# Towards user-definable, semiautomated workflows for curating biodiversity data

P.J. Morris R.A. Morris B. Ludäscher D. Lowery J.A. Macklin T. Song T. McPhillips J. Hanken









Agriculture and Agri-Food Canada

# **Webinar Outline**

• Overview (Bertram)

. . .

- Part 1: Presentation [45 min]
  - 1.1 Example QC Spreadsheet & FilteredPush (Bob)
  - 1.2 Data Cleaning for Natural History Collections (Paul)
  - 1.3 Intro to current FP-Akka/post-proc. tools (Paul)
  - 1.4 Scientific Workflow Automation (Bertram)
- Q&A/Transition/Software check (David, Bob, Paul)
- Part 2: Demo & Hands-on (optional) [55 min]
  - Run FP-Akka workflow –w DwCa –a COL

#### **1.1 Example QC Spreadsheet**

Bob (~ 10 min)

# 1.2 Data Cleaning for Natural History Collections

Paul (~ 15min)



With tests for Fitness for Purpose



# Natural Science Collections Data Fit for what purpose?

- Classical Use by Taxonomists visiting collections, and collections managers.
- Electronic Distribution: Many new uses
  - Species range modeling
  - Modeling effects of climate change
  - Many new uses depend on
    - (1) Good georeferences (where).
    - (2) Good identifications (what taxon)
    - (3) Good collecting event dates (when)









Harvard University Herbaria Plant specimens (now fixed)

#### **Internal Inconsistency**

#### Latitude and Longitude: Somewhere in China Country: United States of America



Harvard University Herbaria Plant specimens, color coded by country

000

# Missing Data

Distribution of Rubus (Raspberries, blackberries) Herbarium Records in GBIF, Jan 2009







Internal Sample Look for Outliers Completeness Check Clean Data With Data **External Authorities** Patterns Outliers Measure

Improved Authority Files Improved User Interfaces





#### **Controlled Data Capture**

ld: Formicidae:Myrmicina	e:Pheidole				
Label data: Verbatim Loca	ality Text				
Verbatim Locality	de Kelifely; Madagascar		0		
Elevation	-		▼ 💽		Causse do
Country:	0	Georeference			and boe de
State/Province:				15	Kelifel: 20
Higher			0	10.57	
Geography					30 XT 1:
Specific Locality:			O No Data	a	an, litte
					dry for
another Collecto	or Name			1000	ary forest
collector					A REAL PROPERTY OF THE PARTY OF
Verbatim Date:				100	
ISO Date yyyy/mm/dd:				Sec.	
Date Collected:					
Collection:					
CollectingMethod.		er en			I CONTRACTOR OF AN INCOMENT
Other Numbers				•	
Notes				$\Theta$	
Specimen Notes: speci	imennotes		0	100	LOOAD
Features:		00			MADAGASCAK
Habitat:			0	1939	THUMAN
Microhabitat micro	habitat		0		A Pevrieras
AssociatedTaxon:	0		0	1000	A. 101-
ValidDistributionFlag: 🗹 🕔	)				Garast humus
Part	Preparation	Count	Modifier	_	IOTESC name
whole animal pi	nned				11:440x 1974
Attribute	Value	Remarks	Add Attributes		&litter 1),
caste w	orker		Edit		
	Add A Speed	mon Part			

0

SB

CC

### **Controlled Data Capture**

Large controlled vocabularies: Geographic, Taxonomic, etc. authority files, with lookups.

Higher Geography

Dominican Republic, La Altagracia Province, Greater Antilles, Hispaniola

Change Edit

#### Locality

Specific Locality

Loma Quita Espuela Scientific Reserve

Class Class Unit

WIIN. ER	ev.	Wax. El	ev. Ele	v. UI	пц
268	ТО	268	m	-	
Min. De	pth.	Max. D	epth. D	epth	Unit
	ТО		1	Ċ	-

Short controlled vocabularies, with pick lists.

#### Locality Remarks

Not Georeferenced Because





Deelumerke										
Bookmarks X	Welcome Data Query Attachments		Sp7demofish KUFishvoucher 🔻 L Search	_og out						
> Bookmarks Tooldar > E Bookmarks Menu z 📾 Unsorted Bookmarks										
DEMO: Specify Integration	New Collection Object: ######### previous next									
	Cat # ######## Accession #	アキャ Prev/Ex	rch # 1234							
	Cataloger Cat Date MM/DD	)///// <b>Ö</b>								
	Determinations (1) Add Delete									
	Taxon Tetraodon mola	/+/	₽ □In Question							
	Preferred Mola mola	^altCatalogNumber	1234							
	Taxon	Found one match.								
	Determiner Chuck McCallum	au	1004							
	Largest bopy fish	taxon	Tetraodon mola							
	Remarks	preferredTaxon	Mola mola							
		determiner	Chuck McCallum							
	Field 1234: Adriatic Sea		1234: Adriatic Sea							
	No:1 Calify	size	.8m long / 2.5m fin-to-fin							
	Preparations Add	weight	1000kg							
	nothing here	remarks	Largest bony fish							
	Tissues: Cylection Object Collection									
	Col Obj Attribute Digte		Т							
	Size 1.8m lon, Y 2. Sex Male Veight 100	Ok								
	Attachments									
	Image: Comparison of the second secon									
	Inventory Edited by									
	FP-DataEn	try	FP-DataEntry							
			Ouers Cornico							
	UT SETVIC	C	Query service							
				) (1)						









# 1.3 Introduction: Current FP-Akka Workflows & Post-processing software

Paul (~ 5min)

#### Sort: Look at top and bottom

\_ = ×

```
File Edit View Terminal Help
mysgl> select recordedby from omoccurrences where recordedby is not null order by recordedby desc li
mit 5;
 recordedby
 ίο,-C. CytJC. STERLING
  {JViO- 'n
  Östman, Magnus
  Österlind, F., O.
5 rows in set (0.00 sec)
mysgl> select recordedby from omoccurrences where recordedby is not null and trim(recordedby) <>
order by recordedby limit 10;
  recordedby
   Conan Taylor
  Mary J. Allen
  Mary J. Allen
  "Reliquiae Tuckermanianae".
  "Reliquiae Tuckermanianae".
  %. D%%%
  %. D%%%
  %. White & C. W. Laskowski
  %. Wright
  %. Yuyeros
10 rows in set (0.00 sec)
mysql>
```

#### **Examine rare values**

•select count(\*), country from omoccurrences group by country order by count(\*), country asc limit 10;

cour	nt(*)	country
	1	0
	1	1998
	1	a
	1	ABKHAZIA (GEORGIA)
	1	ABYSSINIA
	1	Africa Occidentalis [SÜo Tom_ and PrÕncipe]
	1	Alaska
	1	ANDORRA
	1	Arctic
	1	AUSTIRA

# Look for patterns (in aggregated data)

recordedBy	fieldNumber	eventDate
S. Hammer	7540	1999-01-00
S. Hammer	7476	1999-01-00
S. Hammer	7279	1999-01-00
S. Hammer	8237	2000-01-00
S. Hammer	8321	2000-01-00
S. Hammer	8204	2000-01-00
S. Hammer	7843	2000-07-00
S. Hammer	7851	2000-07-00
S. Hammer	7853	2000-07-00
S. Hammer	7849	2000-07-00
S. Hammer	7904	2000-08-09
S. Hammer	7909	2000-08-09
S. Hammer	8002	2000-12-00
S. Hammer	8028	2000-12-00

Integer field number that gets larger over the life of the collector got smaller here.





| sign changed coordinates are on the Earth's surface. | Coordinates not inside country. | transposed/sign changed coordinates to place inside the provided Country UNITED STATES | Transposed/sign changed coordinates are near (within 200.0 km) georeference of locality from the Geolocate service.



### Check names against nomenclators and taxonomic authority sources

Placopecten magellanicus

WAS: Gmelin, 1791; CHANGED TO: (Gmelin, 1791)







Internal Sample Look for Outliers Completeness Check Clean Data With Data **External Authorities** Patterns Outliers Measure

Improved Authority Files Improved User Interfaces



#### **1.4 Scientific Workflow Automation:**

- FilteredPush curation workflows
- Towards Kurator/P

#### Bertram (~ 10min)

#### What problems are we trying to solve?

- Detect and flag data quality issues
- Repair if possible
  - ... ask human curators as needed
- Keep track of **provenance** 
  - (semi-)automatic repairs
  - human curators' edits
- Employ workflow (semi-)automation
  - Scientific workflow systems:
    - Kepler/COMAD, Restflow, Galaxy, Biovel/Taverna, Argo, VisTrails, ...
  - Related technologies
    - Akka parallel execution platform
    - Script-based automation (e.g. Python, R), digital notebooks (iPython)

akka

# **Curation Workflows Users**

#### Collection Managers

- ... who are managing the collections databases
- Can run curation workflows periodically
  - ... in the presence of new data and/or new curation services
- (Biodiversity) Researchers
  - To perform an analysis in the presence of (partially) dirty data, researchers need to
    - Clean or fix dirty data
    - Throw out unfixable data
  - Reporting back to the collection managers (.. push)

#### FilteredPush and Kepler Curation Workflows



Dou, Lei., G. Cao, P.J. Morris, R.A. Morris, B. Ludäscher, J.A. Macklin, J. Hanken. 2012. Kurator: A Kepler Package for Data Curation Workflows, Procedia Computer Science, 9:1614-1619, doi:10.1016/j.procs.2012.04.177

- Today: FP-**Akka** workflow:
  - Validation of (1) SciName; (2) GeoReference; (3) CollectionDate
- Oh, and why workflows?? ASAP!

# **Scientific Workflows: ASAP!**

#### Automation

- wfs to **automate** computational aspects of science
- **Scaling** (exploit and optimize *machine cycles*)
  - wfs should make use of parallel compute resources
  - wfs should be able handle large data

➔ Akka dataflow platform

- Abstraction, Evolution, Reuse (human cycles)
  - wfs should be easy to (re-)use, evolve, share

#### Provenance

- wfs should capture processing history, data lineage
- → traceable data- and wf-evolutior
- ➔ Reproducible Science



akka







Trident

Workbench

**XALON** 

# So many choices ...

- Why not just use system X?
  - Askalon, Kepler, Taverna, Trident, Triana, ...
- Works well for:
  - custom libraries and
  - parameterized workflows



- But there are challenges:
  - ... new actors/functionality (extensibility for mere mortals!)
    - powerful but also complex underlying MoCs ("PhD effect")
  - … adopting system X = learning new language X
    - (tool makers already "speak" languages != X)
  - ... sci-wf systems emphasize process
  - ... but curation about data!

# **Scientific Workflows & Scripts**

- New Kurator Approach
  - Custom GUI (for tool users) [Kurator Phase 2]
    - for custom workflows built from a small library
    - technology/system agnostic
    - include data viewer

- Kurator/P (for tool makers) [Kurator Phase 1]

- empower makers of curation tool
- scientific workflow techniques for the rest of us
- meeting (script-, batch-) programmers half way
- easy integration with Custom GUI (Phase 2)

#### • Spectrum of curation technologies:

Scientific workflows ... YesWorkflow ... anyScript

#### YesWorkflow Example (EnviRecon.org)

 Python, MATLAB, Bash, ... R scripts revealed as workflows with YesWorkflow!

Kyle B., (computational **R**-)archaeologist:

"It took me about 20 minutes to comment. Less than an hour to learn and YW-annotate, all-told."







### 2 FP-Akka Workflow Demo & Hands-On Part (optional)

#### Paul (~ 55 min)

#### iDigBioWebinar May2015

Resources for iDigBio Webinar 28 May 2015 demonstrating an FP-Akka workflow for data quality control. Adobe Connect webinar site & 2-4 pm EDT May 28 2015.

Example result: QC report spreadsheet: Output\_demoset.xls & Look around in this after reading the brief text in its introductory sheet. We'll do so quickly at the beginning of the webinar.

Example input: occurrence.txt file from a DwC Archive: <u>occurrence\_demoset.txt</u> &. Running the software against this gives rise to the above human-centric data cleaning outcomes. The webinar will show you how, following the instructions below.

Contents [hide]

1 Preparation

2 Mechanics of Running the software

2.1 Option 1: From a command prompt

2.2 Option 2: By Editing properties files

3 How to tell the software what to do

4 References

http://wiki.datakurator.net/web/iDigBi oWebinar\_May2015

# Acknowledgments



• NSF-DBI Filtered Push: Continuous Quality Control for Distributed Collections and Other Species-Occurrence Data (ending)

• NSF-DBI Kurator: A Provenance-enabled Workflow Platform and Toolkit to Curate Biodiversity Data (started Aug. 2014)

#### **Additional Material**

*...not part of the presentation... (for Q&A as needed)* 

# **Date Validation**

- Check:
  - Collector's life span
  - .. vs. Date-Collected
- Possible outcomes:
  - Valid
  - Corrected
  - Unable to validate
    - Internal inconsistency
       Contradicting dates
    - External inconsistency
      - Lack of date data



# ... Logic Behind Each Step (cont'd)

#### Scientific Name Validation

- Customer-dependent:
  - Collection Managers:
    - Nomenclature
  - Researchers:
    - Taxonomy (current names)
- Several Remote services
  - IPNI, GNI, ...



• .... <your logic here> ...

#### **Example Output ...**

00	0										ASUHIC.xls												
PH	E (2) E	******	· · · · · · · · · · · ·	125% +	0																		Q+ Search in Sheet
	Home L	avout Tables Charts	SmartArt Formulas D	ata Review																			A 0
	Lda	Ford	Aligna	ndert.	Number									Format									Carla Thomas
1	- 18 Fill -	Arial = 10 =	A- A- = = abc	- Wrap Text - G	eneral	1-1 ( <b>11</b> )	Normal	Bad	ord	Neutral 1	Colouistion	NOVEL 1	Explanatory	lineut		Note	Output	Warning Text	Heading 1	Heading 2	Head	ing 3	West They Inthe Lands Hart
	Ser and	190 F . 11 (	CALL INCOME IN LAND		NOT NOT SHE	50 Conditional	and the second second	Title	arel	Mar Arrents	Mar Arrange 2	Change and	Mill Chromet	Mar Annant	Mill Account	and Account	ANN Access	AND Arrowship	ANN Account	-	ALC: ALC:	Accesst	
Paste	Clear ·			Merge *	* 70 7 300	formatting	Heading 4	Thue 1	otal	20% - Accent1	20% - Accentz D	Uni - Accenta	20% - Accente	20% + Accents	20% - Accente	-40% > Accent1	40% + ACCOULT	40% - Accents	- HUTS - ACCOINCE	AUX + ACCE	15. 40%	ACCENTS	Insert Delete Format Themes All*
	\$1	🗧 💿 💿 🥂 fx Basis of Re	cord																				
-	A	В	C	D	E	G	н	1	1	K	1.1.1			4	Q	R		т	U	V	W	x	Y
1 5	ocs.07.05	Record Id	Determiner Des R. Davis, Peter M. Jump	Georeference Source	s Geodetic Datum	Catalog Number	Family	Recorded By	State Provi	nce Start Day	y Scientific Name		Scientific Name	Authorship	Country	Modified 2013-00-20-1	2:17 Decim	nal Latitude	Decimal Longitur	ie Month	Day Ye	ar Loca	lity 5 of Dottai
3 1	965-07-06	SCAN occurrence 797363	H.A. Sculen	SCAN	WGS84	ASUHIC0020716	Crabronidae	J.H. Davidson	Arizona	187	Eucerceris canalicul	lata	(Say, 1823)	Nor A.	USA	2013-04-12 15:3	2:16 31.93	41	-109.117	7	\$ 19	65 2 mi.	NE of Portal
15	993-11-11	SCAN occurrence.9734	Robert A. Johnson	Google Earth	WGS84	ASUHIC0001668	Formicidae	Robert A. Johnson	Arizona	315	Pogonomyrmex rug	osus	Emery, 1895		USA	2012-09-26 02:5	4:34 33.93	8048	112.973563	91	11 19	93 14 m	. W of Wickenburg
2.16	965-04-10	SCAN occurrence 4360215	Mont A. Cazier	Google Earth	WGS84	ASUHIC0052442	Scarabaeidae	J.H. Davidson	Arizona	100	Gymnopyge hopliae	dormis	Linell, 1895		USA	2014-01-15 11:3	8:50 33.61	5386	5114.217126	4	10 19	66 4 mi.	S of Quartzsite
7	968.05.18	SCAN occurrence 13224	EH Biodoe	SCAN	WGS84	ASUHIC0001316	Geometridae	Roman S. Wielous	Anzona	139	Speranza trilipeana	-	(Grossbork 191	0)	USA	2014-01-29 15:1	8.53 33.96	2139	5111 84833	5	13 19	68 Sove	o Socioos
8 1	965-08-27	SCAN.occurrence.927849	W.B. Warner	SCAN	WGS84	ASUHIC0032951	Scarabaeidae	J.H. Davidson	Arizona	239	Euphoria verticalis		Hom, 1880		USA	2013-07-04 14:3	6.58 31.93	41	-109.117	8	27 19	65 2 mi.	NE of Portal
9 1	965-07-12	SCAN.occurrence.491825	Don R. Davis, Peter M. Jump	SCAN	WGS84	ASUHIC0016413	Tineidae	J.H. Davidson	Arizona	193	Acrolophus daviseli	UB .	Beutenmüller, 18	387	USA	2013-09-20 20:1	3.17 31.89	9097	-109.14083	7	12 19	65 1 mi.	S of Portal
10 1	965-07-23	SCAN.occurrence.4358761	Frank F. Hasbrouck	SCAN	WGS84	ASUHIC0050166	Megalopygidae	J.H. Davidson	Arizona	204	Megalopyge bisses	8	Smith & Abbol. 1	797	USA	2014-01-13 11:4	3.38 31.89	9097	-109.14083	7	23 19	65 1 mi.	S of Portal
	964-08-18	SCAN.occurrence.691216	Don R. Davis, Peter M. Jump	SCAN	WGS84	ASUHIC0017964	Tineidae	Jean H. Puckle	Arizona	231	Acrolophus filicornu	3	(Walsingham, 10	(87)	USA	2014-01-19 00:1	5.55 31.914	4255	-109.130297	8	18 19	64 Porta 71 0.5 m	NE of Costin Dock page Dolla Rispon Jaka
11. 19	990-07-13	SCAN occurrence 4351594	ASUHIC	SCAN	WGS84	ASUHIC0049796	Scarabaeidae	Kim Wismann	Arizona	194	isonychus arizonen	cie	Howden 1959		USA	2013-12-28 14:2	5:24 31.38	3342	110.246642	7	13 19	90 Ash 0	Canvon, Huachuca Mountains, 1 mi, W of Highway I
10 19	965-08-03	SCAN occurrence 798832	Don R. Davis, Peter M. Jump	SCAN	WGS84	ASUHIC0021769	Tineidae	J.H. Davidson	Arizona	215	Acrolophus variabili	5	Walsingham, 18	87	USA	2013-09-20 20:1	3:17 31.896	9097	\$109.14083	8	3 19	65 1 mi.	S of Portal
15 1	965-07-13	SCAN.occurrence.6952196	Mont A. Cazier	SCAN	WGS84	ASUHIC0056189	Megalopygidae	J.H. Davidson	Arizona	794	Norape tenera		(Druce, 1897)		USA	2014-02-19 12:0	8:20 31.89	9097	109.14083	7	13 19	65 1 mi.	S of Portal
16 1	924-05-00	SCAN.occurrence.10826	Roman S. Wielgus	GeoLocate	WGS84	ASUHIC0003678	Nymphalidae	E.V. Walter	Arizona		Vanessa cardui	_	(Linnaeus, 1758	1	USA	2013-04-29 19:2	927 33,42	5914	-111.940005	3	19	24 Temp	0
18	994-08-09	SCAN occurrence 11331	ASUMIC	Scronie Earth	WGS84	ASUHIC0037241 ASUHIC0004437	Formicidae	JD Parket	Arizona	5170	Sphaeropenaima un	P	(Blake, 16/9) Mayr 1866		LISA	2013-09-18 10:1	6-34 31.91	4200	109.130297	8	50 50	64 Porta 64 0.26	mi N of Pineon Sorings
19	955-09-27	SCAN.occurrence.432065	ASUHIC	SCAN	WGS84	ASUHIC0013932	Blattellidae	Blomguist	Arizona	270	Blattella germanica		(Linnaeus, 1767	)	USA	2012-10-17 13:2	2.16 33.42	5914	-111.940005	9	27 19	55 Temp	e
20 1	964-08-06	SCAN.occurrence.1021847	W. Ferguson	SCAN	WGS84	ASUHIC0036969	Mutilidae	Jean H. Puckle	Arizona	219	Odontophotopsis er	ebus	(Melander, 1903	3	USA	2013-07-23 07:1	0.53 31.91	4255	-109.130297	8	8 19	64 Ports	1
21 1	964-07-18	SCAN.occurrence.797646	H.A. Scullen	SCAN	WG584	ASUHIC0020935	Crabronidae	Jean H. Puckle	Arizona	200	Eucerceris tricolor		(Cockerell, 1897	3	USA	2013-02-01 17:2	7:39 31.934	41	-109.117	7	18 19	64 2 mi.	NE of Portal
44	966-06-28	SCAN occurrence 805318	R.M. Bohart	SCAN	WGS84	ASUHIC0026158	Vespidae	J.H. Davidson	Arizona	179	Stenodynerus apaci	he	(Bohart, 1949)	1. A.	USA	2013-04-29 14:3	3:07 35.774	664	-110.13611	5	28 19	66 Jedd	to Trading Post
	966-07-21	SCAN occurrence 362399	R M. Bohart	SCAN	WGS84	ASUHIC0010546	Crabronidae	J.M. Davidson	Arizona	202	Civpeadon taunulus		(Cockerell 1895		LISA	2013-01-02 18:4	4.48 32.22	6356	109,795043	7	21 19	66 3 mi	S of Wilcox
25 1	971-06-27	SCAN.occurrence.11989	Roman S. Wielgus	SCAN	WGS84	ASUHIC0004018	Nymphalidae	Irwin Leeuw	Illinois	178	Speyeria aphrodite		Fabricius, 1787		USA	2014-01-30 13:5	3:38 41.66	722	67.83028	8	27 19	71 Palor	Park
26 1	968-06-28	SCAN.occurrence.797370	G.R. Ferguson	SCAN	WGS84	ASUHIC0020723	Crabronidae	N. Leppla	Baja Califor	nia 180	Eucerceris ferrugino	680	Scullen, 1939		Mexico	2013-12-05 19:5	1.51 32.50	1801	117.072789	6	28 19	68 San/	Ingel
27	959-08-01	SCAN.occurrence.924788	R.M. Bohart	SCAN	WGS84	ASUHIC0032064	Vespidae	R.S. Beal	Colorado	213	Vespula arenaria		(Fabricius, 1775	2.6	USA	2013-06-14 18:4	3:24 39.52	111	105.30472	8	7 79	59 Conit	er
28	964-06-10	SCAN occurrence 1030759	W. Perguson Easth E. Mashersuck	SCAN	WGS84	ASUHIC003/116	Mutilidae	Jean H. Puckle	Anzona	223	Acanthopholopsis b	equaertii	Schuster, 1958	707	USA	2013-09-18 00:2	3.30 31.91	4255	-109.130297 5100.14083	7	10 19	64 Porta	8 of Bostol
10	965-08-06	SCAN occurrence, 1103658	ASUHIC	SCAN	WGS84	ASUHIC0046908	Scarabaeidae	J.H. Davidson	Arizona	218	Xvlorvctes iamaicen	a nais	(Drury, 1773)		USA	2013-11-07 23:4	7:54 31.89	9097	-109.14083	8	8 99	65 1 ml.	S of Portal
31 1	994-09-09	SCAN.occurrence.12494	ASUHIC	SCAN	WGS84	ASUHIC0005223	Formicidae	S. Roberts	Arizona	252	Camponotus		Mayr, 1861		USA	2012-09-26 02:5	4:34 34.05	3539	-109.576804	9	3 19	94 Shee	p Springs
32 11	964-07-14	SCAN.occurrence.359986	R.M. Bohart	SCAN	WGS84	ASUHIC0009360	Crabronidae	Jean H. Puckle	Arizona	196	Philanthus gibbosus		(Fabricius, 1775	)	USA	2013-01-02 18:4	0:44 31,93	61	-109.117	7	74 79	64 2 mi.	NE of Portal
- 1 T	994-11-05	SCAN.occurrence.11284	ASUHIC	GeoLocate	WGS84	ASUHIC0004398	Formicidae	R. Barnes	Arizona	309	Pheidole	ALC: NO	Westwood, 1839		USA	2012-09-26 02:5	4:34 33.70	139	-111.3425	11	5 19	94 Lone	Pine Saddle
15 8	964-06-16	SCAN occurrence 1054482	W.B. Warner W. Ferguson	SCAN	WGS84	ASUHIC0042645	Mutilidae	Joan H. Puckle	Anzona	768	Accorbitionals direct	JOI SHO	(Easey, 1915) (Ease, 1899)		LISA	2013-09-06 20.4	7:08 31.91	4255	5109 130297	-	36 39	64 Porta	avi u Ponal
36 1	968-06-10	SCAN.occurrence.5697734	A.R. Hardy	SCAN	WGS84	ASUHIC0054443	Scarabaeidae	N. Leppla	Sonora	162	Diplotaxis corbula		Vaurie, 1960		Mexico	2014-02-03 14:2	3:15 32.35	7234	114.361937	8	10 19	68 26 m	E of San Luis
37 1	993-02-09	SCAN.occurrence.10568	R.R. Snelling, Robert A. Johns-	Label	WGS84	ASUHIC0002911	Formicidae.	Robert A. Johnson	Baja Califor	nia <sup>7</sup> 40	Solenopsis xyloni		McCook, 1879		Mexico	2014-02-24 18:2	0.01 30.02	5	115.533333	2	9 19	93 8 mi.	E of El Rosario Bridge
38 1	966-04-10	SCAN.occurrence.4360210	Mont A. Cazier	Google Earth	WGS84	ASUHIC0052437	Scarabaeidae	J.H. Davidson	Arizona	100	Gymnopyge hopliae	formis	Linell, 1895		USA	2014-01-15 11:3	8:50 33.61	5386	-114.217128	3	10 19	66 4 mi.	S of Quartzsite
39	965-07-10	SCAN accurrence 490977	Don R. Davis, Peter M. Jump	SCAN	WGS84	ASUHIC0013620	Tineidae	J.H. Davidson	Arizona	7191	Acrolophus daviselli Orianhur, elunatio	us	Beutenmüller, 10	\$87.	USA	2013-09-20 20:1	3:17 31.89	9097	-109.14083	7	10 19	65 1 mi. 75 Dentr	S of Portal
41 1	975-06-20	SCAN occurrence 809894	R. Leuschner	SCAN	WG584	ASUHIC0025658	Tortricidae	Roman S. Wielgus	Arizona	171	Bactra verutana chr	VIAR	Heinrich, 1926		USA	2013-03-17 01:4	5.07 33.44	833	-112.07333	8	20 19	75 Phoe	olx.
42 1	978-07-08	SCAN.occurrence.798971	Don R. Davis, Peter M. Jump	SCAN	WGS84	ASUHIC0021844	Tineidae	Roman S. Wielgus	Anzona	189	Acrolophus variabili	\$	Walsingham, 18	87	USA	2013-03-20 02:2	3.09 31.47	5908	-110.269907	7	B 19	78 5131	S Bannock Street, Sierra Vista, Huachuca Mountai
43 1	966-05-07	SCAN.occurrence.802843	R.M. Bohart	SCAN	WGS84	ASUHIC0024827	Crabronidae	Stanley A. Gorodens	iki Arizona	127	Oxybelus pitanta		Pate, 1943		USA	2013-03-02 11:0	4:34 33.85	2572	-113.904232	5	7 19	66 5 mi.	SE of Bouse
	972-07-29	SCAN.occurrence.6947616	W.B. Warner	SCAN	WGS84	ASUHIC0054878	Scarabaeidae	Martin A. Kolner	Arizona	211	Diplotaxis dentella		Fail, 1909		USA	2014-02-12 17:4	6:36 32.33	749	-110.691998	7	29 19	72 Molir	o Basin picnic area, Santa Catalina Mountains
46	996-00-13	SCAN perumanea (362239	WB Wheney	SCAN	WG584	ASUHIC0045735	Scarabaaidaa	J.M. Davidson	Anzona	164	Ophryastes Chosupaethur Swit	nannin .	Germar, 1829		LISA	2013-12-05 13:4	2:15 32.71	4498	111 209113	2	3 19	66 2 mi.	E of Moenkopi E of Tacea
47 1	964-08-06	SCAN.occurrence.1030470	J.R. Zimmerman	SCAN	WGS84	ASUHIC0037071	Mutilidae	Jean H. Puckle	Anzona	219	Odontophotopsis on	ata	(Melander, 1903		USA	2013-09-17 23:4	2.02 31.91	4255	109,130297	8	8 19	84 Porta	
48 20	000-06-11	SCAN.occurrence.922922	Nico M. Franz	SCAN	WGS84	ASUHIC0024111	Curculionidae	J. Ascher	New York	163	Mononychus		Schüppel in Ger	mar. 1824	USA	2013-06-13 16:2	6.51 42.43	2126	-76.484024	8	31 20	00 Six N	lle Creek, Ithaca
49 1	965-06-29	SCAN.occurrence.1027978	W.B. Warner	SCAN	WGS84	ASUHIC0036702	Scarabaeidae	J.H. Davidson	Arizona	180	Cyclocephala pasad	fenae	(Casey, 1915)		USA	2013-09-20 20:1	3:17 31.89	9097	-109.14083	8	29 19	65 1 mi.	S of Portal
50 1	994-11-26	SCAN.occurrence.492349	Nico M. Franz	Google Earth	WGS84	ASUHIC0015255	Curculonidae	R.M. Baranowski	Arizona	330	Artipus		Sahlberg, 1823		Turks and Calo	2012-12-31 16:2	8:51 21.95	1819	-71.95915	91	26 19	94 North	Calcos Island, Horse Stable Road
52	964-07-22	SCAN occurrence 797875	Doo R. Davis, Peter M. Jumo	SCAN	WG\$84	ASUHIC0020121	Tineidae	Jean H. Puckle	Arizona	204	Acrolophus marring	aster macroo	(Walshoham 14	87)	USA	2013-02-22 12-3	6:36 31.91	4255	-109 130297	7	22 40	64 Ports	(
51 1	915-07-07	SCAN.occurrence.821598	Mesa Experiment Station		WGS84	ASUHIC0026827	Tortricidae	H.B. Scammell	New Jersey	188	Ancylis comptana	and manual	Frölich, 1828		USA	2013-04-04 00:1	9:39 40.07	8516	74.653147	7	7 19	15 Geor	getown
54 1	965-07-04	SCAN occurrence 797935	Don R. Davis, Peter M. Jump	SCAN	WGS84	ASUHIC0020181	Tineidae	Jean H. Puckle	Arizona	785	Acrolophus pervipal	pus	Hasbrouck, 1964	1 C	USA	2013-09-20 20:1	3:17 31.895	9097	-109.14083	7	A 19	65 1 mi.	S of Portal
55 19	978-05-00	SCAN.occurrence.431031	Charles W. O'Brien	Google Earth	WGS84	ASUHIC0012014	Curculionidae	B.E. King	Colorado		Cimbocera conspen	50	Fall, H.C., 1907		USA	2012-10-10 19:4	5.33 40.09	483	108.827059	3	19	76 Rio E	lanco Oil Shale Project, South Slope
57	964-08-14	SCAN occurrence 1001585	C.E. Mickel Dop R. Davis, Balar M. Luna	SCAN	WGS84	ASUHIC0032345	Mutilidae	Jean H. Puckle	Arizona	227	Crypholes elevatus		Blake, 1886	187	USA	2013-07-28 17:2	31.91	4255	-109.130297	8 7	14 19	64 Porta	S of Postal
58 1	966-04-02	SCAN occurrence 482992	Charles A Triplehom	SCAN	W0584	ASI (HICO014249	Tenebrionidae	Stanley C. Williams	Arizona	\$22	Centricolera marica	ta .	LeConte 1851		LISA	2012-11-18 17-3	4-27 113 35	6437	5112 074373	14	7 99	68 South	Mountain Regional Park
1	10	Analysis Results Scientifich	NameValidator Details 🦼 DateValidat	or Details 2 GeoRefValid	tor Details / +															_		_	
Summer of the local division of the local di	Normal Vi	www.Ready			Sum=0																		

#### ... close up ...

-	S1	: 🛛 🔿	( fx Basis of Record						-			
	J	K	L			м	Q	R	Т		U	V
1	State Province	Start Day	Scientific Name	Scientif	fic Nam	e Authorship	Country	Modified	Decimal Lat	titude	Decimal Longitude	Mon
2	Arizona	186	Acrolophus davisellus	Beutenr	müller, 1	1887	USA	2013-09-20 20:13:17	31.899097		-109.14083	7
3	Arizona	187	Eucerceris canaliculata	(Say, 18	823)		USA	2013-04-12 15:32:16	31.9341		-109.117	7
4	Arizona	315	Pogonomyrmex rugosus	Emery,	1895		USA	2012-09-26 02:54:34	33.938048		-112.973563	11
5	Arizona	100	Gymnopyge hopliaeformis	Linell, 1	895		USA	2014-01-15 11:38:50	33.615386		-114.217126	4
6	Arizona	226	Euptoieta claudia	Cramer,	, 1775		USA	2013-03-08 15:50:25	34.052139		-109.729612	8
7	Arizona	139	Speranza trilinearia	(Grossb	beck, 19	910)	USA	2014-01-29 15:18:53	33.96222		-111.84833	5
8	Arizona	239	Euphoria verticalis	Horn, 18	880		USA	2013-07-04 14:36:58	31.9341		-109.117	8
9	Arizona	193	Acrolophus davisellus	Beutenr	M ntific Name Authorship enmüller, 1887 , 1823) ry, 1895 I, 1895 I, 1895 her, 1775 ssbeck, 1910) , 1880 enmüller, 1887 h & Abbot, 1797 singham, 1887) nte, 1856 den tingl Z Ce, ScientificNameValida taeL CORRECT cORRECT CURATED CURATED CURATED UNABLE_CURATE CURATED UNABLE_DETERMINE CORRECT CURATED UNABLE_DETERMINE CORRECT UNABLE_DETERMINE CORRECT UNABLE_DETERMINE CORRECT UNABLE_DETERMINE CORRECT UNABLE_DETERMINE CORRECT UNABLE_DETERMINE CORRECT UNABLE_DETERMINE CORRECT UNABLE_DETERMINE CORRECT UNABLE_DETERMINE CORRECT UNABLE_DETERMINE CORRECT UNABLE_DETERMINE CORRECT UNABLE_DETERMINE CORRECT UNABLE_DETERMINE CORRECT UNABLE_DETERMINE CORRECT UNABLE_DETERMINE CORRECT COR		USA	2013-09-20 20:13:17	31.899097		-109.14083	7
10	Arizona	204	Megalopyge bissesa	Smith &	Abbot,	1797	USA	2014-01-13 11:43:38	31.899097		-109.14083	7
11	Arizona	231	Acrolophus filicornus	(Walsing	gham, 1	1887)	USA	2014-01-19 00:15:55	31.914255		-109.130297	8
12	Arizona	235	Oxygrylius ruginasus	Leconte	9, 1856			2013_00_20 08-44-36	21 288/61		111 005525	<b>-</b> 2
13	Arizona	194	Isonychus arizonensis	Howden	n							
14	Arizona	215	Acrolophus variabilis	Walsing			Z	AA			AB	
15	Arizona	194	Norape tenera	(Druce,	Sci	entificNameValio	lator	DateValidator	G	GeoRefValidator		
16	Arizona		Vanessa cardui	(Linnae)	L CO	RRECT		CORRECT	C	ORREO	СТ	
17	Arizona	219	Sphaeropthalma uro	(Blake,	1 CO	RRECT		CORRECT	C	ORREO	СТ	
18	Arizona	170	Leptothorax	Mayr, 18	8 CO	RRECT		CORRECT	C	ORREO	СТ	
10	• •	<b>B</b> 70	Prin a	<i>a</i> .	CU	RATED		CORRECT	C	ORREO	СТ	
					CU	RATED		CORRECT	U	JNABLE	_DETERMINE_VALIDIT	Γ <b>Y</b>
					UN	ABLE_CURATE		CORRECT	C	ORREO	СТ	
					CO	RRECT		CORRECT	C	ORREO	СТ	
					CO	RRECT		CORRECT	C	ORRE	СТ	
					CU	RATED		CORRECT	C	ORRE	СТ	
					UN	ABLE_DETERMI	NE_VALIDITY	CORRECT	C	ORREG	CT	
					CU	RATED		CORRECT	U	INABLE	CURATE	
					2 CO	RRECT		CORRECT	U	JNABLE		
					CU	RATED		CORRECT	0	ORREG		
					UN	ABLE_DETERMI	NE_VALIDITY	CORRECT		ORREG		
									INE_VALIDITY C	CORRECT		
						ABLE_DETERMI	NE_VALIDITY	CORRECT				
					- 00	RRECT		CORRECT				
								CORRECT			יי	
					CU		NE_VALIDITT	CORRECT			ו כ די	
								CORRECT			יידי	
					00	RRECT		CORRECT		ORRE	от Т	
					- 00	RRECT		CORRECT		ORRE	ст Т	
					CU	RATED		CORRECT	0	ORRE	CT	
					00	RRECT		CORRECT	ŭ	INABLE	CURATE	
					UN	ABLE DETERMI	NE VALIDITY	CORRECT	C	ORRE	CT	
					UN	ABLE DETERMI	NE VALIDITY	CORRECT	C C	ORRE	CT	
					CU	RATED		CORRECT	C	ORREO	CT	
					CO	RRECT		CORRECT	Č	ORREO	СТ	
						DDEOT		0000507				

FP2K: Towards Curation Workflows

# Kurator/P & YW: the road ahead ...

#### • YesWorkflow:

- ... finishing support for retrospective provenance without using a runtime provenance recorder!
- Key insight: scientists already leave provenance "bread crumbs" behind! (it's not an accident!)
- Exploit that via annotations: URI-templates
- Kurator[/P]:
  - How far can we go towards ASAP via YW?

### YesWorkflow.org



# **YW-RECON:** Prospective & Retrospective Provenance ... (almost) for free!



 YW annotations in the script (R, Python, Matlab) are used to recreate the workflow view from the script ...





FP2K: Towards Curation Workflows

### Summary: Data Curation with Scientific Workflow Systems

#### **Scientific Workflows**

- [+] Automation
- [+] Scalability
- [+] Abstraction
- [+] Provenance
- •
- [+/0] Easy to use
  - [0] learning a new paradigm
- [-] Teaching resources: learning a new language!
- [-] Special expertise needed for deep changes e.g. new Java actors, shims, ...

#### Kurator/P: Scripts + YesWorkflow ++

Scripts: [+] Automation, [0] Scalability, [-] Abstraction, [0/-] Provenance Now: Scripts + YesWorkflow Annotations

#### • [+] Abstraction

- explain your methods to mere mortals
- => encourage (re-)use
- [+] Provenance:
  - YesWorkflow (prospective and retrospective provenance)
- [+] Language independent (R, Matlab, Python, ...)
- [+] Empower tool makers (script programmers): give them ...
  - ... some immediate benefits (workflow views, retrospective provenance)
  - ... some long term benefits: think about your methods differently
  - => dataflow programming => [+] Scalability