

Reaching Consensus in Crowdsourced Transcription of Biocollections Information

Andréa Matsunaga (ammatsun@ufl.edu), Austin Mast, and José A.B. Fortes

10th IEEE International Conference on e-Science October 23, 2014 Guarujá, SP, Brazil



iDigBio is funded by a grant from the National Science Foundation's Advancing Digitization of Biodiversity Collections Program (Cooperative Agreement EF-1115210). Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the National Science Foundation.



Background and Motivation

- Estimated 1-billion biological specimens in the US
- iDigBio + Thematic Collections Network (TCN) + Partners to Existing Network (PEN)





Background and Motivation

- TCNs and PENs performing digitization work
- Generating images, transcribing information about what/when/where/who



- If digitization took 1 second and if we performed in sequence:
 - 1,000,000,000 seconds > 30 years
- Parallelism:
 - Crowdsourcing!!



Crowdsourcing Transcription Projects

- NotesFromNature (<u>http://www.notesfromnature.org/</u>)
 - Zooniverse platform



- Once the user selects the region with the label, s/he can start transcribing and parsing information to a number of predefined fields
- For a requester, a pre-defined number of transcriptions are returned



Crowdsourcing Transcription Projects

ALA (<u>http://volunteer.ala.org.au/</u>)
 – Platform: Grails

	ľ	T	in the second			01248676 Taxa: Croton pallidulus	Copy from previous task	
				1. Museum details Shri				
		all		4		USNH 25 Number Cultivated	73221	
	ר	Ž	Hardware der Berleherten der B	20100 N. ROAD,		Sheet	of (if noted	d)
			Acta a state 12.00	and the same	4			
2. Collection details Collector(s)	Reitz			2	3. Location details Verbatim Locality	Rio Coitinho, guarapua	ava, Paraná	?
2. Collection details Collector(s)	Reitz Klein				3. Location details Verbatim Locality	Rio Coitinho, guarapua	ava, Paraná	?
2. Collection details Collector(s)	Reitz				3. Location details Verbatim Locality State/Province/Territory	Rio Coitinho, guarapua Paraná	ava, Paraná	?
2. Collection details Collector(s)	Reitz				State/Province/Territory Country	Rio Coitinho, guarapua Paraná Brazil	ava, Paraná	?
2. Collection details Collector(s)	Reitz Klein				3. Location details Verbatim Locality State/Province/Territory Country Elevation	Rio Coitinho, guarapua Paraná Brazil 1000 - To	ava, Paraná //	?
2. Collection details Collector(s)	Reitz Klein 17657 (from)	15	12 196 MM Vog		State/Province/Territory Country Elevation Verbatim Latitude	Rio Coitinho, guarapua Paraná Brazil 1000 - To D°M'S″ V D	ava, Paraná	?

- User zooms in to read the label and parse to the custom predefined terms
- Single worker followed by expert approval



Crowdsourcing Transcription Projects

Symbiota (<u>http://lbcc1.acis.ufl.edu/volunteer</u>)
 – Platform: PHP

Occurrence Data			🖲 Med Res. 🔍 High R
Collector ? /1 899 Associated Collectors ? Exsiccati Title	Number ? Date ?	Long Form ≤≤ Dupes? Auto search Number	HARRINAN ALASEA EXPEDITION. afferioria a arte com 2/ 7.
Scientific Name ? Nephroma antarcticum Note: Full editing permissions ara Country Locality Locality Latitude Longitude Latitude Longitude Levation in Meters Ve Habitat Substrate	reeded to edit an identification VProvince County incertainty ? Verbatim Coordinates vrbatim Elevation	Tools	Clara E. Cummings Harbarium of Wellestoy College (WELC) Deposited at NV in 1988
Save Edits Status Auto-Set: Pending F	teview T	HARRIMAN ALASK be) (/* FREDERICK V. CC ¿L/2-1r. COLLECTORS. . W <- 'THOAIAS' Nephroma arcti Lichen substan Det. G. M. Wet The Lichen Gen America. Clara E. Cummin Notes: Source: ABBYY:2013-02-C Save OCR Edi	A EXPEDITION. OVILLE I r ?? " " - f /1 899 f H. KEARNEY, JR. J U^J/O . cum (L.) Torss. ces extracted: more, 1958, no. H/Q, us Nephroma in North and Middle ngs Herbarium D9 LBCC Parser 1 of 1

- Ability to OCR and parse data
- Single worker
 followed by
 expert
 approval



2

The transcription task





Habitat and description

Collected by

7



Proposed Consensus Approach

- Goal:
 - Reach consensus with minimum number of transcriptions
- Method:
 - Control the number of workers per task
 - Apply lossless and/or lossy algorithms per field





Lossless normalization algorithms

Code Lossless functionality

- b Removes all extra whitespaces
 - Apply specific transformation functions on a per-field basis (e.g., to normalize section/township/range, proper names, and latitude/longitude)
 - t Apply specific translation tables on a per-field basis to expand abbreviations (e.g., hy, hwy, and hiway to highway) or to shorten expansions (e.g., Florida to FL)



Lossy normalization algorithms

Code	Lossy functionality
W	Approximate comparison by ignoring all whitespace (e.g., "0–3" is equivalent to "0 – 3")
С	Case insensitive approximate comparison (e.g., "Road" and "road" are considered equivalent)
S	Consider two sequences equivalent when one is a substring of another or one sequence contains all words from the other sequence
р	Punctuation insensitive approximate comparison (e.g., separation of sentences with comma, semi-colon or period are considered equivalent)
f	Approximate fingerprint comparison ignoring the order of words in sentences
I	Approximate equivalency when sequences have Levenshtein distance within a configurable threshold (I2 indicates a maximum distance of 2)



Alternative voting and consensus output

Code Voting and consensus

- Consensus is reached when there is a single group (set of matching answers) that has the most votes, instead of requiring strict majority vote among all answers
- a Outputs best available answers when consensus is not achieved

n=4

Majority voting requires 3 matching answers $\left[\frac{n}{2}+1\right]$ \rightarrow Consensus not reached v: Blue set has most votes

 \rightarrow Consensus reached



Experimental Setup

- Notes from Nature
- Herbarium specimens from a single institution
- Configured to require 10 workers per task that yielded close-tolinear distribution due to empty tasks and skips
- 23,557 total transcriptions completed by at least 1,089 distinct workers

Field	Uniq#	Field	Uniq#
Country	39	Location	16,161
State/Province	288	Habitat	15,134
County	655	Collected by	3,380
Scientific name	5,941	Collector Number	3,665
Scientific author	4,088	Collection date	2,287



Transcriptions Performed by Individual Workers







- Full consensus improvement from 1.8% to 84.2%
- Confirms intuition that *country, state/province, county, collector number* and *collector date* are "easy"
- Lossless algorithms have small impact except for *scientific author* and *collected by*
- Being insensitive to whitespace, punctuation, and letter case as well as considering substrings, provide the greatest improvement when including lossy algorithms in "difficult" fields



Additional Verification



- Consensus reached mainly with lossless algorithms
- Low percentage of blank responses



Consensus Accuracy

- 300 labels were transcribed by an expert
- Expert had access to information across labels that workers did not have
- Effect on the overall accuracy is minimal
 - 0.9% drop for accepting cases that did not reach consensus
 - 2.3% drop for minimizing the needed workforce

$\frac{2 * matches}{len(s1) + len(s2)}$







Workforce Savings

- Can be as high as 55.8% for the distribution in the studied dataset
- Good for a fixed setting: 3 workers
- Controller advantage: 3 workers is just a good average, and our results show that there are cases where up to 9 workers were needed to reach overall consensus



Task Design Improvement recommendations

- Restricted user interface improves consensus and accuracy
 - Caution to not restrict valid scenarios (e.g., partial dates, range of dates)
 - Broadly defined fields could be engineered to capture more parsed data (e.g., lat/long, TRS)
- Exploring relationships between tasks
 - Enter collector number and collection date first
 - Update related record to have the same information
- Additional training
 - Problems pronounced in separating scientific names from its authorship

Kral

Godfrey

Cacalia lanceolata Nutt., var. elliottii (Harper)



Additional Improvements to Consensus Algorithms

- Code is modular and open; thus, opportunity for:
 - Custom dictionaries could be applied (general dictionaries led to a high number of false positives due to the amount of abbreviation and names)
 - Scientific name parsers
 - External contributions
 - <u>https://github.com/idigbio-citsci-hackathon/CrowdConsensus</u>
- Merge matched outputs after lossy algorithms are applied
 - R. E. Perdue, Jr. and K. Blum
 - Re Perdue Jr and K. Blum
 - R. R. Perdue Jr, K. Blum
- Additional validation across fields (consistency)
- Apply consensus controller on a per-field basis



Recommendations Beyond Crowdsourcing

- Leveraging and improving:
 - Optical Character Recognition (OCR)
 - Natural Language Processing (NLP)
- 2-way street scenarios:
 - Use crowdsourcing to select clean text for OCR
 - Use even poor OCR to guide tasks to the right crowd by creating clusters of tasks
 - Use NLP to parse verbatim data from the crowd
 - Improve NLP and OCR training with additional data from the crowd



Additional parsed data

PhiloJIVE: http://phylojive.acis.ufl.edu

iDigBio Portal PhyloJIVE Home	OpenTree 👻	Sample Trees	Tutorial	Research Tools	
iDigBio Portal PhyloJIVE Home Existing Tree: Helianthus Helianthus tree by Joe Miller Select another tree: Helianthus T • Click the top button to get the navigation aid • Click nodes to get maps and external services • Try choosing characters (if available) to plot on the tree; • Align-names feature; search; set-root; rotate, etc.	OpenTree -	Sample Trees	Tutorial	Research Tools	Atus: done
Create New Tree	- Helia	Helianthus anon Helianthus deserver the anthus bolander Attl. anthus exilis	aflet I Man data Å coverLife Wo as of Living A 조산자조 Helian	© OpenStreetMan contributor 1d Map ust. Spat. portal thus petiolaris	······································



Questions?







- facebook.com/iDigBio
- twitter.com/iDigBio
- vimeo.com/idigbio



V

idigbio.org/rss-feed.xml

webcal://www.idigbio.org/events-calendar/export.ics

Thank you!



iDigBio is funded by a grant from the National Science Foundation's Advancing Digitization of Biodiversity Collections Program (Cooperative Agreement EF-1115210). Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the National Science Foundation.

