



#pacdigi

The New York Botanical Garden

LICHENS OF VERMONT, U.S.A.

fyrenidium gregatum K. Knudsen & Kocourk.

on Phaeophyscia rubropulchra (Degeer.) Essl.

Essex County: Town of Victory, Victory State Forest,

Umpire Mountain, 44°34'27"N, 71°50'57W, 550 m;

northern hardwoods and forested cold air, granitic talus slope.

16 May 2008

.AA[

Richard C. Harris 54358 1111 III

O107574

Colk,tcd on the 5 Crum Bryophyte /urkho

# Uses for Optical Character Recognition (OCR) Output

Data Discovery and **Doer Happiness!**

Deborah Paul (et al), @idbdeb @iDigBio

Symposium IV: Digitization Workflows, Tools, and Techniques

Biological Digitization in the Pacific, March 23-27, 2014

East-West Center, Bishop Museum, University of Hawaii,

and the Pacific Science Association



Trend

# Minimal Data Capture



- “filed as” name
- higher geography
- barcode
- image
- all sheets in folder get the same initial data
- only the barcode differs

How do we get these minimal records *completed*, or more complete?

# Would you like to...?

- enter records **faster**?
- use the **ditto** feature often?
- find **duplicates** quickly?
- find the labels
- find the labels with lots of **handwriting**?
- create **your own** record sets to transcribe?
  - by collector
  - by country or county
  - by your Great Aunt Susan
  - by taxon
  - by **language**
- create cogent sets to **speed up validation** and **database updates**?
- increase **happiness of transcribers** (paid and volunteer)?

Meaningful datasets make transcription

\* faster

\* less error prone

\* *more fun!*

Got Text?  
Got Handwriting?

*P. alaskanum* *Hulten* No. 7138

National Herbarium of Canada

*Yukon*  
FLORA OF NORTHWEST TERRITORIES

*Papaver nudicaule* L.

Hab. and Loc., Arctic Coast west of Mackenzie River delta:  
Between King Pt. and Kay Pt., 69° 12' N., and 138° to  
138° 30' W.

*Semi-barren ridges*

Collector, A. E. Porsild

July 23-25, 1934

Label

OCR

No. ....2L31.

National Herbarium of Canada  
FLORA OF 'T TERRITORIES

Hab. and Loc., Arctic Coast west of Mackenzie River  
delta:  
Between King Pt. and Kay Pt., 69° 12' N., and 138° to  
138° 30' W.

Collector, A. E. Porsild July 23-25, 1934

*P. alaskanum* *Hultén* No. .... 7138 .....

National Herbarium of Canada

FLORA OF *Yukon* NORTHWEST TERRITORIES

..... *Papaver nudicaule* L. ....

Hab. and Loc., Arctic Coast west of Mackenzie River delta:  
Between King Pt. and Kay Pt., 69° 12' N., and 138° to  
138° 30' W.

..... *Semi-barren ridges* .....

Collector, A. E. Porsild July 23-25, 1934

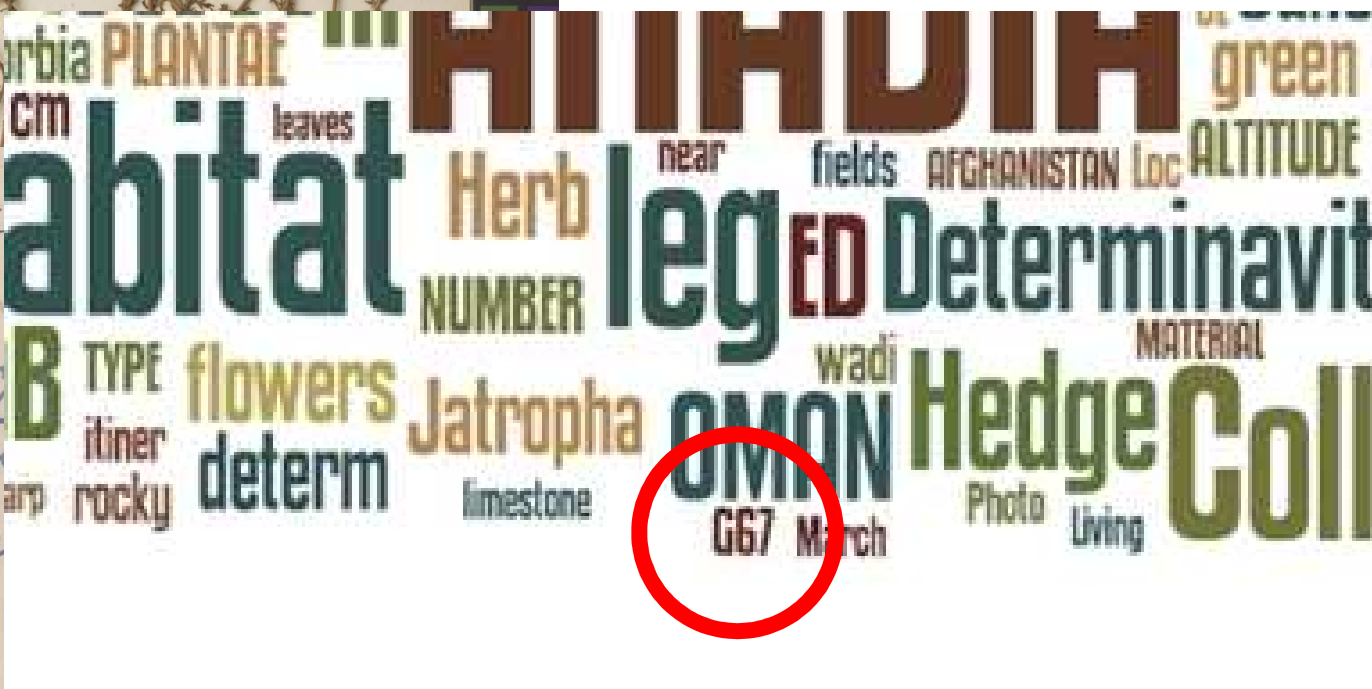
Next imagine  
output from  
1000s of  
labels!



- It's surprising what can be used to help filter specimens – the black art of search terms!



FL  
NAME  
LOCALITY  
HABITAT  
REMARKS  
DATE

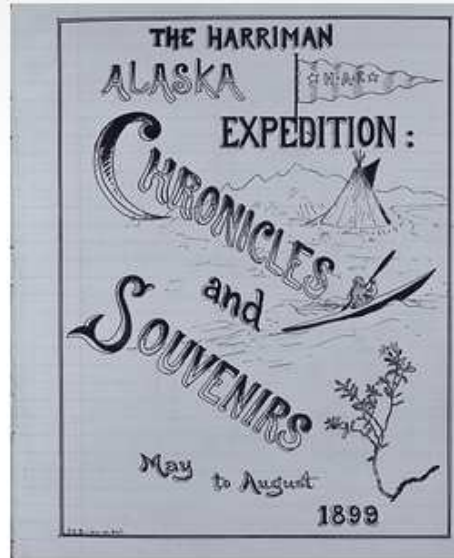


(645) Wt.74459/74 10m(4) 5/53 L.&S. G67.



## Lichen Record Sets

### Harriman Alaska Expedition



Front page of the expedition album created by its members.

In 1899 an interdisciplinary team of scientists and artists traveled from Seattle to Siberia and back exploring the Alaska coastline. Railroad magnate Edward Harriman commissioned the voyage as a hunting trip but it turned into a scientific expedition. The team he brought along to document the trip included some of the best and brightest minds in the country at the time.

Botanists on the trip included:

- William Trelease - director of the Missouri Botanical Garden
- De Alton Saunders - phycologist (algae)
- Frederick Coville - USDA Chief Botanist and Founder of the United States National Arboretum
- Thomas Kearney - research associate in botany at the California Academy of Sciences

Other notable members of the expedition:

- John Muir - naturalist
- Clinton Hart Merriam - zoologist
- John Burroughs - author



Expedition party on a dock in Alaska. Photograph by Edward Curtis. Courtesy of the Missouri Botanical Garden.

### Create Your Own Expedition or Record Set Link

Country  State/Province

Family  Genus

Collection

OCR Fragment

[Click To Go To Records](#)

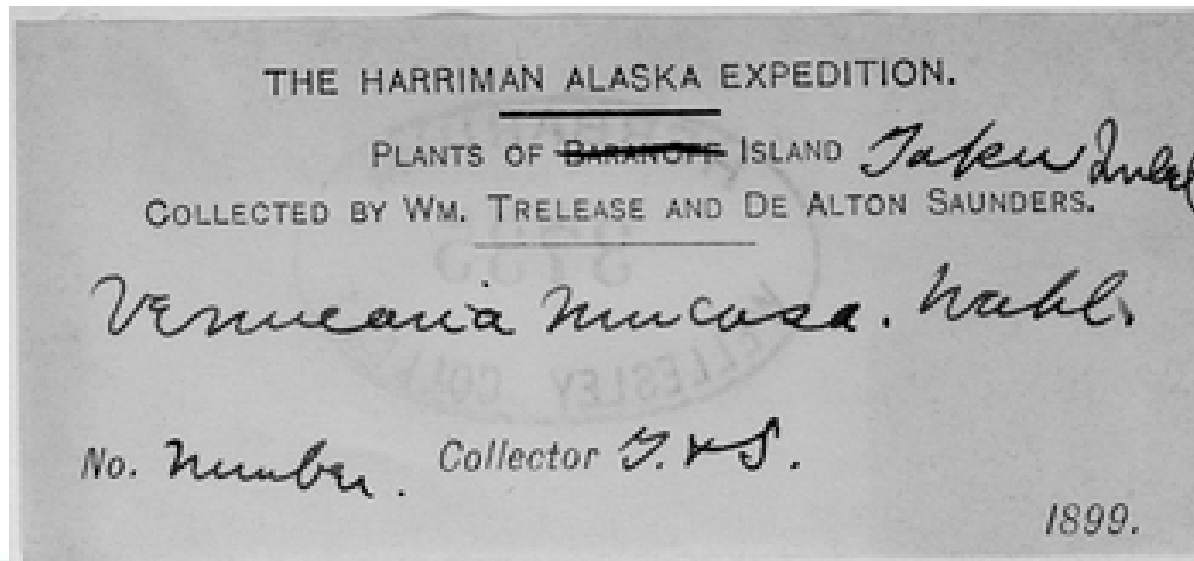
[Reset](#)



# Inside the 1899 Harriman Expedition

Home >> [Crowd Sourcing Central](#) >> Occurrence Record Table View << | 1-416 of 416 records | >>

Symbiote ID	Catalog Number	Family	Scientific name	Author	Date Last Modified
<a href="#">964616</a>	01394706	Verrucariaceae	<i>Wahlenbergiella mucosa</i>	(Wahlenb.) Gueidan & Thüs	2012-03-22 19:08:49
<a href="#">964634</a>	01407608	Nephromataceae	<i>Nephroma antarcticum</i>	(Jacq.) Nyl.	2012-03-22 19:08:49
<a href="#">964643</a>	01407614	Nephromataceae	<i>Nephroma antarcticum</i>	(Jacq.) Nyl.	2012-03-22 19:08:49
<a href="#">964645</a>	01407669	Nephromataceae	<i>Nephroma bellum</i>	(Sprengel) Tuck.	2012-03-22 19:08:49
<a href="#">964649</a>	01407673	Nephromataceae	<i>Nephroma bellum</i>	(Sprengel) Tuck.	2012-03-22 19:08:49
<a href="#">964654</a>	01407671	Nephromataceae	<i>Nephroma bellum</i>	(Sprengel) Tuck.	2012-03-22 19:08:49
<a href="#">964678</a>	01407616	Nephromataceae	<i>Nephroma antarcticum</i>	(Jacq.) Nyl.	2012-03-22 19:08:49



# New York Botanical Garden (NY)

**Occurrence Data**

Collector?  Number?  Date?  Dupes?  Auto search

Associated Collectors?  Verbatim Date?

Exsiccati Title  Number

Scientific Name?

*Note: Full editing permissions are needed to edit an identification*

Country  State/Province  County

Locality

Latitude  Longitude  Uncertainty?  Verbatim Coordinates  Tools

Elevation in Meters  -  Verbatim Elevation

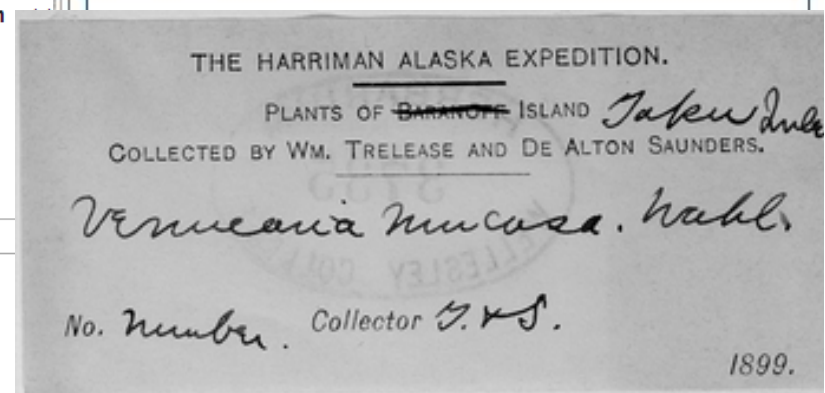
Habitat

Substrate

Notes

Status Auto-Set:  ▼

## Label Processing



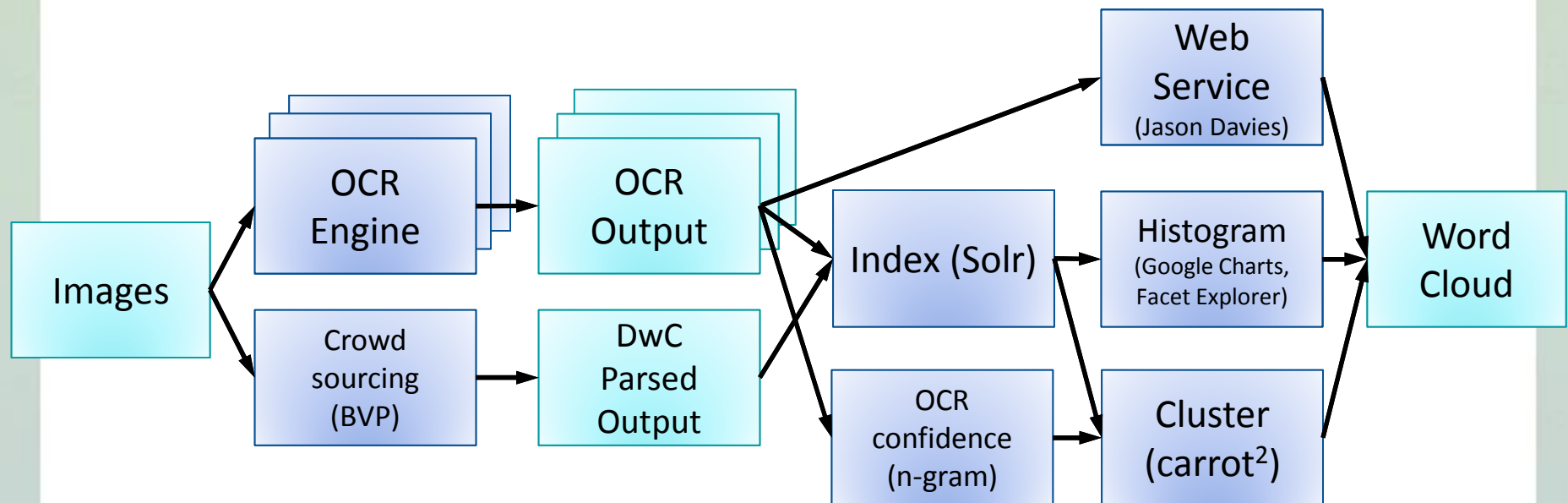
Notes:

Source:

1 of 1

Some work from the recent iDigBio CITSCribe Hackathon

# Overall Word Cloud Workflow



Google Charts: <http://developers.google.com/chart/interactive/docs/gallery>

N-gram: <http://github.com/idigbio-citsci-hackathon/OCR-Error-Estimation>

Facet explorer: <http://github.com/idigbio-citsci-hackathon/facet-explorer>

Jason Davies WC: <http://www.jasondavies.com/wordcloud/>

Apache Solr: <http://lucene.apache.org/solr/>

carrot<sup>2</sup>: <http://project.carrot2.org/>

# Word Clouds using N-gram Scoring, Faceting, Solr + Carrot<sup>2</sup>



[Hide options](#)

Download  Cluster with

Title field name

Summary field name

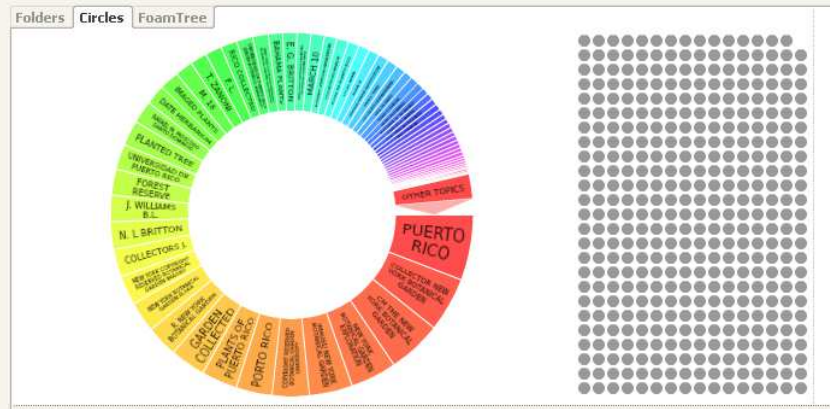
URL field name

ID field name

Read Solr clusters if present

Use highlighter output if present

[Hide advanced options](#)



6 7 8 9 10 " XT "

The New York copyright reserved botanical Garden  
 NEW YORK BOTANICAL GARDEN  
 00040931  
 BOTANICAL  
 GARDE\*\*  
 REPORT ANY REIDENTIFICATION OF THIS VOUCHER  
 TO THE INSTITUTE OF ECONOMIC BOTANY, NY  
 The Dev York Botanical Garden INSTITUTE OF ECONOMIC BOTANY Plants of Conaonfreath of Doiinica  
 APIACEAE  
 Eryngium foetidum L.  
 I >AMS. T.B«Oc«n  
 WEST INDIES, Dominica, Carib Territory,  
 Bataka, 300m up a trail heading W of feeder  
 near field in Galaback. 15°13'N, 61°16'W.  
 Flat ground,  
 road. -----  
 100m. Maintained in a garden, shady, moist soil.  
 Herb to 0.3m with bitter odor; leaves 6-8 at base; flowers with sharp spiked petals.  
 n.v.: Chadron Beni [Crcole-Patois].  
 DSE: Untitnssive lvs. Infnsion with Pluchea inti tussive for colds ileavesfflowers | ade with 3 leaves and 1 flower >C| Sample codes  
 )  
 Janes Higgins 12 with Prosper Paris  
 synphytifolia Drink Infusion  
 August 8, 1992  
 Fieldwork supported by the national Cancer Institute and Berck Research  
 laboratories  
 00040931

**Legend** - Level of confidence that token is an accurately-transcribed word  
 ■ extremely low ■ very low ■ low ■ undetermined ■ medium ■ high ■ very high

Top 1000 results of about 5000 for \*:

1 <http://ammatsun.acis.ufl.edu:5901/carrot2-webapp-3.8.1/herballsilvertrigram/00040931.txt.html>

6 7 8 9 10 " XT " The New York copyright reserved botanical Garden NEW YORK BOTANICAL GARDEN 00040931 BOTANICAL GARDE\*\* REPORT ANY REIDENTIFICATION OF THIS VOUCHER TO THE INSTITUTE OF ECONOMIC BOTANY, NY The Dev York Botanical Garden INSTITUTE OF ECONOMIC BOTANY Plants of Conaonfreath of Doiinica APIACEAE Eryngium foetidum L. - I >AMS. T.B«Oc«n WEST INDIES, Dominica, Carib Territory, Bataka, 300m up a trail heading W of feeder near field in Galaback. 15°13'N, 61°16'W. Flat ground, road. ----- 100m. Maintained in a garden, shady, moist soil. Herb to 0.3m with bitter odor; leaves 6-8 at base; flowers with sharp spiked petals. n.v.: Chadron Beni [Crcole-Patois]. DSE: Untitnssive lvs. Infnsion with Pluchea inti tussive for colds ileavesfflowers | ade with 3 leaves and 1 flower >C| Sample codes: IMOCUHO03-O) Janes Higgins 12 with Prosper Paris synphytifolia Drink Infusion August 8, 1992 Fieldwork supported by the

# Imagine Integration with current software

GOT IT! LET ME TRANSCRIBE

GOT IT! LET ME CHOOSE MY TRANSCRIPTION GROUP



Choose your group:

Country

- Use for initial sort or validation



# Working Group Collaboration

## aOCR + DROID1 OCR workflow, ...

Task ID	Task Name	Explanations and Comments	Resources (human, physical tool, software)
		Get help from folks already doing OCR (iDigBio aOCR wg!)	
	Select OCR software	<p>Factors to use in making selection might include:</p> <ul style="list-style-type: none"> <li>• Cost</li> <li>• Batch processing capability</li> <li>• Output types (XML, Text, etc.)</li> <li>• Language Support</li> <li>• SDK availability</li> <li>• OS platform (Windows, Linux, Mac etc)</li> <li>• APIs</li> <li>• Authority File Support</li> <li>• Watched Folders - allowing use across multiple users.</li> </ul> <p>It will be helpful to review the next task to ensure adequate hardware/OS support.</p>	<p>Available OCR software and tools include:</p> <p>Proprietary</p> <ul style="list-style-type: none"> <li>• ABBYY Finereader</li> <li>• ABBYY Finereader Pro Edition</li> <li>• ABBYY Finereader Corporate Edition</li> <li>• ABBYY Recognition Server</li> <li>• OmniPage</li> <li>• Prime Recognition</li> </ul> <p>Open Source</p> <ul style="list-style-type: none"> <li>• Tesseract</li> <li>• J(G)OCR</li> <li>• OCRopus</li> </ul>
	Install OCR software	<p>Hardware/OS issues</p> <ul style="list-style-type: none"> <li>• Ensure that hardware meets requirements of the OCR software</li> <li>• Ensure sufficient memory capacity</li> </ul>	

- Setting up OCR
- Running OCR
- Machine Learning
- Natural Language Processing

# OCR use, a bit more...

Got Text?  
*Got Handwriting?*

- aOCR WG, JRA Synthesys3, ...
- Who is using OCR output? ...some examples
  - New York Botanical Garden
    - Caribbean [Field Notebook](#) Project
  - Royal Botanic Garden Edinburgh
  - Kew
  - Lichen, Bryophyte , Climate Change (LBCC) TCN
  - Biodiversity Heritage Library
- N-gramming (confidence scores)
- parsing algorithm optimization
- user-interface group
- exemplar ML and NLP workflows

Work presented here  
made possible by many  
and especially...

# Mahalo!

- Andrea Matsunaga, Researcher, iDigBio
- Miao Chen, Indiana University, Data to Insight Center
  - Jason Best, Botanical Research Institute of Texas
- Sylvia Orli, IT Head, Smithsonian Botany Department
  - William Ulate, Technical Director, BHL
  - Reed Beaman, Informatics Specialist, iDigBio
- iDigBio Augmenting Optical Character Recognition (aOCR)  
Working Group

