

Herbarium Digitization Workshop



iDigBio
Integrated Digitized Biocollections

Database Tools & Techniques

Gil Nelson

September 16-18, 2012

Valdosta State University



iDigBio's Biological Collections Databases, Tools, and Data Publication Portals

<https://www.idigbio.org/content/biological-collections-databases>

(On the Wiki under Database Resources)

If there is something you'd like reviewed, let us know!

Herbarium Digitization Workshop

Spread Sheets: The Scientist's Buddy!

Microsoft Excel spreadsheet showing a table of herbarium records. The columns include RecordID, accid, Taxon, fedstatus, statestatu, country, stapro, coupur, site, township, rangevar, section, quarter, sixteenth, quadnum, notes, collector, geolat, geolon, datum, georefme, geoprecis, link, and technician. The data rows show various plant species like *Sarracenia flava*, *Drosera brevifolia*, *Drosera capillaris*, *Drosera tracyi*, *Pinguicula pumila*, and *Utricularia* species, along with their collection locations and dates.

- Not relational (flat, not normalized)
- Has a mind of its own!
- Data quality issues
- Accepts various data types in same column
- Useful as a tool for download/upload

Herbarium Digitization Workshop

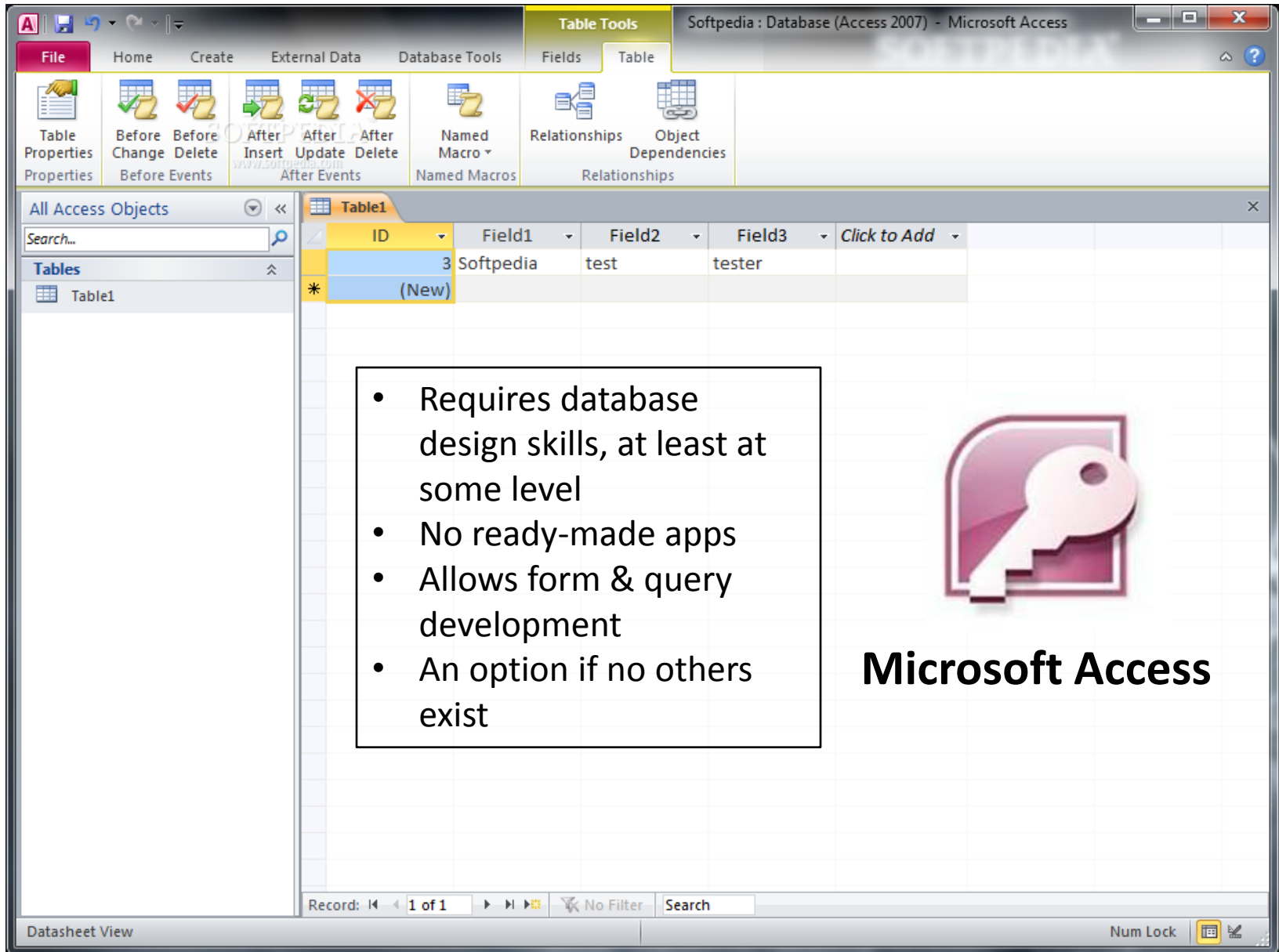


Table Tools

File Home Create External Data Database Tools Fields Table

Table Properties Properties Before Change Before Delete Before Events After Insert After Update After Delete After Events Named Macro Named Macros Relationships Object Dependencies Relationships

All Access Objects


Search...

Tables

Table1

ID	Field1	Field2	Field3	Click to Add
3	Softpedia	test	tester	
(New)				

- Requires database design skills, at least at some level
- No ready-made apps
- Allows form & query development
- An option if no others exist



Microsoft Access

Record: 1 of 1 No Filter Search

Datasheet View Num Lock

Herbarium Digitization Workshop

Botanical Research and Herbarium Management System Department of Plant Sciences, University of Oxford, UK



BOTANICAL RESEARCH AND HERBARIUM MANAGEMENT SYSTEM

BRAHMS is used by research institutions globally supporting collection management in herbaria, botanic gardens and seed banks; taxonomic study; botanical survey; diversity analysis and many further categories of botanical research initiative.

Downloads

[Version 7 update notes](#)

The largest single database to date with ca. 2 million specimens runs at the National Herbarium of the Netherlands, accessed via terminal services by all of the main Dutch herbaria and also published online. The country with the highest number of individual projects (2012) is Brazil.

The BRAHMS project is part of the plant diversity and systematics research group at the University of Oxford.

BRAHMS is a powerful database management system for botanical researchers and herbaria. It provides wide-ranging and innovative functionality to gather, edit, analyse and publish botanical data, optimizing its use for the widest possible range of curation services and research outputs. Find out more in [English](#), [French](#), [Spanish](#), [Portuguese](#) and [Russian](#).



- FoxPro Files
- Mostly European
- Fairly easy to use and setup
- Good training manual
- Links to IPNI

Herbarium Digitization Workshop



THE FLORIDA STATE UNIVERSITY
Biology Department
Robert K. Godfrey Herbarium

SEARCH FLORIDA STATE
FSU Biology Search GO

- Herbarium Main Page
- People
- Projects
- Contact Information
- Loans and Exchanges
- Visits
- Volunteer!
- Friends
- Robert K. Godfrey
- Links
- Database
- Login



The Robert K. Godfrey Herbarium at Florida State University

Florida State University's Robert K. Godfrey Herbarium is a museum-quality collection of over 200,000 plant and microalgae specimens. These document the distribution and natural variation of the 2,400 species of flowering plants, ferns, conifers, and cycads found in northern Florida—one of North America's biodiversity hotspots—and the microalgae of Florida's Gulf and Atlantic coasts. Each plant specimen is carefully identified, pressed, dried, and mounted to archival standards, with accompanying data on where and when it was collected. The specimens are a valued resource to local, state, national, and international biologists studying plant and microalgae systematics, ecology, evolution, biogeography, conservation biology, anatomy, and morphology. New specimens are added to the collection each week. In August 2003, [Austin Mast](#) became the new director of the herbarium. Though now retired, [Loran Anderson](#) remains active in the herbarium. The current curator is Chris Buddenhagen.



Photo: *Hymenocallis henryae*, Photo Credit: Gary Knight

database links on our menu or by clicking [here](#).

In August of 2003 the Robert K. Godfrey Herbarium set up a digital imaging system and SQL database. Currently 68,716 of our 200,000+ specimens have been entered into the herbarium's database. In late 2010, Tall Timbers Research Station's Herbarium began serving specimen data and images on this site as well. Currently, data and images for 10,347 of Tall Timbers Research Station's 10,000+ herbarium specimens are being served. Specimen images and data from both institutions can be searched by following the

The Florida State University gratefully acknowledges funding from the National Science Foundation (Awards 0956372 and 0646222), the Florida Department of Environmental Protection, and the Florida Fish and Wildlife Conservation Commission to support the digitization of specimens in the Robert K. Godfrey Herbarium. The Tall Timbers Research Station gratefully acknowledges funding from the National Science Foundation (Award 0956343) to support the digitization of specimens in its herbarium.

[Biological Science](#) | [Ecology and Evolution](#) | [Austin Mast](#)

STRENGTH SKILL CHARACTER STRENGTH SKILL CHARACTER STRENGTH SKILL CHARACTER STRENGTH SKILL CHARACTER STRENGTH SKILL CHARACTER STRENGTH SKILL CHARACTER STRENGTH SKILL CHARACTER

© The Florida State University, Department of Biological Science, 319 Stadium Drive, Tallahassee, FL 32306-4295
Phone: (850) 644-3700 Fax: (850) 645-8447

“Build Your Own”
OpenHerbarium
at FSU

Herbarium Digitization Workshop

The screenshot shows a web browser window with the URL `openherbarium.bio.fsu.edu`. The page features a dark red header with a navigation menu on the left and a search bar on the right. The main content area has a yellow background with the title "The OpenHerbarium Project at Florida State University" and a brief description of the project. A navigation menu on the left includes links for "Herbarium Main Page", "People", "Projects", "Contact Information", "Loans and Exchanges", "Visits", "Volunteer!", "Friends", "Links", "Database", and "Login".

Open Herbarium

FSU Biology Search GO

Herbarium Main Page
People
Projects
Contact Information
Loans and Exchanges
Visits
Volunteer!
Friends
Links
Database
Login

**The OpenHerbarium Project
at Florida State University**

This is the beginning of the OpenHerbarium project. Software code for front end and back end to run a digital herbarium.

[More Information](#) | [Open Herbarium Project](#) | [Alexander Stuy](#)

© The Florida State University, Open Herbarium Project, 111 Stadium Drive, Tallahassee, FL 32306
Phone: (850) 644-1006 Fax: (850) 644-1006

Herbarium Digitization Workshop



Specimen Database Search

Please cite!

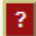
Submit

Undo Changes

Clear All

Search Criteria

Show: Hide:

Institution:	<input type="text" value="Any Institution"/>	
Family:	<input type="text"/>	Spacebar yields full listing. 
Genus:	<input type="text"/>	
Species:	<input type="text"/>	The scientific name (e.g., Pinus palustris).
Common Name:	<input type="text"/>	
Collection Date:	<input type="text" value="="/> <input type="text"/>	(YYYY-MM-DD or MM-DD)
Collection Date:	<input type="text" value="="/> <input type="text"/>	
Collector Name:	<input type="text"/>	
Collector's Identifier:	<input type="text"/>	
Barcode:	<input type="text"/>	
Country:	<input type="text"/>	
State:	<input type="text"/>	
County:	<input type="text"/>	
Nearest Named Place:	<input type="text"/>	
Flowers present:	<input type="text"/>	
Fruit present:	<input type="text"/>	
Habitat:	<input type="text"/>	

Herbarium Main Page

People

Projects

Contact Information

Loans and Exchanges

Visits

Volunteer!

Friends

Robert K. Godfrey

Links

Database

Login

Herbarium Digitization Workshop

- Herbarium Main Page
- People
- Projects
- Contact Information
- Loans and Exchanges
- Visits
- Volunteer
- Friends
- Robert K. Godfrey
- Links
- Database
- Login

Specimen Database Search

Please cite!

Search Criteria

Show: Hide:

Sort By:

Species

To modify your search choose the "Show" radio button above. If you would like to sort your results differently, change the field in "Sort by" above.

Search Results

Results 1-35 of 177 - Page 1 of 6



Herbarium Digitization Workshop

Specimen Database Search

Please cite!

Search Criteria

Show: Hide:

Sort By:

Species

To modify your search choose the "Show" radio button above. If you would like to sort your results differently, change the field in "Sort by" above.

Search Results

Results 1-177 of 177 Page 1 of 1 [Enlarge Map](#)



[Herbarium Main Page](#)

[People](#)

[Projects](#)

[Contact Information](#)

[Loans and Exchanges](#)

[Visits](#)

[Volunteer!](#)

[Friends](#)

[Robert K. Godfrey](#)

[Links](#)

[Database](#)

[Login](#)

Herbarium Digitization Workshop

Table	Action	Records	Type	Collation	Size	Overhead
AccessionRecords		455	MyISAM	latin1_swedish_ci	22.2 KiB	-
CollectorRecords		99,058	MyISAM	latin1_swedish_ci	4.3 MiB	-
Counties		6,138	MyISAM	latin1_swedish_ci	619.5 KiB	200 B
Countries		251	MyISAM	latin1_swedish_ci	18.9 KiB	-
FIPS		3,222	MyISAM	latin1_swedish_ci	146.6 KiB	-
FLGNIS		54,037	MyISAM	latin1_swedish_ci	6.1 MiB	-
FLMA		2,075	MyISAM	utf8_general_ci	233.3 KiB	-
FloridaSections		56,588	MyISAM	latin1_swedish_ci	5.7 MiB	-
geolocaterecords		1,668	MyISAM	latin1_swedish_ci	601.8 KiB	-
geolocateuploads		814	MyISAM	latin1_swedish_ci	268.4 KiB	-
Images		3	MyISAM	latin1_swedish_ci	4.2 KiB	-
iptTestView		~02	View	---	-	-
MasterSpeciesList		10,283	MyISAM	latin1_swedish_ci	1.9 MiB	-
ModRecords		457,124	MyISAM	latin1_swedish_ci	27.3 MiB	-
Notes		9,147	MyISAM	latin1_swedish_ci	848.0 KiB	-
People		3,883	MyISAM	latin1_swedish_ci	467.3 KiB	-
Places		29,514	MyISAM	latin1_swedish_ci	2.0 MiB	-
PLSS		~84,492	InnoDB	utf8_general_ci	12.5 MiB	-
Projects		35	MyISAM	latin1_swedish_ci	7.8 KiB	-
Project_species		903	MyISAM	latin1_swedish_ci	38.4 KiB	-
QuadListFlorida		53,801	MyISAM	latin1_swedish_ci	4.9 MiB	-
Quads		53,765	MyISAM	utf8_general_ci	2.7 MiB	-
SelectedCounties		3,040	MyISAM	latin1_swedish_ci	293.8 KiB	-
SpecimenRecords		79,063	MyISAM	latin1_swedish_ci	33.2 MiB	-
SpecimenRecordsTmp		49,453	MyISAM	latin1_swedish_ci	18.9 MiB	-
SpecimenRecordsTmp2		49,397	MyISAM	latin1_swedish_ci	18.9 MiB	-
SpecimenRecordsTmp3		49,497	MyISAM	latin1_swedish_ci	18.9 MiB	-
SpecimenRecordsTmp4		49,397	MyISAM	latin1_swedish_ci	18.9 MiB	-
testIPTview		~02	View	---	-	-
test_CollectorRecords		80,249	MyISAM	latin1_swedish_ci	3.5 MiB	-
test_SpecimenRecords		66,796	MyISAM	latin1_swedish_ci	30.4 MiB	-
test_VerificationPeople		79,727	MyISAM	latin1_swedish_ci	5.0 MiB	-
test_VerificationTable		74,126	MyISAM	latin1_swedish_ci	8.6 MiB	-
tmplaton		5,981	MyISAM	latin1_swedish_ci	674.6 KiB	-
tmplaton2		0	MyISAM	latin1_swedish_ci	1.0 KiB	-
tmpnative		114	MyISAM	latin1_swedish_ci	6.6 KiB	-
TTRS1FSUCollectorExport		100	MyISAM	latin1_swedish_ci	3.4 KiB	-
TTRS1FSSpecimenExport		100	MyISAM	latin1_swedish_ci	27.3 KiB	-



- Open source
- Apache/IIS
- PHP
- Enterprise level



- Can be installed on a workstation
- Requires database knowledge and skills

Herbarium Digitization Workshop

<http://www.youtube.com/watch?v=UXvzZUlaB7I&feature=plcp>

<http://www.youtube.com/watch?v=faCP15wjc4g&feature=plcp>

Silver

BIOLOGY



Herbarium Digitization Workshop

Data Capture/Enrichment Techniques

(See link on Wiki to Workflow Modules and Tasks: Data Capture)

Keystroking:

- From images
- From specimen sheets
- Long vs. short (skeleton) records
- May be the quickest, most efficient method, especially if recording skeleton records



Optical Character Recognition (OCR)

Scanning electronic images with software designed to extract and make readable embedded text.

OCR Software



ABBYY Finereader 11, Corporate

- **Converts to Word or text, single files or multiple**
- **Provides a user interface**
- **Includes batch processing options**
- **Supports training to specific data sets**
- **Relatively inexpensive**
- **Relatively easy to configure**



tesseract-ocr

Tesseract open source OCR

Originally developed by HP in the 1980s

Now owned by Google

Focus of iDigBio OCR working group

Optical Character Recognition (OCR)

Potential Uses

Ingesting unedited OCR: Specify

Building robust searches of unedited text: VSU

Use as part of other software tools: Apiary, Symbiota



tesseract-ocr

The Apiary Project:

A collaborative workflow for extraction of herbarium label data

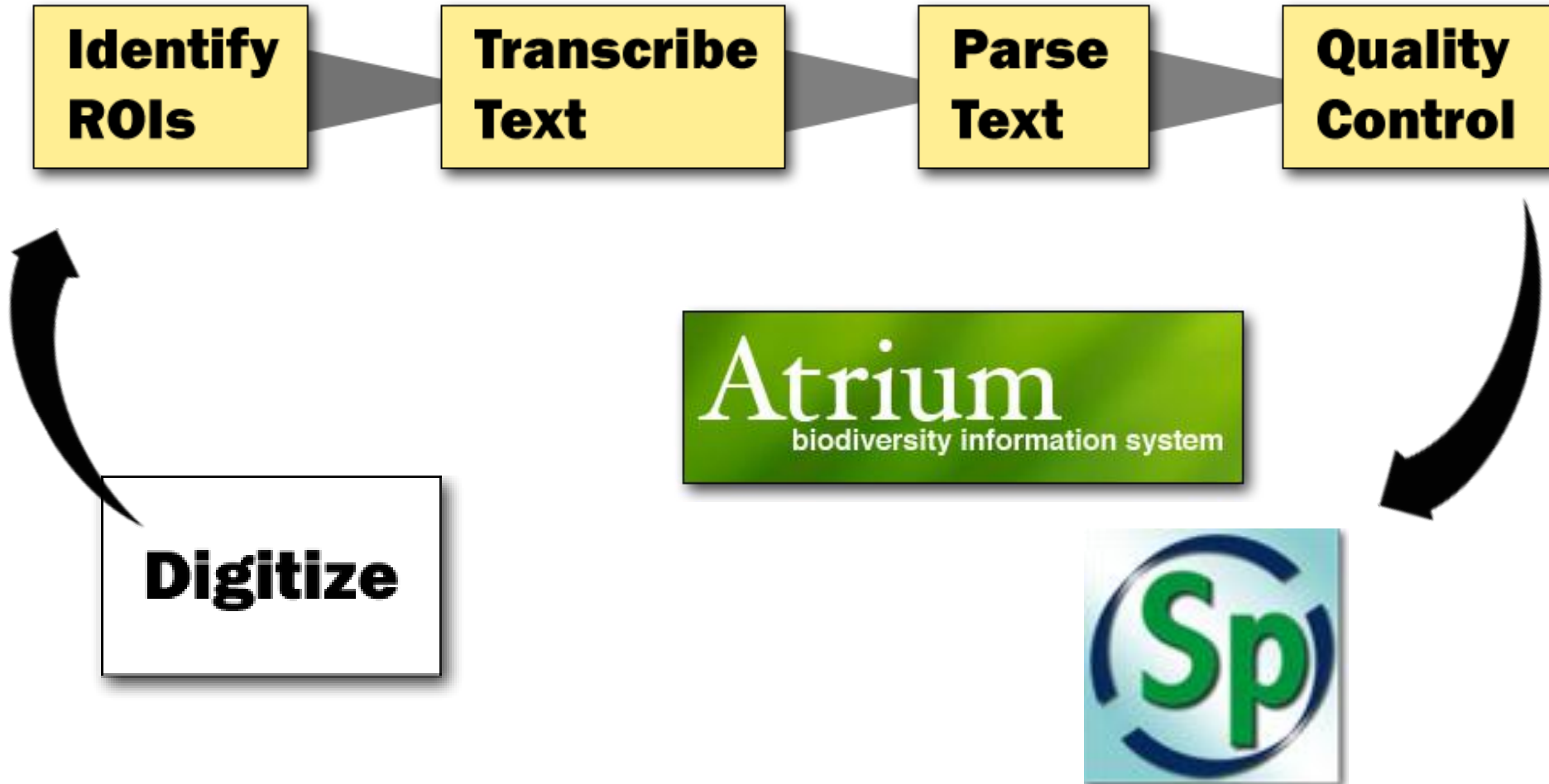
A project of BRIT and UNT's Texas Center for Digital Knowledge





Collector Name	Number	Scientific Name
J. C Taylor	31007	Solidago altiplanities
W. Hess	7283	Spiraea densiflora
H. St. John	14846	Sophora manejarevaensis
C.G. Pringle	13949	Polygala minutifolia
R.F. Hoover	3614	Senecio clevelandii
O. Degener	33,680	Wikstroemia perdita

The Technology and Workflow



Digitize



**Identify
ROIs**

**Transcribe
Text**

**Parse
Text**

**Quality
Control**

Finding Regions of Interest

The screenshot displays the 'Apiary Workflow' web application. The browser address bar shows the URL: http://demo.apiaryproject.org/drupal/modules/apiary_project/workflow/index.php?workflow_id=1#. The navigation menu includes: Home Page | BRIT Apiary | ANALYZE SPECIMEN | TRANSCRIBE TEXT | PARSE TEXT. The main image area shows a botanical specimen with three ROIs: a red box around the primary label, a blue box around a handwritten note, and a green box around a barcode. The primary label text reads: 'PLANTS OF TEXAS', 'COLLECTED FOR THE HERBARIUM OF THE UNIVERSITY OF MICHIGAN BY ROGERS MCVAUGH AND A. M. HARTVILL, JR.', 'Lesquerella Mcvaughiana Rollins (Paratype)', 'Rocky (limestone) slopes, main canyon on northeast side of Sierra Madera, about 25 miles south of Ft. Stockton, Pecos Co.', 'Abundant in shaded canyon bottoms; flowers white, drying purplish; pods ovoid to globose.', 'ROGERS MCVAUGH, NO. 7912', 'APRIL 12, 1947'. The barcode is labeled 'BOTANICAL RESEARCH INSTITUTE OF TEXAS' and '24434'. The right-hand sidebar contains three ROI control panels: 'ROI Type: Primary Label Edit', 'ROI Type: Annotation/Other Edit', and 'ROI Type: Barcode Edit'. The bottom panel includes a 'My Queue' sidebar, a 'My Queue Summary' (5 specimens, 26 ROIs), an 'ap-specimen: Specimen-1 - Analyzing' status panel, an 'ROI Legend' (Primary Label, Annotation/Other, Barcode, Undefined), and 'Session Statistics' and 'End Apiary Session' buttons.

Transcription or OCR

The screenshot displays the 'Apiary Workflow' web application interface. The browser address bar shows the URL: http://demo.apiaryproject.org/drupal/modules/apiary_project/workflow/index.php?workflow_id=1#. The navigation menu includes 'ANALYZE SPECIMEN', 'TRANSCRIBE TEXT', and 'PARSE TEXT'. The main content area is split into two panels. The left panel shows a scanned image of a botanical specimen label with the following text:

PLANTS OF TEXAS
COLLECTED FOR THE HERBARIUM OF THE UNIVERSITY OF MICHIGAN
BY ROGERS MCVAUGH AND A. M. HARVILL, JR.

Lesquerella Mcvaughiana Rollins
(Paratype)

Rocky (limestone) slopes, main canyon on northeast side of Sierra Madera, about 25 miles south of Ft. Stockton, Pecos Co.

Abundant in shaded canyon bottoms; flowers white, drying purplish; pods ovoid to globose.

ROGERS MCVAUGH, NO. 7912 APRIL 12, 1947

The right panel, titled 'TEXT TRANSCRIPTION', contains the following text:

Texas
Lesquerella mcvaughiana
Rocky (limestone) slopes, main canyon on northeast side of Sierra Madera, about 25 miles south of Ft. Stockton, Pecos Co.
Abundant in shaded canyon bottoms; flowers white, drying purplish; pods ovoid to globose.
Rogers McVaugh, No. 7912
April 12, 1947

Below the transcription area is a 'Path: p' field and a 'Save text' button. At the bottom of the interface, there is a 'My Queue' sidebar with a 'My Queue Summary' showing 5 specimens and 26 ROIs, and a 'ap-specimen:Specimen-1 - Transcribing' panel with ROI-83 and a status of 'not started'. On the far right, there are buttons for 'Session Statistics' and 'End Apiary Session'.

Herbarium Digitization Workshop

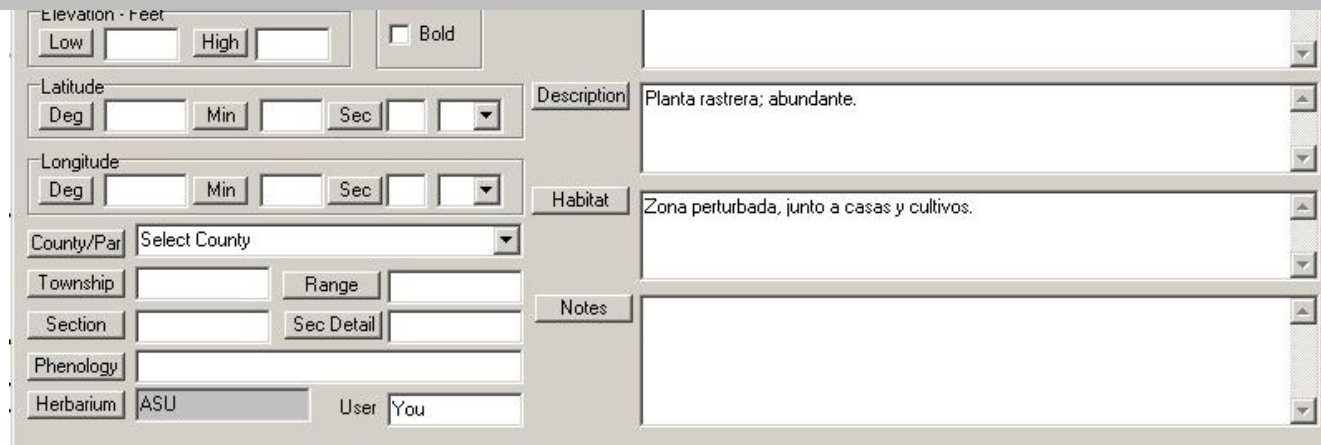
Uploading a CSV in Salix: <http://vimeo.com/42586885>



Salix software download: <http://daryllafferty.com/salix/>

Salix documentation: <http://nhc.asu.edu/vpherbarium/canotia/SALIX3.pdf>

These links are on the Wiki under Database Resources and Tools



Voice/Speech Recognition



Dragon Naturally Speaking
Nuance (now owns IBM's ViaVoice)
Mac & PC
Works better with a single user(?)
~\$200.00 for premium version



Speech to text
Training
BRIT project (Windows API)
Included with Windows

Capturing Bar Code Values

Barcode scanning

- Linear
- 2D
- Avoid data other than catalog number

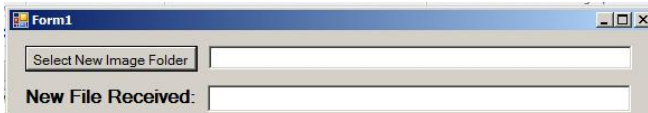


Sync barcode value with camera-named files

Herbarium Digitization Workshop

Capturing Bar Code Values

FNIntercept

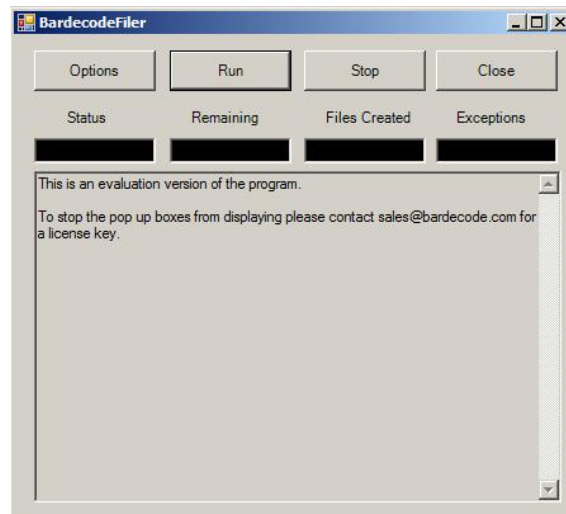


SilverImage



Barcode values can be capture at more than one place in the workflow.

- Pre-digitization curation
- Data capture
- Image capture



Barcodefiler

BCRename

Renaming files to the barcode value



Thank You!

Herbarium Digitization Workshop

Herbarium Digitization Workshop

Herbarium Digitization Workshop