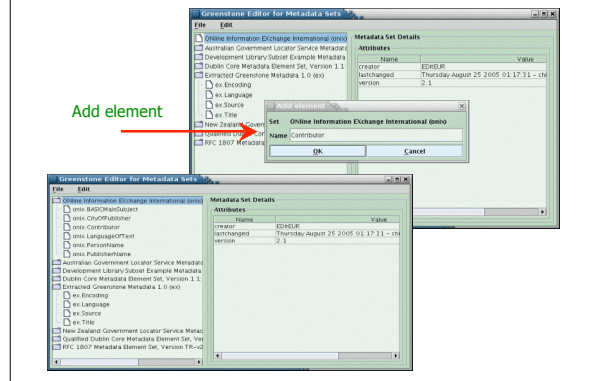
## Adding a New Metadata Set



## Adding Metadata Attribute



## Adding Metadata Element



## Using an existing metadata set

- ❖ Using a pre-defined metadata set
  - Must comply with Greenstone metadata XML file format
    - ❖ Write a script to convert?
  - Copy to /gsdl/gli/metadata directory