

### The Case for Adaptive Hierarchical Metadata

Steve Worth

#### Henrique Nogueira

Director, Global Innovation Operations, EMC

SW Engineer Eldorado Research Institute

#### Arthur Guerra

SW Engineer Eldorado Research Institute

EMC<sup>2</sup>

iRODS User Group Meeting 2016 6/9/2016





### Find out what your DNA says about you and your family.

- Learn what percent of your DNA is from populations around the world
- Contact your DNA relatives across continents or across the street
- Build your family tree and enhance your experience with relatives

order now



Copyright 2016 EMC Corporation. All rights reserved.

# 5% = 71M = 20 ExB = \$2.8 Billion/Yr

**EMC**<sup>2</sup>

Copyright 2016 EMC Corporation. All rights reserved



Copyright 2016 EMC Corporation. All rights reserved.



Copyright 2016 EMC Corporation. All rights reserved.

```
usxxpaeglem1:ino paegle$ ls -l
total 0
-rw-r--r-- 1 paegle CORP\Domain Users 0 Sep 23 12:29 my.bam
-rw-r--r-- 1 paegle CORP\Domain Users 0 Sep 23 12:29 my.log
usxxpaeglem1:ino paegle$ []
```



### **Coping Behaviors**



#### Advanced Metadata is Critical

- Metadata can be stored/transported more cheaply
  - An object with 1000 tags, 4096 character limit, 3 tag AVU model requires < 12MB for the tag set.
- Maximize storage assets find what's valuable, no matter where it's located
  - Eliminate data that is not valuable (so called "dark data")
- Automate movement and processing of data
- Securely share data with collaborators
  - Easier to separate PHI from non-PHI data
  - Can be used to ensure data is authentic and unaltered
- Metadata supports indexing / fast search methods
  - Ideal for non-SQL techniques

### Working With iRODS...

- User commands (subset)
  - iinit
  - iput
  - iget
  - imkdir
  - ichmod
  - icp
  - irm
  - ils
  - ipwd
  - icd
  - irepl
  - iexit
  - ipasswd
  - ichksum
  - imv
  - iphymv
  - ireg
  - irmtrash
  - irsync

Copyright 2016 EMC Corporation. All rights reserved.

- ladmin commands (subset)
  - It
  - Ir
  - Is
  - Iz
  - Ig
  - mkuser
  - moduser
  - aua
  - rua
  - rpp
  - rmuser
  - mkdir
  - rmdir
  - mkresc
  - modresc
  - modrescdatapaths
  - rmresc
  - mkzone
  - modzone



#### MetaLnx – iRODS Simplified

Organized

sworth@swdebian64-1: ~ File Ecit View Search Terminal Help (5).jpg sha2:ympAnVjSGckaHMpJhnGxAGJ92WxhvC82oaf1inInxKU= demoResc generic /var/lib/irods/iRODS/Vault/home/sworth/lots of files/2012 Spring 0055 (5).jpg sworth 0 demoResc 191005 2015-09-23.01:05 & 2012 Spring 0055 (6).jpc sha2:ympAnVjSGckaHMpJhnGxAGJ92WxhvC82oaf1inInxKU= demoResc generic /var/lib/irods/iRODS/Vault/home/sworth/lots of files/2012 Spring 0055 (6).jpg 0 demoResc 191005 2015-09-23.01:05 & 2012 Spring 0055 sworth (7).jpg sha2:ympAnVjSGckaHMpJhnGxAGJ92WxhvC82oaflinInxKU= demoResc generic /var/lib/irods/iRODS/Vault/home/sworth/lots of files/2012 Spring 0055 (7).jpg 191005 2015-09-23.01:05 & 2012\_Spring\_0055 sworth 0 demoResc (8).jpg sha2:ympAnVjSGckaHMpJhnGxAGJ92WxhvC82oaflinInxKU= demoResc generic /var/lib/irods/iRODS/Vault/home/sworth/lots\_of\_files/2012\_Spring\_0055 (8).jpg 0 demoResc 191005 2015-09-23.01:05 & 2012 Spring 0055 sworth (9).jpg sha2:ympAnVjSGckaHMpJhnGxAGJ92WxhvC82oaflinInxKU= demoResc aeneric /var/lib/irods/iRODS/Vault/home/sworth/lots of files/2012 Spring 0055 (9).jpg sworth 0 demoResc 191005 2015-09-23.01:05 & 2012 Spring 0055. jpg sha2:ympAnVjSGckaHMpJhnGxAGJ92WxhvC82oaflinInxKU= demoResc aeneric /var/lib/irods/iRODS/Vault/home/sworth/lots of files/2012 Spring 0055.pg sworth@swdebian64-1:~\$ 📕 Complex

Collections			Search Q			
🗅 iRods / t	empZone / home / s	worth / lots_of_fi	les			
Browse In	fo Metadata Per	missions	ᆂ Upl	load 🖿 Add Collection		
10 items per pag	ge 👻 Previous	5 1 2 3	1503 Next	Page# Go		
	Displ	aying items <b>1 - 10</b> of	15021 total			
Select Action 👻						
Name	Owner	Kind	Modified	Size		
2000	Spring sworth	file	Sep 22 2015, 19:			
□ 🗳 2000	Spring sworth	file	Sep 22 2015, 19:			
2000	Spring sworth	file	Sep 22 2015, 19:			
2000	Spring sworth	file	Sep 22 2015, 19:			
🗆 🔛 2000	Spring sworth	file	Sep 22 2015, 19:			
2000	Spring sworth	file	Sep 22 2015, 19:			
2000	Spring sworth	file	Sep 22 2015, 19:			
2000	Spring sworth	file	Sep 22 2015, 19:			
🗆 📓 2000	Spring sworth	file	Sep 22 2015, 19:			

#### Simplified Admin



#### Metadata Automation



- 1. FASTQ\_FILE, path to fastq file on ftp site
- 2. MD5, md5sum of file
- 3. RUN\_ID, SRA/ERA run accession
- 4. STUDY\_ID, SRA/ERA study accession 5. STUDY\_NAME, Name of stury
- 6. CENTER\_NAME, Submission centre name
- 7. SUBMISSION\_ID, SRA/ERA submission accession
- SUBMISSION\_DATE, Date sequence submitted, YYYY-MM-DAY
   SAMPLE\_ID, SRA/ERA sample accession
- 10. SAMPLE NAME, Sample name
- 11. POPULATION, Sample population, this is a 3 letter code and it is defined in README.populations

- 12. EXPERIMENT\_ID, Experiment accession 13. INSTRUMENT\_PLATFORM, Type of sequencing machine 14. INSTRUMENT\_MODEL, Model of sequencing machine
- 15. LIBRARY\_NAME, Library name
- 16. RUN\_NAME, Name of machine run
- 17. RUN\_BLOCK\_NAME, Name of machine run sector (This is no longer recorded so this column is entirely null, it was left in so as not to disrupt existing sequence index parsers)
- 18. INSERT\_SIZE, Submitter specifed insert size
- 19. LIBRARY\_LAYOUT, Library layout, this can be either PAIRED or SINGLE
- 20. PAIRED\_FASTQ, Name of mate pair file if exists (Runs with failed mates will have a library layout of PAIRED but no paired fastg file)
- 21. WITHDRAWN, 0/1 to indicate if the file has been withdrawn, only present if a file has been withdrawn
- 22. WITHDRAWN\_DATE This is generally the date the file is generated on
- 23. COMMENT, comment about reason for withdrawal 24. READ\_COUNT, read count for the file
- 25. BASE\_COUNT, basepair count for the file

26. ANALYSIS\_GROUP, the analysis group of the sequence, this reflects sequencing strategy. Currently this includes low coverage, high coverage, exon targetted and exome to reflect the 2 non low coverage pilot sequencing stratergies and the 2 main project sequencing stratergies used by the 1000 genomes project.

### Metadata Templates



#### Metadata Searching



#### MetaLnx Value Proposition: Simplified Metadata Management

- Simple, easy iRODS grid administration
- Simple, intuitive tool for non-IT researchers:
  - Collection management
  - Automated embedded metadata extraction
  - Metadata searching / editing
  - Metadata templates consistent tagging of data
- Scale out design supports PCs and tablets
- Extensible to any data vertical

### EMC MetaLnx Architecture





© Copyright 2016 EMC Corporation. All rights reserved.

### EMC iRODS Efforts

#### Metalnx

- An iRODS Administration and Metadata Interface tool
  - Rev 1.0 available on Github at:
    - <u>https://github.com/sgworth/metalnx-web</u> (Metalnx application)
    - <u>https://github.com/sgworth/metalnx-rmd</u> (Remote Monitor Daemon)
    - <u>https://github.com/sgworth/metalnx-msi</u> (iRODS microservices for Genomes)
    - will be available via www.emccode.com soon
  - Supports iRODS 4.x
    - Git has instructions & tools to build RPM, DEB, and Docker images
  - Please try. If you like join the community!
- Storage Resource Drivers for iRODS
  - Isilon HDFS interface
  - ECS Atmos object interface
    - Will be posted this summer

#### Roadmap

- Metadata Validation
- Turnkey rule set deployment
- Self registration
- Faster, large file transfer support
- Federation/Quota support
- Better data in place linking

#### Data Validation – Two Models

#### "Lazy" Validation

- Validation can be added (or not) for each AVU added on each object.
- Enforced through the UI
- Can be used with Metadata templates
- Type options:
  - String (default)
  - Integer
  - Float
  - List (single selection)
  - List (multiple selection)
  - Boolean
  - Date
  - List of Integers
  - List of Float

- Strong Validation
  - Attribute names and validation rules defined in a system vocabulary
  - Users can also have personal vocabularies
  - System settings, forced via UI for "none, system first, user first, user selection" Also can lock system.
  - In "locked" mode an AVU must come from the vocabulary (per rule) and the value must match the validation rule
  - Options same as "Lazy" validation

### Validation Example

<b>626</b>	Collections			
Dashboard	Concentrations	Add Metadata ×		
Resources	← → ← □ iRods / te	Attribute 31.j	jpg	
Users		Flash_level		
<b>101</b>	Metadata	Type		
Groups	Select Action 👻	Text (default)		
Collections	10 Page # >	Integer Float List (single selection) List (multiple selection) Boolean		
Search	D 🗄 Attribute	Date Range of integers	41	
8	ColorSpace	Range of floats		
Templates	Contrast	Close Save changes		
4	CustomRendered			
Shared Links	DateTime	2009:05:04 23:09:02 -		

### Validation Example

DateTime	2009:05:04 23:09:02	-		Details
DateTimeDigitized	2008:09:04 16:40:01	-		0 Details
DateTimeOriginal	2008:09:04 16:40:01	-		0 Details
ExposureMode	0	-		O Details
ExposureProgram	0	-	Attribute: Flash_level	0 Details
Flash	0	-	Value:	O Details
Flash_level	4	-	Unit:	1 Details
FocalLengthIn35mmFilm	-13824	-	-	O Details
GainControl	0	-	Integer	O Details
ISOSpeedRatings	-29696	-		0 Details
LightSource	0	-		0 Details

#### MetaLnx Live Demo



#### Metadata Hierarchy is Real

```
"Individual": {
                                                                 "phenotypes": [{
                                                                         "ontologySource": "ICD10",
     "id": "538764ae-5bc8-4f2b-b691-d99d6e1f754a",
     "groupIds": ["8c7409e2-d8be-4998-b00d-f7d5bee7c164"],
                                                                         "id": "I12.9",
     "name": "John Smith",
                                                                         "name": "hypertensive renal disease"
                                    Partial example from
     "description": null,
                                                                      }],
                                                                      "stagingSystem": "Dukes C",
     "created": 1416416735.
                                    an old GA4GH
     "updated": 1416416735,
                                                                      "clinicalTreatment": null,
                                    metadata use case
     "species": {
                                                                      "strain": null,
       "ontologySource":
                                                                      "info": { }
"http://purl.obolibrary.org/obo/NCBITaxon_9606",
                                                                   },
        "id": "NCBITaxon: 9606",
                                                                 "Experiment": {
       "name": "Homo sapiens"
                                                                      "id": "a93ba556-1e31-4d8c-bfe8-575a562a961c",
                                                                      "name": "Whole Genome Sequencing of Charge-S ARIC
     },
     "sex": "MALE",
                                                                 sample A09994",
     "developmentalStage": {
                                                                      "description": null,
       "ontologySource": "WHO:ICD",
                                                                      "created": 1416413182,
       "id": "C18.5",
                                                                      "updated": 1416413182,
       "name": "Malignant neoplasm of splenic flexure"
                                                                      "runDate": 1416413182,
                                                                      "molecule": "genomics DNA",
     },
     "dateOfBirth": 1400324354,
                                                                      "strategy": "whole genome sequencing",
     "diseases": [{
                                                                      "selection": "random",
                                                                      "library": "IWG_IND-ARXREB.A09994_1pA",
       "ontologySource": "WHO: ICD-O",
                                                                      "libraryLayout": "Paired",
       "id": "8140/3",
        "name": "adenocarcinoma, NOS"
                                                                      "instrumentModel": "Illumina HiSeg 2000"
     }],
                                                                 }
                                                                                                               EMC<sup>2</sup>
```

#### With complex associations...

Cross reference example from the the Gene Ontology Consortium



#### But changing....

#### • From a later GA4GH Metadata team report:

- No single ontology works for all use cases
- Complex nesting needed e.g. disease state and disease stage relationships
- Multiple ontologies may be associated for aspects of a disease

 Recent work is shifting to a metadata framework which ties to other tools for ontology specifics

#### ON THE ABANDONMENT OF THE DEWEY DECIMAL SYSTEM

2012-11-15 · by Erin Lee Barsan · in Libraries & Librarians, Trends & Technologies.



Gone are the days of the card catalog. Is Dewey next? [Princeton University Archives]

While many public libraries stopped using the Dewey Decimal Classification (DDC) system to group fiction books long ago, the last five years have seen an emerging trend of libraries choosing to go "Dewey Free" for non-fiction books as well. In lieu of the industry standard system, these libraries are utilizing more so-called, "user-friendly" bookstore-style layouts. The 3,000 categories in this bookstore system model, Book Industry Standards and Communications (BISAC), are far less expansive and specific when compared to the 27,000 categories used to classify books with Dewey. The main features of BISAC include displaying books by their covers whenever possible, utilizing signage to spotlight current popular categories and grouping books by subject

#### Dewey Decimal Classification System General Guidelines Below 000 Generalities 200 Religion 010 Bibliography 020 Library & information sciences 110 Metaphysics COMPUTERS 120 Epistemology, causatio 030 General encyclopedic works 040 Unassigned 050 General serials & their indexes 130 Paranormal phenomeni 140 Specific philosophical si 150 Psychology Use "COMPUTERS / Electronic Commerce" for works on the computer skills and technology needed to facilitate electronic commerce and use subjects beginning with "BUSINESS & ECONOMICS / E-Commerce" for works discussing the business 060 General organizations & museology 160 Logic 170 Ethics (moral philosoph 180 Ancient, medieval, Orier 190 Modern Western philosof 190 Modern Western philosof 000 General organizations & moseology 070 News media, journalism, publishing 080 General collections 090 Manuscripts & rare books COM000000 COMPUTERS / General 300 Social sciences COM082000 COMPUTERS / Bioinformatics 410 Linguistics 420 English & Old English 430 Germanic languages Gl 440 Romance languages Fi 450 Italian, Romanian langu 460 Spanish & Portuguese I 470 Italic Languages, Latin 490 Lidicia Languages, Cat 310 General statistics 320 Political science COMPUTERS / Business Software see Enterprise Applications / Business Intelligence Tools 330 Economics COM006000 COMPUTERS / Buver's Guides 0 Law 0 Public administration 0 Social services; associations COM007000 COMPUTERS / CAD-CAM 400 Hellenic languages, Latin 480 Hellenic languages, Cla 490 Other languages 70 Education 380 Commerce, communications, transport 390 Customs, etiquette, folklore COM009000 COMPUTERS / CD-DVD Technology COM055000 COMPUTERS / Certification Guides / General 710 Civic & landscape art 610 Medical sciences and medicine T20 Archited users are T20 Archited users COM055010 COMPUTERS / Certification Guides / A+ T20 praise arts, scuipture T40 Drawing & decorative at COM055020 COMPUTERS / Certification Guides / MCSE 620 Engineering & allied operations 630 Agriculture 640 Home economics & family living 650 Management & auxiliary services 660 Chemical engineering 670 Manufacturing 680 Manufacture for specific uses 690 Buildings 760 Graphic arts, printmakin 770 Photography & photogra COM061000 COMPUTERS / Client-Server Computing 780 Music 790 Recreational & perform/ COM091000 COMPUTERS / Cloud Computing COM010000 COMPUTERS / Compiler COM059000 COMPUTERS / Computer Engineering COM012000 COMPUTERS / Computer Graphics

Technology

DDS which dates from 1876 is slowly being replaced by BISAC (Book Industry Standards And Communication)

COMPUTERS / Computer Industry see BUSINESS & ECONOMICS / Industries / Computers & Information



#### Metadata Browsing



#### Tailored Metadata Browsing



## Questions

Copyright 2016 EMC Corporation. All rights reserved.