# The Emergence of
# Computational Archival Science (CAS)

Richard **MARCIANO**

marciano@umd.edu
**University of Maryland iSchool**

UNIVERSITY OF
MARYLAND

UNIVERSITY OF
MARYLAND

COLLEGE OF
INFORMATION
STUDIES
*iSchool*

dcic digital curation
innovation center

APGridPMA/IGTF Mee... Biomedicine & Life S... Biomedicine & Life S... Closing Keynote & C...

Coffee Break  Coffee Break  Coffee Break  Coffee Break  Coffee Break  Coffee Break & Poste...

Coffee Break & Poste...  Cryo-EM Workshop  Data Management & ...  ECAI Workshop

ECAI Workshop  Earth, Environmental...  Earth, Environmental...  Environmental Comp...

GDB Meeting  Humanities, Arts & S...  Humanities, Arts & S...  ICT-Enhanced Educa...

Infrastructure Clouds...  Infrastructure Clouds...  Keynote Session II  Keynote Session III  Lunch

Massively Distributed...  Netwok, Scurity, Infr...  Network, Security, In...  Network, Security, I...

Network, Security, I...  Opening Ceremony &...  Physics & Engineering I  Physics & Engineerin...

Poster Session  Security Workshop  Supercomputing, Hig...  VRE  Workshop on Linux C...

e-Science Activities i...  e-Science Activities i...  e-Science Actvities i...

less...

## Tue 7/3

| Time | | |
|---|---|---|
| 09:00 | **Opening Remarks** | |
| | Conf. Room 2, BHSS, Academia Sinica | 09:00 - 09:30 |
| | **On-the-fly Capacity Planning in Support of High Throughput Workloads** | Dr. Miron LIVNY |
| 10:00 | Conf. Room 2, BHSS, Academia Sinica | 09:30 - 10:30 |
| | **Coffee Break & Photo-taking** | |
| | BHSS, Academia Sinica | 10:30 - 11:00 |

| | | |
|---|---|---|
| 11:00 | **Image Processing in cryoEM: Open problems and current perspectives** | **e-Science Activities in Japan - Building Academic Inter-Cloud Infrastructure** Dr. Kento AIDA |
| | | **eScience Activities in China** Dr. Gang CHEN |
| | Conf. Room 1, BHSS, Academia Sinica 11:00 - 11:45 | Conf. Room 2, BHSS, Academia Sinica 11:15 - 11:30 |
| | | **GSDC activities for scientific computing** Dr. Sang-Un AHN |
| | **Applications of cryo-electron microscopy to understand complex structures** Dr. Sunny WU | **e-Science Activities in Taiwan** Dr. Eric YEN 11:45 - 12:00 |
| 12:00 | | **e-Sciences Activities in Mongolia** Mr. Batzaya E. 12:00 - 12:15 |
| | Conf. Room 1, BHSS, Academia Sinica 11:45 - 12:30 | **Q&A** Conf. Room 2, BHSS, Academia Sinica 12:15 - 12:30 |

| | | |
|---|---|---|
| | **Lunch** | |
| 13:00 | 4F Recreation Hall, BHSS, Academia Sinica | 12:30 - 13:30 |

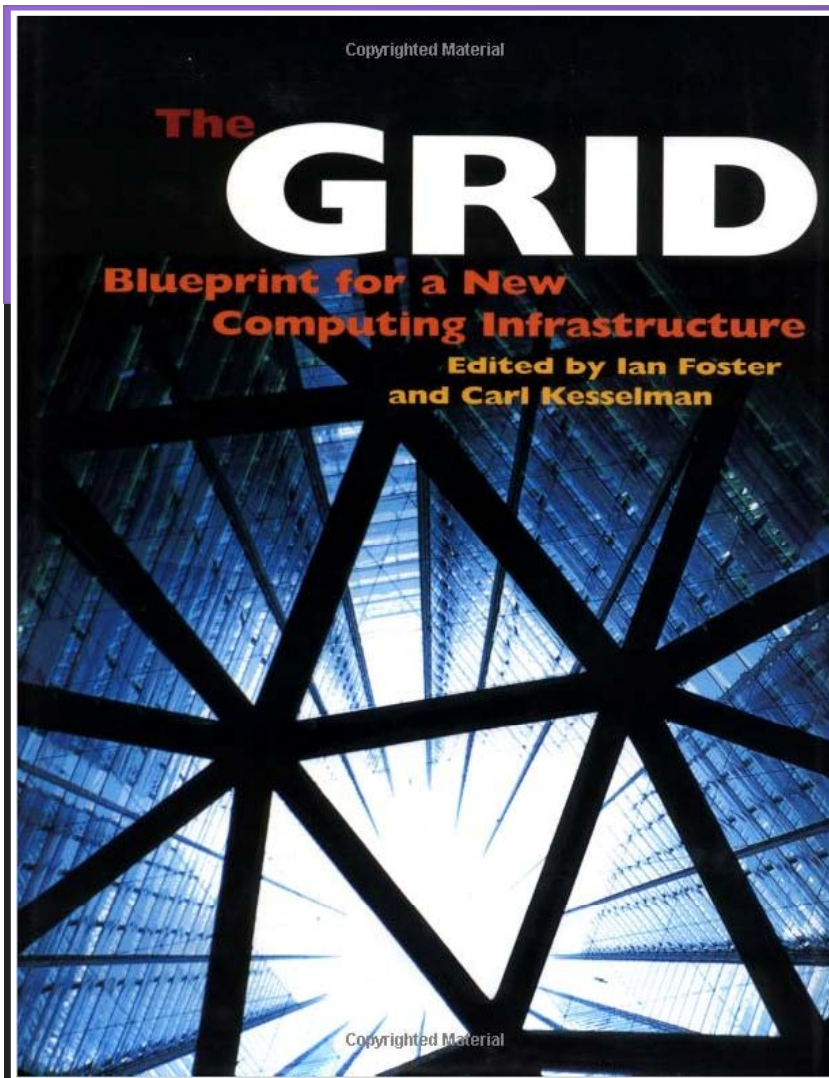| | | | | |
|---|---|---|---|---|
| | **EMAN Dr. Sunny WU 2 (Part 1)** | | | |
| 14:00 | | **Can R&E federations trust Research** | **Towards a cloud-based computing and analysis** | Thomas HAHN |
| | | **WLCG Security Operations Centres Working Group** | **Data storage accounting at RAL** | **Stopping the flow - The Yellow River and China's Grand Geographic Information and Arches** |
| | Conf. Room 1, BHSS, Academia Sinica | **Collaborating for WISEr Information Security** | **dCache, managing Quality of Service in Cloud** | **Volunteered Conf. Room 2, Information and BHSS, Academia Arches** |
| 15:00 | **Machine Learning analysis of CMS data transfers** | **EGI-CSIRT: Coordinating Operational** | **Coffee Break** | **Digital** Janet TAN **Economy and Asian Production Network – A Reality Check for** |
| | **Q&A** | **Q&A** | BHSS, Academia Sinica | |
| | **EMAN Dr. Sunny WU 2 (Part 2)** | **Coffee Break** BHSS, Academia Sinica | **Coffee Break** BHSS, Academia Sinica | **Coffee Break** BHSS, Academia Sinica |
| 16:00 | | **Identifying Suspicious Network Activities in Grid** | **Data Provenance Tracking as the Basis for a** | **Endangered Languages and Flow of Identities:** |
| | | **Modern Monitoring Systems** | **Design and Implementation of Portal System for** | **History and Projections of Tombs Research in** |
| 17:00 | Conf. Room 1, BHSS, Academia Sinica | **Status** Dr. Tian YAN **of Network Security Operations at...** | **Framework for Developing Cloud enabled** | **Earth Deity Mapping and Community** |

## Wed 8/3

| Time | | |
|---|---|---|
| 09:00 | **Caches all the way down: Infrastructure for Data Science** | Prof. David ABRAMSON |
| | Conf. Room 2, BHSS, Academia Sinica | 09:00 - 09:45 |
| 10:00 | **High-resolution Integrative Modelling of Biomolecular Complexes from Fuzzy Data** Dr. Alexandre M.J.J. BONVIN | **Poster Session** Conf. Room 2, BHSS, Academia Sinica |
| | Conf. Room 2, BHSS, Academia Sinica 09:45 - 10:30 | 09:45 - 10:30 |

| | | | |
|---|---|---|---|
| | **POWERFIT and DISVIS (part 1)** | **Coffee Break & Poster Session** BHSS, Academia Sinica | **Coffee Break & Poster Session** BHSS, Academia Sinica |
| 11:00 | Conf. Room 1, BHSS, Academia Sinica | **e-Science Activities in Thailand** Conf. Room 2, BHSS, Academia Sinica 10:30 - 11:00 | **GDB** Mr. Ian COLLIER **Introduction** |
| | **Coffee Break** BHSS, Academia Sinica | **e-Science Activities in Indonesia** Dr. Suhaimi NAPIS | **ASGC Report** Mr. Felix LEE Media Conf. Room, BHSS, Academia Sinica |
| | **POWERFIT and DISVIS (part 2)** | **eScience Activities in Malaysia** Dr. Nam THOAI | **Asian Tier Forum report** Dr. Sang Un AHN Media Conf. Room, BHSS, Academia Sinica |
| 12:00 | | **eScience Activities in Vietnam** Dr. Peter BANZON | **Asian Network Status** Dr. Hsin-Yen CHEN |
| | | **eScience Activities in Philippine** | **Q&A** Media Conf. Room, BHSS, Academia Sinica |
| | | **Q&A** | |
| 13:00 | Conf. Room 1, BHSS, Academia Sinica | **Lunch** | **PC Meeting** |
| | **Lunch** BHSS, Academia Sinica 13:00 - 13:45 | BHSS, Academia Sinica 12:30 - 14:00 | Room 901, BHSS, Academia Sinica |
| 14:00 | **Scipion (part 1)** | **Platform for Humanities Open Data** BHSS, Academia Sinica | **Traceability & Isolation Working Group update** |
| | | **Data Science as a Foundation Toward Open** Dr. Ji-Ping LIN | **IPv6 Rollout** Dr. David KELSEY |
| 15:00 | Conf. Room 1, BHSS, Academia Sinica | **A Preliminary Study on Reconstructing Faded Color by Spectral Estimation Method for Heritage Object** | **Sirtfi - Incident response for for Identity Federation** Media Conf. Room, BHSS, Academia Sinica |
| | **Coffee Break** BHSS, Academia Sinica | | |
| | **Scipion (part 2)** | **Coffee Break & Poster Session** BHSS, Academia Sinica | **Coffee Break & Poster Session** BHSS, Academia Sinica |
| 16:00 | | **Using Advanced e-Systems for Community-Engaged Research** | **Workload Management Trends in WLCG** Maarten LITMAATH |
| | | **A Proposal: ePortfolio for enhancing active learning for the future generation** | **Tier 1 Configuration Evolution & Options** Dr. Josep FLIX Conf. Room, BHSS, Academia Sinica |
| 17:00 | Conf. Room 1, BHSS, Academia Sinica | **Occupation recommendation with major programs for adolescents** | **Wrap Up & WLCG Workshop update** Mr. Ian COLLIER |
| | **Closing remark/round table discussion** | | |

## Thu 9/3

| Time | | |
|---|---|---|
| 09:00 | **Big Data-Driven Drug Discovery** | Prof. Jung-Hsin LIN |
| | Conf. Room 2, BHSS, Academia Sinica | 09:00 - 09:45 |
| 10:00 | **High Performance Computing Environment and Applications in CAS** | Dr. Xuebin CHI |
| | Conf. Room 2, BHSS, Academia Sinica | 09:45 - 10:30 |
| | **Coffee Break & Poster Session** | |
| | BHSS, Academia Sinica | 10:30 - 11:00 |
| 11:00 | **eScience Activities in Singapore** | Dr. John KAN |

## Fri 10/3

| Time | | | |
|---|---|---|---|
| 09:00 | **The 'Cloud Area Padovana': lessons learned after two years of a** | **Examination of dynamic partitioning for multi-core jobs in the Tokyo Tier-2** | **Listening to the ecosystem: the integration of machine** |
| | **Synergy, a new approach for optimizing the resource usage in OpenStack** | **The High Throughput Strategy of IHEP** Dr. Jiaheng ZOU | **Collaboration on monitoring Asian soundscape and the** |
| 10:00 | **Efficiency Improvement on Distributed Cloud System** Mr. Felix LEE et al. | **The Billing System of IHEP Data Center** | **Revealing Philippine Climate Type Using Remotely-sensed Rainfall** |
| | **Coffee Break** | | |
| | BHSS, Academia Sinica | | 10:30 - 10:50 |
| 11:00 | **VCondor - an implementation of dynamic virtual computing cluster** | **Framework for distributing Radio Astronomy processing across Clusters and Clouds** | **NeIC** Dr. John WHITE **EISCAT_3D support project: Nordic computing challenge** |
| | **GUOCCI – The Entryway to Federated CloudIPStore: A Cloud SaaS Repository for your Intellectual Properties** Mr. Radim JANČA | **Exploiting clouds for smart cities applications - The Coelho2020 project** Academia Sinica | **Towards Environmental Computing Compendium** Media Conf. Room, BHSS, Academia Sinica |
| 12:00 | **Supporting Open Science with the EGI Federated Cloud - Experiences, success stories, lessons** Yin CHEN | **Future warming scenario and impacts study over Taiwan: Results from ECHAM5/MPIOM-WRF dynamical downscaling** Dr. Chuan Yao LIN Media Conf. Room |
| | **Q&A** Conf. Room 1, BHSS, Academia Sinica 12:10 - 12:20 | 11:50 - 12:20 |
| | **The Emergence of Computational Archival Science** | Prof. Richard MARCIANO |
| | Conf. Room 2, BHSS, Academia Sinica | 12:20 - 13:00 |

(Thu 9/3 continued)

| | | | |
|---|---|---|---|
| | **infrastructure - latest developments and** | **proteins in the cloud** Conf. Room 2, BHSS, Academia Sinica | **CNGrid As a HPC Application Cloud Service Provider** |
| | **A solution for secure use of Kibana dCache, towards Federated Identities and** Dr. Paul MILLAR | **Investigating community detection algorithms and their capacity as markers** Prof. Eva HLADKA | **EGI federated platforms supporting accelerated computing** Conf. Room, BHSS, Academia Sinica |
| 17:00 | **A Method for Remote Initial Vetting** Dr. Eisaku SAKANE | **2D and 3D Medical Images for Anatomy Education using a cloud computing platform** Conf. Room 2, BHSS, Academia Sinica | **A Novel Architecture towards Exascale Computing** Media Conf. Room, BHSS, Academia Sinica |
| | **Q&A** | | |

# ISGC Topics

1. Applications from the Virtual Research Communities and Industry
   1. Physics & Engineering applications
   2. Biomedicine & Life Sciences applications
   3. Earth & Environmental Sciences & Biodiversity applications
   4. **Humanities, Arts, and Social Sciences applications**

2. Technologies that Provide Access and Exploitation of Different Site Resources and Infrastructures
   5. Virtual Research Environment (including Middleware, tools, services, workflow, etc.
   6. **Big Data & Data Management**

3. Infrastructure for Research
   7. Networking, Security, Infrastructure & Operations
   8. **Infrastructure Clouds and Virtualization**
   9. **Business Models, Policy, and Sustainability**
   10. Massively Distributed Computing and Citizen Sciences
   11. Supercomputing, High Throughput, Accelerator Technologies and Integration

洲際永久電子文檔管理方案

**Transcontinental Persistent Archive Prototype**

Reagan W. Moore
Richard Marciano
Arcot Rajasekar
Chien-Yi Hou
University of North Carolina, Chapel Hill

Mike Wan
Bing Zhu
Wayne Schroeder
University of California, San Diego

The Grid

**Blueprint for a New Computing Infrastructure**

Edited by Ian Foster and Carl Kesselman

with: Chien-Yi HOU, Sheau-Yen CHEN

# What are Data Grids?

Data Grids are middleware services
- Sitting between the applications and data providers
- Providing transparent and uniform access to diverse types of digital assets
  - Files, databases, streams, web, programs,...
  - Documents, images, data, sensor packets, tables,...
- From heterogeneous resources
  - File Systems, tape archives, sensor streams,...
- Distributed over a wide area network
  - Multiple administrative and security domains
- With users unaware of physical attributes of the data access
  - System addresses, paths, protocols, ...

# Using a Data Grid – *in Abstract*

*Data Grid*

Ask for data

Data delivered

- User asks for data from the data grid
  - The data is found and returned
  - Where & how details are hidden

8

# Digital Curation Innovation Center (DCIC) @ U. Maryland (USA)

**Mission:**

- Be a leader in the digital curation research and educational fields, and foster interdisciplinary partnerships using **Big Records and archival analytics** through public / industry / government collaborations.

- Sponsor interdisciplinary projects that explore the integration of archival research data, user-contributed data, and technology to **generate new forms of analysis and historical research engagements.**

## ARC: Archives Research & Collaboration Lab

*Director*: Ricky **Punzalan**

ARC studies and develops innovative approaches, systems, strategies, and tools to foster sustainable futures for archives, preservation, and digital curation.

http://archivescollaboratory.umd.edu/

## SALT: Sustainable Archives & Leveraging Technologies

*Director*: Richard Marciano

SALT is an interdisciplinary lab, which focuses on the long-term preservation of digital cultural and research assets at scale. SALT is an acronym for SustainA-biLiTy and uses as its logo the two thousand year-old ancient Chinese pictograph for salt ("yan") which is a metaphor for the integration of policy, governance, infrastructure, and content.

http://salt.umd.edu

---

## curatelab

**Hornbake South 4110**
Digital lab for group learning, collaborative design, and hands-on digital curation project development (23 seats, 3 interactive screens, 12 workstations with 12TB of storage).

## digitizationlab

**Hornbake South 4110D**
Document scanning, image manipulation, and archival ingestion facility for group projects.

## serverfarm

**UMD Computer & Space Sci. Bldg**
On-campus virtual machine farm for research data processing, storage, and hosting (15TB storage, 2 Dell servers, VMWare-powered).

## cloudlab

**Amazon Cloud**
Dashboard-enabled virtual computing lab in the cloud for creating Windows/Ubuntu instances using Amazon Web Services (AWS).

## dataCave

**UMD Cyberinfrastructure Center at the Rivertech Bldg**

Building **DRAS-TIC**

Digital Repository At Scale That Invites Computation (To Improve Collections): a peta-scale archival storage and preservation repository (based on DRAS-TIC open-source software (NoSQL Cassandra database) and computational infrastructure (4 Dell nodes).

---

# dcic

## digital curation
## innovation center

http://dcic.umd.edu

## Mission

Be a leader in the digital curation research and educational fields, and foster interdisciplinary collaborations using Big Records and archival analytics with public / industry / government partnerships.

Sponsor interdisciplinary projects that explore the integration of archival research data, user-contributed data, and technology to generate new forms of analysis and historical research engagements, particularly in the arenas of social justice, human rights, and cultural heritage

COLLEGE OF
**INFORMATION**
STUDIES

# Projects

## Cyberinfrastructure for the curation & management of digital assets at scale:

### "Brown Dog"

A CIC Big-10 $10.5M NSF/DIBBs-funded collaboration with U. of Illinois NCSA Supercomputing Center and industry partners (NetApp and Archive Analytics Solutions). This project aims to help accelerate the development of digital curation processes and services and create a data observatory to provide access to Big Records training sets and teach students practical digital curation skills.

### "Curate Cloud"

A $300K IMLS-funded project that helped launch a new online professional education certificate for digital curation professionals, the Curation and Management of Digital Assets (CMDA). Curate Cloud is also developing an open-source research and educational platform, the VCL (Virtual Computing Lab), to remove barriers to access for curation tools and resources.

## Digital Curation training:

### Digital Curation Fellowships

The iSchool has several Fellowship opportunities for students in digital curation and archives. These include a collaboration with the **National Agricultural Library** (NAL); extensive project work with the **National Park Service** (NPS); and a scholarship established in honor of **Bruce Ambacher**, retired senior archivist and iSchool faculty member.

### Interdisciplinary Research Teams

Gain new digital skills, conduct interdisciplinary research, explore professional development opportunities at the intersection of archives, big data, and analytics through a number of project themes: *Refugee Narratives, Community Displacement, Racial Zoning, Cyberinfrastructure for Digital Curation, Movement of People, Citizen Internment.*

# People

## Research Staff:

| | |
|---|---|
| Richard Marciano | Director & SALT Lab Director |
| Michael Kurtz | Associate Director |
| Ricardo Punzalan | Research Associate & ARC Lab Director |
| Ken Heger | Research Associate & DigitizationLab Director |
| Greg Jansen | Research Software Architect |
| Maria Esteva | Affiliate Professor |
| Victoria Lemieux | Affiliate Professor |
| William Underwood | Affiliate Professor |



Mary Kendig / Myeong Lee    Graduate Research Assistants



## Research Affiliates:

### U. Maryland
Tammy Clegg, Nick Diakopoulos, Jesse Johnston, Trevor Owens, Jenny Preece, Katie Shilton

### External
Bruce Ambacher, Natalie Baur, John Burns, Andrew Lau, Scott Madry

## Postdoctoral Fellows:
Morgan Daniels, Kathryn Gucer, Adam Kriesberg
(Advisor: Punzalan)

## Students (Undergraduate, Master's [MLIS, MIM, HCIM], Doctoral):
Maddie Allen, Saba Al-Dughaither, Vinila Atre, Myuresh Amdekar, Richard Bool, Arpit Chandra, Shiyun Chen, John Dela Cruz, Anne Dempsey, Shaina Destine, Kelsey Diemand, Pai Doshi, Erin Durham, Will Frolikdong, Alicia Geller, Karishma Ghiya, Janet Glazier, Rajesh Gnanasekaran, Rhett Greenfield, Allison Gunn, Ashley Huddix, Scott Harkless, Torra Hausmann, Eric Hung, Hardik Jhaveri, Ruchira Kapoor, Amar Kurane, Yuting Liao, Zhenye Ma, Sheryl Mathias, Paridhi Mathur, Martin Moreno, Jennifer Proctor, Brian Redford, Darlene Reyes, Benjamin Sagey, Sohan Shah, Jay Sheth, Niraj Shirame, Edel Spencer, Akash Udani, Sydney Valle, Jennifer Wachtel, Melissa Wertheimer, Meaghan Wilson, Jiahui Wu, David Zhang, Xinyun Zhang

## Doctoral Students:

| | |
|---|---|
| Andrew Casertano, Will Thomas | (Advisor: Marciano) |
| Diane Travis | (Advisors: Butler/Marciano) |
| Edward Summers, Amy Wickner | (Advisor: Punzalan) |

# Projects

## Justice, Human Rights, & Cultural Heritage:

### Overseas Pension Project

A student- and professional society-driven project to collect information documenting payment of pensions to American veterans living overseas. The project creates datasets documenting migration patterns, the flow of money, health conditions, and family connections prior to World War I.

### International Research Portal Project (IRP²)

This project will improve access to an important tool which identifies and locates looted art and other cultural assets found on the *International Research Portal for Records Related to Nazi-Era Cultural Property.*

### Mapping the Voyage of the St. Louis

In 1939, 937 passengers (mostly Jews) fled Germany aboard the SS St. Louis ship, heading to Cuba, where they were turned away and forced to return to Europe where 254 were killed during the Holocaust. The project looks at mapping individual and collective stories through graph database techniques.

### Japanese-American WWII Camps

Building on a UMD FIA Seed Grant, the project explores the integration of archival and user-contributed data using social networking graphs to link people, places, and events. Using WWII Camp data.

### Mapping Inequality

A project with Johns Hopkins, Virginia Tech, and U. of Richmond where a national collection of New Deal redlining records is being crowdsourced (these unique records capture racial, ethnic, and economic conditions).

### The Human Face of Big Data

A student-led project that will create access and collaborative opportunities around historically and socially- significant heterogeneous datasets rooted in urban renewal housing records for a number of cities.

# RACE, SPACE, and PLACE

## Refugee Narratives

### The Journey of Refugees: What Happened After St. Louis Voyage
Sohan Shah, Yuting Liao, Ruchira Kapoor, Pal Doshi, and Mary Kendig

COLLEGE OF INFORMATION STUDIES

## Racial Zoning

### Mapping Inequality: Redlining in New Deal America

## Citizen Internment

### Computational Linguistics & Graph Analytics: Archives in the Age of Big Data

---

## Human Face of Big Data: Urban Renewal in Asheville, NC

## Reusability, Digital Reunification and Analysis of US Overseas Pension Records

### Community Displacement

### Movement of People

---

# dcic digital curation innovation center

## INNOVATIVE CYBER-INF.

### Virtual Learning

Brown Dog: Cloud Services for Auto-Curation

DRAS-TIC

Building Virtual Learning Environments: Virtual Computing Lab (VCL)

**Curation At Scale**

---

# DIGITAL CURATION

## Local Gov.

### Historic Survey of the Eastport Community: a Partnership for Action Learning in Sustainability (PALS)

COLLEGE OF INFORMATION STUDIES

## Federal Gov.

### Digital Curation at the National Agricultural Library

USDA NAL

COLLEGE OF INFORMATION STUDIES

## Health Care

### RAPTOR: Radiology Protocol Tool & Recorder
Andrew Casertano, Richard Marciano

COLLEGE OF INFORMATION STUDIES

# Pursuing a CAS Training / Teaching Agenda

There is a need to :

- create innovative classes that emphasize new modes of collaboration, and interdisciplinary work.
- blend elements of archival thinking and computational thinking:
    - problem solving that uses modeling, decomposition, pattern recognition, abstraction, algorithm design, and scale.
- develop inter-disciplinary iSchools with faculty from Computer Science, Archival Science, and Data Science.
- develop extensive hands-on experience working with cyberinfrastructure to carry out archival functions.

# How Each Project Is Related to Computational Archival Science (CAS) Themes:

| Project | Computational Linguistics | Data Modeling & Evolutionary Prototyping | Graph Analytics | Crowdsourcing | GIS |
|---|---|---|---|---|---|
| 1. Human Face of Big Data [Community Displacement] | | X | | X | X |
| 2. Mapping Inequality [Racial Zoning] | X | | | X | X |
| 3. St. Louis Voyage [Refugee Narratives] | | X | X | | X |
| 4. World War II Japanese Camps [Citizen Internment] | X | X | X | X | X |

# IEEE Big Data 2016
## "Computational Archival Science: digital records in the age of big data

**http://dcicblog.umd.edu/cas/ieee_big_data_2016_cas-workshop/**
Dec. 2016 workshop
**http://dcicblog.umd.edu/cas**
April 2016 workshop

**Upcoming… IEEE Big Data 2017 in BOSTON**

- Mark Hedges, Tobias Blanke, KCL
- Bill Underwood, GTRI (now UMD)
- Victoria Lemieux, UBC
- Maria Esteva, TACC
- Richard Marciano, Michael Kurtz, UMD

- **Application of analytics to archival material**, including text-mining, data-mining, sentiment analysis, network analysis.
- **Analytics in support of archival processing**, including appraisal, arrangement and description.
- **Scalable services for archives**, including identification, preservation, metadata generation, integrity checking, normalization, reconciliation, linked data, entity extraction, anonymization and reduction.
- **New forms of archives**, including Web, social media, audiovisual archives, and blockchain.
- **Cyber-infrastructures for archive-based research** and for development and hosting of collections
- **Big data and archival theory** and practice
- Digital curation and preservation
- Crowdsourcing and archives
- **Big data and the construction of memory** and identity
- Specific big data technologies (e.g. NoSQL databases) and their applications
- Corpora and reference collections of big archival data
- Linked data and archives
- Big data and **provenance**
- Constructing big data research objects from archives

## Our working definition of Archival Computational Science (CAS):

*An interdisciplinary field concerned with the application of computational methods and resources to large-scale records/archives processing, analysis, storage, long-term preservation, and access, with aim of improving efficiency, productivity and precision in support of appraisal, arrangement and description, preservation and access decisions, and engaging and undertaking research with archival material.*

*e.g.: NSF/SBE RIDIR, LOC National Digital Initiative, IMLS Always Already Computational, etc.*

**DataNet: $62.6M**

**SEAD:** $8.1M (Michigan): data curation software and services for the "long tail" of small- and medium-scale data producers in sustainability science.
**Terra Populus:** $8.2M (Minnesota) -- Build tools for data integration across the domains of social science and environmental data
**DFC:** $8.3M (North Carolina) -- Use the integrated Rule-Oriented Data System (iRODS) to provide data grid infrastructure for science and engineering.
**DataONE:** $27.9M (New Mexico) -- $20M + $7.9M Oct. 2014 --  platform for collaborative environmental and ecological science
**DataConservancy:** $10M – Johns Hopkins U. – 2009-2014

**DIBBS: 3 1/2 years: $115 M**

**====== Fall 2013:  $32.8M**

**NSF Office of Advanced CI (part of CISE)**
**DataNet &DIBBs:**
**~ $178M of National Investments ($115M for DIBBs)**

| Brown Dog | -- University of Illinois at Urbana-Champaign / U. Maryland | -- $10.5M |
|---|---|---|
| Data Exacell | -- Carnegie Mellon University | -- $8.9M |
| SkyServer | -- Johns Hopkins University | -- $8.9M |
| GABBs | -- Purdue University | -- $4.5M |

**====== Fall 2014:  $20.8M**
-**Building a Modular Cyber-Platform for Systematic Collection, Curation, and Preservation of Large Engineering and Science Data** -- Purdue University -- $1.5M
-**User Driven Architecture for Data Discovery** -- Corporation for National Research Initiatives (NRI) -- $1.5M
-**Collaborative Research: Cyberinfrastructure for Interpreting and Archiving U-series Geochronologic Data** -- College of Charleston -- $580K
-**T2-C2: Timely and Trusted Curator and Coordinator Data Building Blocks** -- University of Illinois at Urbana-Champaign -- $1.5M
-**Scalable Capabilities for Spatial Data Synthesis** -- University of Illinois at Urbana-Champaign -- $1.5M
-**Domain-Aware Management of Heterogeneous Workflows: Active Data Management for Gravitational-Wave Science Workflows** -- Syracuse University -- $750K
-**SPIDAL: Middleware and High Performance Analytics Libraries for Scalable Data Science** -- Indiana University -- $5.1M
-**Ubiquitous Access to Transient Data and Preliminary Results via the SeedMe Platform** -- University of California-San Diego -- $1.3M
-**DIBBs for Intelligence and Security Informatics Research Community** -- University of Arizona -- $1.5M
-**STORM: Spatio-Temporal Online Reasoning and Management of Large Data** -- University of Utah -- $1.2M
-**Systematic Data-Driven Analysis and Tools for Spatiotemporal Solar Astronomy Data** -- Georgia State University Research Foundation -- $1.5M
-**An Infrastructure for Computer Aided Discovery in Geoscience** -- Massachusetts Institute of Technology -- $1.4M
-**Porting Practical Natural Language Processing (NLP) and Machine Learning (ML) Semantics** -- University of Colorado at Boulder -- $1.5M

**====== 2015:  $27.5**
-**Tripal Gateway, a Platform for Next-Generation Data Analysis and Sharing** -- Washington State University – $1.5M
-**An Integrated System for Public/Private Access to Large-Scale, Confidential Social Science Data** --Duke University -- $1.5M
-**LearnSphere: Building a Scalable Infrastructure for Data-Driven Discovery and Innovation in Education** -- Carnegie-Mellon University -- $4.8M
-**An Infrastructure Supporting Collaborative Data Analytics Workflow Design and Management** -- Carnegie-Mellon University -- $1M
-**DNI: Give Your Data the Edge: A Scalable Data Delivery Platform** -- University of Arizona -- $3.8M
-**DNI: Multi-Institutional Open Storage Research InfraStructure (MI-OSiRIS)** --University of Michigan Ann Arbor -- $4.9M
-**DNI: Data Analysis and Management Building Blocks for Multi-Campus Cyberinfrastructure through Cloud Federation** -- Cornell University -- $5M
-**DNI: The Pacific Research Platform** -- University of California-San Diego -- $5M

**====== 2016:  $31.4M**
-**EI: Virtual Data Collaboratory: A Regional Cyberinfrastructure for Collaborative Data Intensive Science** – Rutgers -- $4M
-**EI: Data Laboratory for Materials Engineering** – SUNY at Buffalo-- $2.9M
-**EI: mProv: Provenance-based Data Analytics cyberinfrastructure for High-frequency Movile Sensor data** – U. Memphis -- $4M
-**EI: Merging Science and Cyberinfrastructure Pathways: The Whole Tale** -- University of Illinois at Urbana-Champaign -- $5M
-**PD: Ontology-Enabled Polymer Nanocomposite Open Community Data Resource** -- Rensselaer Polytechnic Institute -- $500K
-**EI: The Local Spectroscopy Data Infrastructure (LSDI)** – UC Berkeley -- $3.9M
-**EI: VIFI:Virtual Information-Fabric Infrastructure (VIFI) for Data-Driven Decisions from Distributed Data** --  UNC Charlotte -- $4M
-**PD: Metadata Toolkits for Building Multi-Faceted Data – Relationship Models** – MIT -- $500K
-**EI: Continuous Capture of Metadata for Statistical Data** – U. Michigan -- $2.6M
-**EI: North East Storage Exchange** – Harvard U. -- $4M

**====== 2017:  $3.2M so far…**
-**EI: Vizier, Streamlined Data Curation** – SUNY at Buffalo -- $2.7M
-**PD: Accelerating Comparative Metagenomics through an Ocean Cloud Commons** – U. Arizona -- $500K

This solicitation includes two classes of science data pilot awards:
- **Early Implementations** are large "at scale" evaluations, building upon cyberinfrastructure capabilities of existing research communities or recognized community data collections, and extending those data-focused cyberinfrastructure capabilities to additional research communities and domains with broad community engagement.

- **Pilot Demonstrations** address advanced cyberinfrastructure challenges across emerging research communities, building upon recognized community data collections and disciplinary research interests, to address specific challenges in science and engineering research.
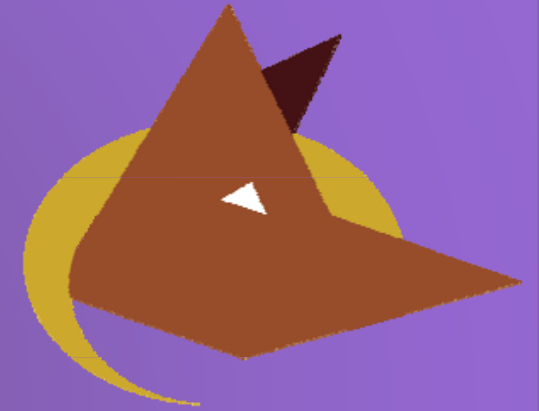
# Projects with Preservation Focus

**Integration Pilots**
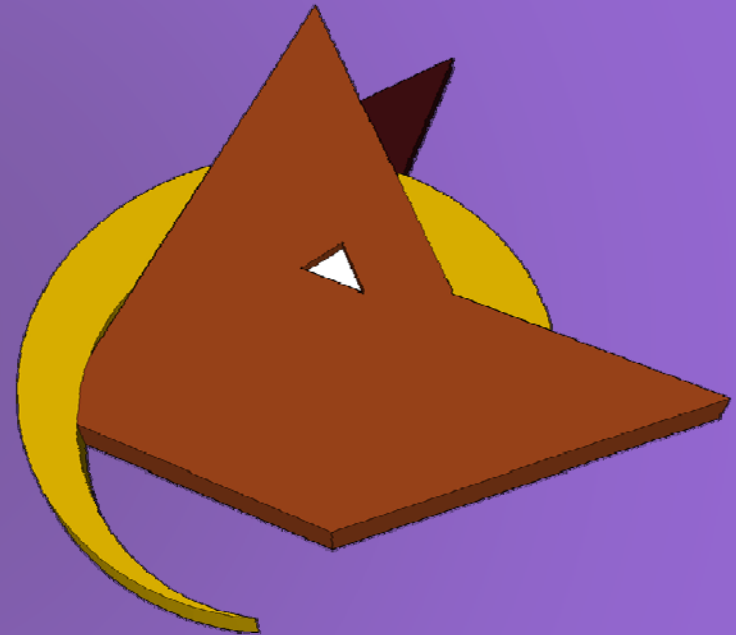
**Syndicate**
*Larry Peterson*
$3,804,911
2015-2019
THE UNIVERSITY OF ARIZONA
**I/O workflows**
*Software: Cyverse, iRODS, OpenCloud, Hadoop*

**Whole Tale**
*Bertram Ludaescher*
$4,986,951
2016-2021
**Interfaces to data/software**
*Software: Globus, Jupyter, iRODS, ownCloud, DataONE, Brown Dog*

executable papers, provenance preservation
data transformation, software preservation / publishing

**Tool Interface**

social curation, auto curation, publishing data

**Aristotle**
*David Lifka*
$4,975,094
2015-2020
**Cloud resource broker**

**SPIDAL**
*Geoffrey Fox*
$5,000,000
2014-2019
**HPC optimized tools, frameworks and analysis tools**
*Software: Apache Software*

**GABBs**
*Xiaohui Carol Song*
$3,409,029
2013-2018
PURDUE
**Geospatial analysis and tools**
*Software: HUBzero*

**BROWN DOG**
*Kenton McHenry*
$10,519,716
2013-2018
**Data transformation and tool publishing**
*Software: Clowder, Polyglot, Versus, Daffodil*

**SEAD**
*Margaret Hedstrom*
$8,000,000
2011-2016
UNIVERSITY OF MICHIGAN
**Data sharing, curation, publishing**
*Software: Clowder, Virtual Archive*

iRODs, data management

data curation and sharing

**Data Interface**

**Data Providers**

**The Data Exacell**
Carnegie Mellon University
*Michael Levine*
$4,902,601
2013-2018
**Wide area file system**
*Software: Slash2*

**DFC** DataNet FEDERATION CONSORTIUM
*Reagan Moore*
$8,300,992
2011-2016
**Policy management leaving data stays where its at**
*Software: iRODS*

**DataONE**
UNM
*Bill Michener*
$21,194,548
2009-2014
**Earth science data**
*Software for: federated storage, tools to curate data*

**TERRA POPULUS**
*Steven Ruggles*
$7,993,266
2011-2016
**Census & raster data integration**

**SciServer**
*Alex Szalay*
$7,603,723
2013-2018
JOHNS HOPKINS
**Astronomy data**
*Software for: dropbox like sharing*

**Storage Abstraction**

**Resource Interface**

**XSEDE**
Extreme Science and Engineering Discovery Environment
**Computation**

# The Problem Addressed by Brown Dog

○ Large collections of **un-curated** and/or **unstructured** digital data ("long-tail" data)

  ○ Many file formats
  ○ No metadata
  ○ No useful filenames
  ○ No useful directory structure
  ○ No textual contents

# What Is Needed

○ Means of indexing data contents so that large collections of data can be searched and desired data found

    ○ An ability to compare data

# Brown Dog Data Transformation Services

- The Data Access Proxy (DAP)
  - http://dap.ncsa.illinous.edu/conversion/:output/:file
  - File in, File out

- The Data Tilling Service (DTS)
  - http://dts.ncsa.illinois.edu/extraction/:domain/:file
  - File in, JSON out
  - JSON can contain metadata, tags, signatures, links to derived data products, etc…

# Brown Dog Use Cases

○ Addressed specifically here:

1. Biology/Ecology
2. Civil and Environmental Engineering
3. Social Science

○ Testbed data:

4. UMD CI-BER testbed, at the U. Maryland iSchool

# Brown Dog

## The Data Tilling Service (DTS)

# Data Tilling

○ Data Tilling (v): To prepare and cultivate (*data*) for *analysis*

○ Data Tillage (n): Is the *computational* preparation of *data* by *algorithmic* agitation of various types, such as digging, stirring, and overturning

# (Pre) Data Analysis

○ Not necessary *data cleaning*

○ More like *metadata extraction*

○ Not full analysis / Not perfect results

○ Apply as many methods as possible

○ Support the user in finding the metadata they need

# Extracting Information from Raw Files

File / URL → DTS →

- Arbitrary metadata
- Previews
- Tags
- Content based signatures

# Extractors

# [M3d] User Case Extractors

○ Floodplain extraction

○ Pond extraction from aerial photos

○ Gap filled versions of the data

○ Text extraction from digitized documents (in particular numerical values)

○ River locations from hand drawn maps

○ River locations from aerial photos

○ Route/image extraction

○ Geolocation

○ Green Index extractor

○ Human preference from images

○ Sentiment analysis from text

○ Data extraction from articles (e.g. tables)

# [M3d] Other Extractors

- OpenCV
  - Faces, eyes
- Tika
  - Language detection
- Simple Summary
  - Summaries
- Cell Profiler
  - Human, yeast, fly, tumor, …

- Tesseract
  - Text extraction from images
- CMU Sphinx
  - Speech recognition
- VLFeat
  - Plane, motorcycle, …

**Data Collection**
*URL, File System, …*

**Native Byte Encoding**
*File Formats, Data Bases, Websites, Documents*

**DAP**

**Data Structures**
*Arrays, Strings, Images, Videos, Audio, 3D Models, …*

**DTS**

**Derived Data/ Metadata**
*Tags, Signatures*

**Applications**
Search, Relate, View, Process, Use

**Usable Data**

*M. Dietze, Ecology, Boston University*
*A. Desai, Ecology, University of Wisconsin*
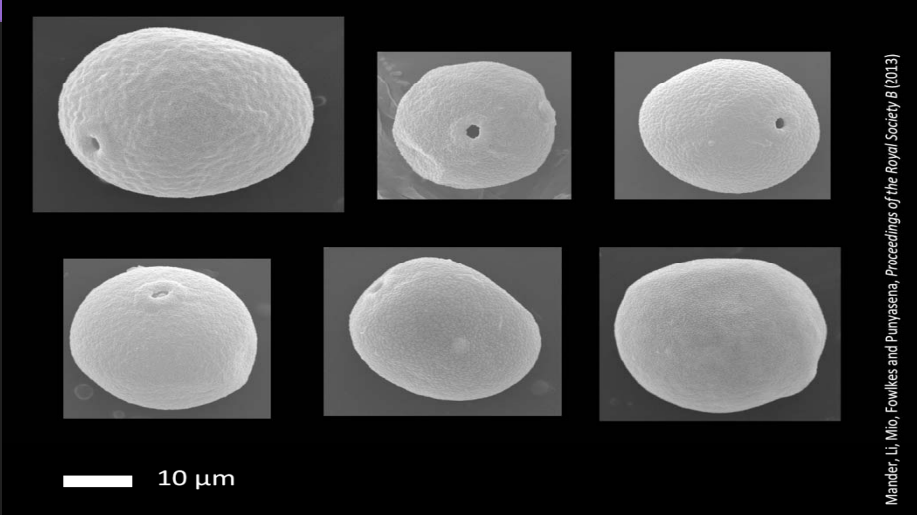*D. LeBauer, Ecology, University of Illinois Urbana-Champaign*

**Data Collection**
*Handwritten*
*Settlement Vegetation*
*Data*

**Native Byte Encoding**
*Various Image Formats*

**DAP**

**Data Structures**
*Image*

**DTS**

**Derived Data/ Metadata**
*Text, Number Values*

**Applications**
Climate Modeling

**Usable Data**

T.37.N. R.11.E. 2nd. Mer. Indiana.

Chs. lks.

a B. Oak 12 in. diam. bs. N.85.E. 15 lks. dis.
a Do 12 " " S.40.W. 20 " "
79.97 to Sec. Cor. Apud 24th

North Between Section 10 & 11
15.00 Left Prairie
22.75 a Gum 10 in. diam.
28.50 Road to Indiana
30.00 a Marsh
36.30 Canal Line
40.00 Set 2r Sec. post pr. ch.
a W. Oak 10 in. diam. bs. S.10.E. 365 lks. dis.
a B. Oak 15 " " S.35.W. 310 "
42.50 Pigeon River 40 lks. wide co. N.W.
47.00 Left Marsh
72.70 a W. Oak 12 in. diam.
80.00 Set post Cor. to Sec. 2.3.10 & 11 pr. ch.
a W. Oak 16 in. diam. bs. S.15.W. 200 lks. di.
a Do 12 " S.65.W. 144 "
Land &c. Same

East On Random Betw. Sec. 2 & 11
14.00 Set Temp'y post
45.00 Wet Prairie
58.00 Left Do
79.90 Intersect. N. & S. Line 16 lks. N. of post
Land &c. Same

West On True Line
39.94 Set 2r Sec. post pr. ch.
a W. Oak 12 in. diam. S.55.E. 22 lks. di.
a Do 8 " " S.55.W. 32 " "
79.88 to Sec. Cor.

North Between Section 2 & 3
23.10 a B. Oak 14 in. diam.

T.37.N. R.11.E. 2nd. Mer. Indiana.

Chs. lks.
25.00 a Marsh
26.00 Left Do
40.00 Set 2r Sec. post pr. ch.
a W. Oak 18 in. diam. bs. N.2.W. 75 lks. di.
a R. Oak 20 " " S.38.E. 40 " "
51.14 Intersect. R. Boundy. 75 lks. of post
Set post at intersection pr. ch.
a Hickory 10 in. diam. bs. N.9.E. 230 dis.
a B. Oak 12 " " S.55.W. 275 "
Land &c. Same Apud 25th

North Between Section 33 & 34
33.50 a W. Oak 10 in. diam.
40.00 Set 2r Sec. post pr. ch.
a B. Oak 20 in. diam. bs. N.5.W. 20 lks. di.
a Do 20 " " S.15.E. 18 "
69.80 a W. Oak 18 "
80.00 Set post Cor. to Sec. 27.28.33 & 34 pr. ch.
a B. Oak 12 in. diam. bs. N.9.E. 90 lks. di.
a Do 40 " " N.70.W. 70 "
Land rolling 1st Rate — Timber
Ash, Walnut, Poplar, Oak &c —

East On Random Betw. Sec. 27 & 34
36.00 a Creek 50 lks. wide co. N.W.
41.00 Set Temp'y post
80.00 Intersect. N. & S. Line 10 lks. S. of post
Land 2nd Rate Timber
B. & W. Oak. Sassafras & Grapevine

West On True Line
40.00 Set 2r Sec. post pr. ch.
an L. Ash 20 in. diam. bs. S.40.E. 22 lks. di.
a Sycamore 20 " " N.29.W. 10 "
80.00 to Sec. Cor.

*B. Minsker, Civil & Env. Engineering, University of Illinois Urbana-Champaign*
*A. Schmidt, Civil & Env. Engineering, University of Illinois Urbana-Champaign*
*B. Sullivan, Landscape Architecture, University of Illinois Urbana-Champaign*

**Data Collection**
*Satellite/Aerial Data*

**Native Byte Encoding**
*Various Image Formats*

**DAP**

**Data Structures**
*Image*

**DTS**

**Derived Data/ Metadata**
*Land Cover/Usage/*

**Applications**
Human Access to Green Infrastructure

**Usable Data**

*P. Kumar, Civil & Env. Engineering, University of Illinois Urbana-Champaign*

**Data Collection**
*LiDAR Data*

**Native Byte Encoding**
*LAS*

**DAP**

**Data Structures**
*Depth Data*

**DTS**

**Derived Data/ Metadata**
*Floodplains*

**Applications**
Flood Plain Analysis

**Usable Data**

*P. Kumar, Civil & Env. Engineering, University of Illinois Urbana-Champaign*

**Data Collection**
*LiDAR Data*

**Native Byte Encoding**
*LAS*

**DAP**

**Data Structures**
*Depth, Polygons*

**DTS**

**Derived Data/Metadata**
*Floodplains, Depth Distribution*

**Applications**
Flood Plain Analysis

**Usable Data**

*P. Kumar, Civil & Env. Engineering, University of Illinois Urbana-Champaign*

**Data Collection**
*Digitized 19th Century Maps*

**Native Byte Encoding**
*Various Image Formats*

**DAP**

**Data Structures**
*Image*

**DTS**

**Derived Data/ Metadata**
*River/Stream Locations*
**Applications**
River Meander

**Usable Data**

Township No 19 North of the Baseline, Range No 6 East of the 3rd principal Meridian

Office of the Surveyor General for Illinois and Missouri
St. Louis March 13th 1857

*B. Minsker, Civil & Env. Engineering, University of Illinois Urbana-Champaign*
*A. Schmidt, Civil & Env. Engineering, University of Illinois Urbana-Champaign*
*B. Sullivan, Landscape Architecture, University of Illinois Urbana-Champaign*

**Data Collection**
*Landscape/Architecture/Design Images*

**Native Byte Encoding**
*Various Image Formats*

**DAP**

**Data Structures**
*Image*

**DTS**

**Derived Data/Metadata**
*Human Preference Score*

**Applications**
Green Infrastructure Design

**Usable Data**

*M. Poole, Social Science, University of Illinois Urbana-Champaign*
*F. Pena-Mora, Civil & Env. Engineering, Columbia University*
*D. Espelage, Education, University of Illinois Urbana-Champaign*

**Data Collection**
*Groupscope*

**Native Byte Encoding**
*Various Video Formats*

**DAP**

**Data Structures**
*Video*

**DTS**

**Derived Data/ Metadata**
*People Locations/ Interactions*

**Applications**
Large Dynamic Group Behavior

**Usable Data**

Filename: current_results.mp4
Resolution: 1920×1080
Duration: 0:57

*K. Nahrstedt, Computer Science,  University of Illinois Urbana-Champaign*
*J. Rogers, Material Science, University of Illinois Urbana-Champaign*

**Data Collection**
*Material Fabrication Data*

**Native Byte Encoding**
*SEM Images*

**DAP**

**Data Structures**
*Image*

**DTS**

**Derived Data/ Metadata**
*Successful/Failed Experiments*

**Applications**
Data Mining towards future Materials Development

**Usable Data**

1.495μm

5.0kV 7.6mm x45.0k

10.0kV 7.4mm x9.00k                    5.00um

*Industry*

**Data Collection**
*Digitized Technical Drawings*

**Native Byte Encoding**
*Various Image Formats*

**DAP**

**Data Structures**
*Image*

**DTS**

**Derived Data/ Metadata**
*Shape Descriptors, Metadata*

**Applications**
Indexing

**Usable Data**

| | | *Croton hirtus* | | *Mabea occidentalis* | | *Agropyron repens* | |
|---|---|---|---|---|---|---|---|
| | | Raw | Deconvolved | Raw | Deconvolved | Raw | Deconvolved |
| Widefield | | | | | | | |
| Apotome (raw) | | | | | | | |
| Confocal 405 nm | | | | | | | |
| Confocal 561 nm | | | | | | | |
| Two-photon 780 nm | | | | | | | |

*E. Spalding, Botany, University of Wisconsin*

**Data Collection**
*Seedling Images*

**Native Byte Encoding**
*Various images formats*

**DAP**

**Data Structures**
*Image*

**DTS**

**Derived Data/Metadata**
*Root tip, Root tip characteristics*

**Applications**
Phenomics

**Usable Data**

Stephan Joslyn, Veterinary Medicine, University of Illinois Urbana-Champaign

**Data Collection**
*Aerial photos of cattle feedlots*

**Native Byte Encoding**
*Images formats*

DAP

**Data Structures**
*Image*

DTS

**Derived Data/Metadata**
*Cattle detection, cattle clustering*

**Applications**
Disease detection, Food supply

**Usable Data**

*David LeBauer, Ecology, University of Illinois Urbana-Champaign*

**Data Collection**
*Biofuel Research Publications*

**Native Byte Encoding**
*Document formats*

**DAP**

**Data Structures**
*Tables, Images*

**DTS**

**Derived Data/metadata**
*Plant species, locations, yields*

**Applications**
Biofuel production, Environmental Impact

**Usable Data**

| Indicators | $N_0$ | $N_{60}$ | $N_{120}$ | $LSD_{05}$ | $N_0$ | $N_{60}$ | $N_{120}$ | $LSD_{05}$ |
|---|---|---|---|---|---|---|---|---|
| | | | Second year of growth | | | | | |
| | | | Annual biomass | | | | | |
| t ha$^{-1}$ | 15.8 | 20.0 | 24.7* | 5.63 | 4.62 | 6.55* | 6.82* | 1.85 |
| % | 100 | 127 | 157 | 40.5 | 100 | 142 | 148 | 42.9 |
| | | | Biomass weight per plant | | | | | |
| kg | 1.66 | 2.26* | 2.53* | 0.627 | 0.49 | 0.74* | 0.70* | 0.208 |
| | | | Third year of growth | | | | | |
| | | | Annual biomass | | | | | |
| t ha$^{-1}$ | 27.0 | 28.5 | 29.7 | 5.31 | 10.5 | 10.7 | 11.5 | 2.47 |
| % | 100 | 105.6 | 110.1 | 18.68 | 100 | 102.3 | 110.2 | 22.61 |
| | | | Biomass weight per plant | | | | | |
| kg | 2.05 | 2.18 | 2.25 | 0.396 | 0.79 | 0.81 | 0.87 | 0.187 |

| | A | B | C | D | E | F | G | H | I | J | K |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | Site | Date | Species | Genotype | Type | Block | IRGA | Curve | | Topt | PAR |
| 2 | NY | 5/24/13 | Willow | SX61 | Sun | 1 | Dietze | Temp | | 27.25 | 1500 |
| 3 | NY | 5/25/13 | Willow | SX61 | Sun | 1 | Dietze | Temp | | Bad Curve | 1500 |
| 4 | NY | 5/24/13 | Willow | SX61 | Sun | 2 | Dietze | Temp | | Bad Curve | 1500 |
| 5 | NY | 5/25/13 | Willow | SX61 | Sun | 2 | Dietze | Temp | | 26.332 | 1500 |
| 6 | NY | 5/24/13 | Willow | SX61 | Sun | 3 | Dietze | Temp | | Bad Curve | 1500 |
| 7 | NY | 5/24/13 | Willow | FC | Sun | 1 | USDA | Temp | | Bad Curve | 1500 |
| 8 | NY | 5/25/13 | Willow | FC | Sun | 1 | USDA | Temp | | Bad Curve | 1500 |
| 9 | NY | 5/24/13 | Willow | FC | Sun | 2 | USDA | Temp | | 23.6 | 1500 |
| 10 | NY | 5/25/13 | Willow | FC | Sun | 2 | USDA | Temp | | 25.9 | 1500 |
| 11 | NY | 5/24/13 | Willow | FC | Sun | 3 | USDA | Temp | | Bad Curve | 1500 |

| | A | B | C | D | E | F | G | H | I | J | K | L | M |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | Miscanthus | MCB | SD | MNB | SD | MCM | SD | MNM | SD | MCT | SD | MNT | SD |
| 2 | Jun | 19.73 | 7.48 | 21.28 | 2.38 | | | | | 33.02 | 6.44 | 40.35 | 7.26 |
| 3 | Jul | 11.45 | 2.96 | 18.88 | 13.13 | 15.99 | 0.70 | 9.69 | 1.74 | 27.83 | 3.46 | 31.45 | 3.87 |
| 4 | Aug | 10.37 | 3.68 | 13.00 | 2.14 | 21.25 | 7.18 | 15.81 | 3.76 | 27.61 | 8.68 | 22.01 | 7.21 |
| 5 | Sep | 12.37 | 5.06 | 12.24 | 4.15 | 15.85 | 2.07 | 22.80 | 3.19 | 23.48 | 3.63 | 21.93 | 3.58 |
| 6 | Oct | 23.16 | 10.35 | 20.27 | 11.75 | 17.06 | 2.08 | 18.55 | 9.11 | 31.51 | 4.08 | 30.41 | 2.19 |
| 7 | | | | | | | | | | | | | |
| 8 | | | | | | | | | | | | | |
| 9 | | | | | | | | | | | | | |
| 10 | Switchgrass | CB | SD | NB | SD | CM | SD | NM | SD | CT | SD | NT | SD |
| 11 | Jun | | | | | | | | | 25.62 | 4.38 | 29.43 | 5.30 |
| 12 | Jul | 10.53 | 0.89 | 11.44 | 2.42 | | | | | 23.05 | 2.09 | 20.41 | 4.14 |
| 13 | Aug | 13.37 | 2.97 | 11.96 | 3.94 | | | | | 12.61 | 3.41 | 23.91 | 5.92 |
| 14 | Sep | 10.71 | 1.01 | 11.97 | 3.37 | | | | | 12.48 | 4.95 | 8.83 | 1.51 |

Batteries

Stabilizer 6V 1Amp

12V
2.6 Amp

280 Ohm

1000 Ohm

Halog

Flash

Com

Structure, species of individual trees

Maps of leaf area index (LAI); Landscape biomass, carbon

*Tim Gernat, Entomology, University of Illinois Urbana-Champaign*
*Gene Robinson, Entomology, University of Illinois Urbana-Champaign*

**Data Collection**
*Bee Hive Images/Video*

**Native Byte Encoding**
*Image and Video formats*

**DAP**

**Data Structures**
*Image, Image Sequences*

**DTS**

**Derived Data/ Metadata**
*Bee locations, movement, interactions*

**Applications**
*Social Interactions of Bees*

**Usable Data**

*Amelia Bartholomew, Medicine, University of Illinois Chicago*

**Data Collection**
*Kidney Biopsy*
*Microscopy Images*

**Native Byte Encoding**
*Image formats*

**DAP**

**Data Structures**
*Image*

**DTS**

**Derived Data/Metadata**
*Tissue classification, Changes in tissue*

**Applications**
*Population Obesity, Renal Failure*

**Usable Data**

*Amelia Bartholomew, Medicine, University of Illinois Chicago*

**Data Collection**
*Publications*

**Native Byte Encoding**
*Document formats*

**DAP**

**Data Structures**
*Tables, Images*

**DTS**

**Derived Data/metadata**
*Demography data, gene/manifestation correlations*

**Applications**
*Population Obesity, Renal Failure*

**Usable Data**

# Prevalence of Renal Insufficiency in Individuals with Hypertension and Obesity/Overweight: The FATH Study

**Table 1.**

Characteristics of the patients[a]

| Variable | Overweight (BMI 25 to 29.9 kg/m$^2$)($n = 2060$) | Obesity (BMI ≥30 kg/m$^2$)($n = 2525$) | P |
|---|---|---|---|
| Age (yr) | 61.9 (10.5) | 61.9 (10.7) | NS |
| Male (%) | 51.8 | 45.0 | <0.0001 |
| BMI (kg/m$^2$; mean [SD] | 27.8 (1.3) | 35.1 (4.1) | <0.0001 |
| Waist (cm; mean [SD]) | | | |
| male | 101.1 (10.5) | 113.6 (11.5) | <0.0001 |
| female | 94.2 (10.7) | 107.9 (12.9) | <0.0001 |
| SBP (mmHg; mean [SD] | 145.75 (17.4) | 145.84 (18.2) | NS |
| DBP (mmHg; mean [SD] | 85.01 (10.3) | 85.5 (10.8) | NS |
| Glucose (mg/dl; mean [SD] | 110.0 (28.9) | 117.7 (34.1) | <0.0001 |
| HDL cholesterol (mg/dl; mean [SD] | 53.6 (15) | 51.3 (13.2) | <0.0001 |
| Triglycerides (mg/dl; mean [SD] | 148.0 (68) | 161.7 (78) | <0.0001 |
| Diabetes | 26.04 (24.1 to 27.9) | 37.03 (10.9 to 13.5) | <0.0001 |
| MS 1 (% [95% CI]) | 80.2 (78.0 to 82.2) | 92.8 (91.5 to 94.0) | <0.0001 |
| MS 2 (% [95% CI]) | 85.4 (83.4 to 87.2) | 95.1 (94.0 to 96.0) | <0.0001 |

- ↵[a] BMI, body mass index; CI, confidence interval; DBP, diastolic BP; MS 1, metabolic syndrome (Adult Treatment Panel III criteria); MS 2, metabolic syndrome (International Diabetes Federation criteria); SBP, systolic BP.

*Ian Brooks, Biochemistry, University of Illinois Urbana-Champaign*

**Data Collection**
*Community Data*

**Native Byte Encoding**
*Documents, …*

**DAP**

**Data Structures**
*Numerical, Temporal Feature Vectors, Classifications*

**DTS**

**Derived Data/ Metadata**
*Correlations in data*

**Applications**
*Disease Spread/Containment*

**Usable Data**

As of: 12/3/2012 10:26:08 AM

| | BIRD | CAT | DOG | OTHER | Total |
|---|---|---|---|---|---|
| ALOPECIA | 3 | 15 | 11 | 3 | 32 |
| COCCIDIA | | | | | |
| COUGH | | | | | |
| DEMODEX | | | | | |
| DERMATITIS | | | | | |
| DIARRHEA | | | | | |
| EAR INFECTION | | | | | |
| EAR MITES | | | | | |
| FIV | | | | | |
| FLEA DERMATITIS | | | | | |
| FRACTURE | | | | | |
| GIARDIA | | | | | |
| HEARTWORM | | | | | |
| KENNEL COUGH | | | | | |
| PARVO | | | | | |
| RINGWORM | | | | | |
| ROUNDWORMS | | | | | |
| TAPEWORMS | | | | | |
| VOMITING | | | | | |
| WHIPWORMS | | | | | |
| Total | | | | | |

BIRD
A070187
A070187
A070375
CAT
A069732
A064727
A070012
A070042
A070043
A056111
A068864
A069858
A070011
A069106

| Codes: E | Elementary A | Elementary B | Elementary C | High School B | High School | High School C |
|---|---|---|---|---|---|---|
| 07/24/2008 | 2 | | | | | |
| 07/25/2008 | 5 | | | | | |
| 07/28/2008 | 3 | | | | | |
| 07/29/2008 | 3 | | | | | |
| 07/30/2008 | 4 | | | | | |
| 07/31/2008 | 12 | | | | | |
| 08/01/2008 | 18 | | | | | |
| 08/04/2008 | 8 | | | | | |
| 08/05/2008 | 13 | | | | | |
| 08/06/2008 | 9 | | | | | |
| 08/07/2008 | 9 | | | | | |
| 08/08/2008 | 14 | | | | | |
| 08/11/2008 | 20 | | | | | |
| 08/12/2008 | 17 | | | | | |
| 08/13/2008 | 9 | | | | | |
| 08/14/2008 | 10 | | | | | |
| 08/15/2008 | 24 | | | | | |
| 08/18/2008 | 9 | | | | | |
| 08/19/2008 | 10 | | | | | |
| 08/20/2008 | 6 | | | | | |
| 08/21/2008 | 8 | 1 | 6 | 14 | 3 | 1 |
| 08/22/2008 | 17 | 2 | 11 | 44 | 37 | |
| 08/25/2008 | 10 | 5 | 11 | 42 | 36 | |
| 08/26/2008 | 11 | 10 | 13 | 35 | 43 | |
| 08/27/2008 | 8 | 12 | 4 | 60 | 47 | 1 |
| 08/28/2008 | 9 | 10 | 15 | 57 | 56 | 2 |

DIARRHEA                    09/19/2012                    10/25/2012                    BOMBAY/MIX

**Christopher Lynberg, Computer Science, Center for Disease Control (CDC)**

**Data Collection**
*Pathogen Cards*

**Native Byte Encoding**
*Image formats*

**DAP**

**Data Structures**
*Image*

**DTS**

**Derived Data/metadata**
*Pathogen attributes from blood work*

**Applications**
*Disease*

**Usable Data**

**GENUS** _Paenibacillus_   **SPECIES** _species_   **CARD NO.** 1

| Cult No. | Glu | Xyl | Man | Lact | Suc | Malt | Bk | Aes | Pen | Pig | O/F | 10%L | TSI S/BH₂S | Hem | Cat | O | AGT / SS | Cit/Cet | Urea | Nit | Ind | MR/VP | Gel | Milk | NaCl Tol. | LAO/PPA | S/A | Temp | Mot | Notes & Source | Sender |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 4E | A | A | A | A | A | A | – | + | | – | | | A/A | 2+ | LS/Steb | 2+ | 0,C Sw | – | – | +? | – | – | +M | | OL 6N | – | F 92 | + | Peri + | | |
| 4EA | g? | – | A | – | A | + | | – | | | A/A | 4+ 2+ Gv | | – | +/0,L Sw | – | – | +? | – | OL | | OL 6– | –/– | A? 16 | + | + | L. foot | OH |
| 3E | A | – | – | A | – | A | – | | – | | | a/iv | 2+ | – | 2+ | +/0,C Sw | – | – | + | – | + | 14 iv | OL L– | +/– | A? 4 | + | + | blood | MS |
| 3E | A | A | A | A | A | A | + | | – | | | A/A | 1+ | – | 2+ | –/0,T Sw | – | – | +/– | – | + | 7 iv | A/iv 0M 6– | – | | +/– | + | blood | GA |
| 2E | g? | A? | g? | A? | A | A? | + | | (Pink) | | | N/N | – | – | + | +/◯ | – | – | – | – | – | 7 | – | OL 6M | –/– | M 43 | –/+ | Peri + | blood | FL |
| E | A | A?/a Na | A?/a | A | A | A?/a | + | | – | | | A/A | +/– | | +/wt Sw | – | – | + | – | + | +IR | L | | F | –/+ | +? | peri blood | Colo. |
| 5E | A | A | A | A | A | A | + | | – | | | A/A | 4 | ly | 2 | 7/0T Sw | – | – | + | – | + | 14 IR(A) | m | – | F 80 | 3/4 | peri + | sev | OH |
| 2E | A | A | A | A | A | A | + | | – | | | A/A | 4 | ly | 2 | 7/0T Sw | – | – | + | – | + | 14 IR(A) | m | – | A? 28 | 3/4 | peri + | blood | NH |
| 5E | A | A | A | A | A | A | + | | – | | | A/A | 4 | ly | 1 | –/0T Sw | – | – | +/– | – | + | 14 IR(A) | m | – | N 73 | 3/4 | peri + | blood | ARK |
| 10E | A | – | – | – | – | + | | – | | | C/C | 3 | – | 2+ | 0 ST/Sw | – | – | + | | – | IR | m | M | ±/3 | peri + | blood | TN |
| 0F | A | A | A | A | A | A | + | | – | | | a/N | 1 | | 1 | +/0 St Sw | – | – | – | | – | IR | L | M 69 | 1/2 | peri + | blood | MS |
| F | A | A | – | A | – | A | + | | – | | | A/A | 2 | ly | 1 | +/0c Sw | – | – | + | | – | IR | m | 84 | ±/2 | peri + | blood | OL |
| 0F | A | A?/a2 | – | A | A | A | + | | – | | | A/A | 3 | | 2 | –/0 free | – | – | – | | – | ADR/Clot | L? | A? 10 | m/m | peri + | Rt Palm wound | ID |
| 0F | A | A | A | A | – | A | + | | – | | | A/A | 1 | – | 1 | +/◯ | – | – | + | | – | IR | m | M 83 | ±/– | peri + | blood | NY |
| 3 0F | A? | A? | A? | A? | A? | A? | K | + | | – | | | A/A? | 3+ | – | 2 | –/0 St Sw | – | – | – | | – | IR | m | | 2/3 | – | blood | TN |
| 3E | A | A | – | A | – | A | + | | – | | | A/A | 3+ | – | – | –/◯ | – | – | +/+ | – | – | IR/A | m | F 76 | –/– | + | blood | OK |
| 1E | A | A | – | A | A? | A | + | | – | | | P/A | 3+ | – | – | –/◯ | – | – | +/– | – | – | IR/A | m | F 76 | ±/– | + | blood | OK |

52.81 (f. 3.365) 12-91   **GENUS CROSS FILE**   (See Reverse)

**Sonia Giovinazzi, Civil & Natural Resource Engineering, University of Canterbury**

**Data Collection**
*Sewer/Water System Videos*

**Native Byte Encoding**
*Video formats*

DAP

**Data Structures**
*Image, Image Sequences*

DTS

**Derived Data/ Metadata**
*Structural defects*

**Applications**
*Urban Infrastructure*

**Usable Data**

*Chris German, Marine Geochemistry, Woods Hole Oceanographic Institute (WHOI)*
*Scott Gallager, Biology, Woods Hole Oceanographic Institute (WHOI)*
*James Kinsey, Mechanical Engineering, Woods Hole Oceanographic Institute (WHOI)*
*Joe Futrelle, Computer Science, Woods Hole Oceanographic Institute (WHOI)*

**Data Collection**
*UAV Data*

**Native Byte Encoding**
*Image Formats, Video Formats*

**DAP**

**Data Structures**
*Image*

**DTS**

**Derived Data/ Metadata**
*Color Correction, Event detection, Species classification/counting*

**Applications**
*Marine Biology, Fishing Industry*

**Usable Data**

*David Zeppa, Computer Science, University of California Santa Cruz*

**Data Collection**
*Publications*

**Native Byte Encoding**
*Document formats*

**DAP**

**Data Structures**
*Image*

**DTS**

**Derived Data/metadata**
*Polymers*

**Applications**

**Usable Data**

$$CH_3-O\left(CH_2CH_2-O\right)_n H$$

$$\left[CH_2-CH\right]_n$$

with phenyl group

*Robert Markley, English, University of Illinois Urbana-Champaign*

**Data Collection**
*Historical Maps*

**Native Byte Encoding**
*Image formats*

**DAP**

**Data Structures**
*Image*

**DTS**

**Derived Data/ Metadata**
*Lake locations, boundaries, signatures for CBIR*

**Applications**
Coastline/climate changes over time

**Usable Data**

*Ann MacNeil, Music, University of North Carolina at Chapel Hill*

**Data Collection**
*Sheet Music*

**Native Byte Encoding**
*Image formats*

**DAP**

**Data Structures**
*Image*

**DTS**

**Derived Data/ Metadata**
*Notes*

**Applications**

**Usable Data**

# NEW! DRAS-TIC

## Digital Repository At Scale - That Invites Computation
## [ To Improve Collections ]

**GOAL:** Build out the open source DRAS-TIC platform into a horizontally scalable archives

framework **serving the national library, archives, and scientific data management communities**

- **Product** of a 2-year startup by partners, Archive Analytics Solutions Ltd.
- **Scaling** to billions of files and beyond
- **Interfaces:**
    - Web client
    - Command-line client
    - REST storage API (CDMI) industry standard
- **Key-value** metadata
- **Listener** mechanism
- **Python** source on GitHub (Open AGPL license)
- **Apache Cassandra** database (CERN, eBay, GitHUB, Hulu, Instagram, Netflix, Twitter…)

- **Computational Finding Aids**

# The Rest of the Talk..

- Approaching 1 Billion files

- New DRAS-TIC Repository

- NCSA's Brown Dog Service

- Automatic Feature Extraction & Curation

- Digging into Collections with Elasticsearch

- Projects & Opportunities

**DRAS·TIC**

# RG 029 - Records of the Bureau of Census   Edit  Delete

Search ...                                Go!

Archive
Users
Groups
Activity

Add new collection    Add new item

⊗ 📂 2006 Census Operational Photos

⊗ 📂 A Profile Of Older Workers In West Virginia

⊗ 📂 acs

⊗ 📂 acs2002

⊗ 📂 acs2003

⊗ 📂 acs2004

# Workflow for a Digital Object



**PDF**

File Name
Directory
File Size

elasticsearch

**REPOSITORY**

**SERVICES**

# Text Format Conversion (PDF to TXT)

**PDF**

**TXT**

File Name
Directory
File Size

elasticsearch

**REPOSITORY**

**SERVICES**

# Now we have a full text index..

**TXT**

**PDF**

**Full Text**
File Name
Directory
File Size

elasticsearch

**REPOSITORY**

**SERVICES**

# Optical Character Recognition (OCR) Extractor



**PDF**

**PNG**

**OCR**

**OCR Text**
File Name
Directory
File Size

elasticsearch

**REPOSITORY**

**SERVICES**

# Format Recognition (Siegfried PRONOM Extractor)

PDF

PNG

OCR

PUID
PDF
1.4.2

**Format**
OCR Text
File Name
Directory
File Size

elasticsearch

**REPOSITORY**

**SERVICES**

# Facial Recognition (Computer Vision Extractors)



**PDF**

**PNG**

**OCR**

**6 FACES**

**PUID PDF 1.4.2**

elasticsearch

**# Faces**
Format
OCR Text
File Name
Directory
File Size

**REPOSITORY**

**SERVICES**

# Facial Recognition (Computer Vision Extractor)



**REPOSITORY**

**SERVICES**

PDF

PNG

OCR

6 Faces

12 Eyes

1 Close Up

3 in Profile

PUID PDF 1.4.2

# Faces
# Eyes
# Close Ups
# Profiles
Format
OCR Text
File Name
Directory
File Size

elasticsearch

# PDF Object Enhanced with Extracted Metadata

# DON'T PANIC

# ElasticSearch + Kibana

- Free plugin for Elasticsearch

- Gives shape to an Elasticsearch index

- Write queries visually and interactively

Lots of ways to explore the data

Files Formats
○ Concentric Pie Chart
○ Inner: Mimetype
○ Outer: PRONOM PUID

| field | value | Count |
| --- | --- | --- |
| mimetype | application/pdf | 1,208 (49.31%) |
| puid | info:pronom/fmt/18 | 411 (34.11%) |

# Charts can be added to data dashboards..

# Arrangement can be used as a Facet

As you browse the hierarchy...

The entire dashboard is redrawn to reflect the particular record group, series or folder under study.
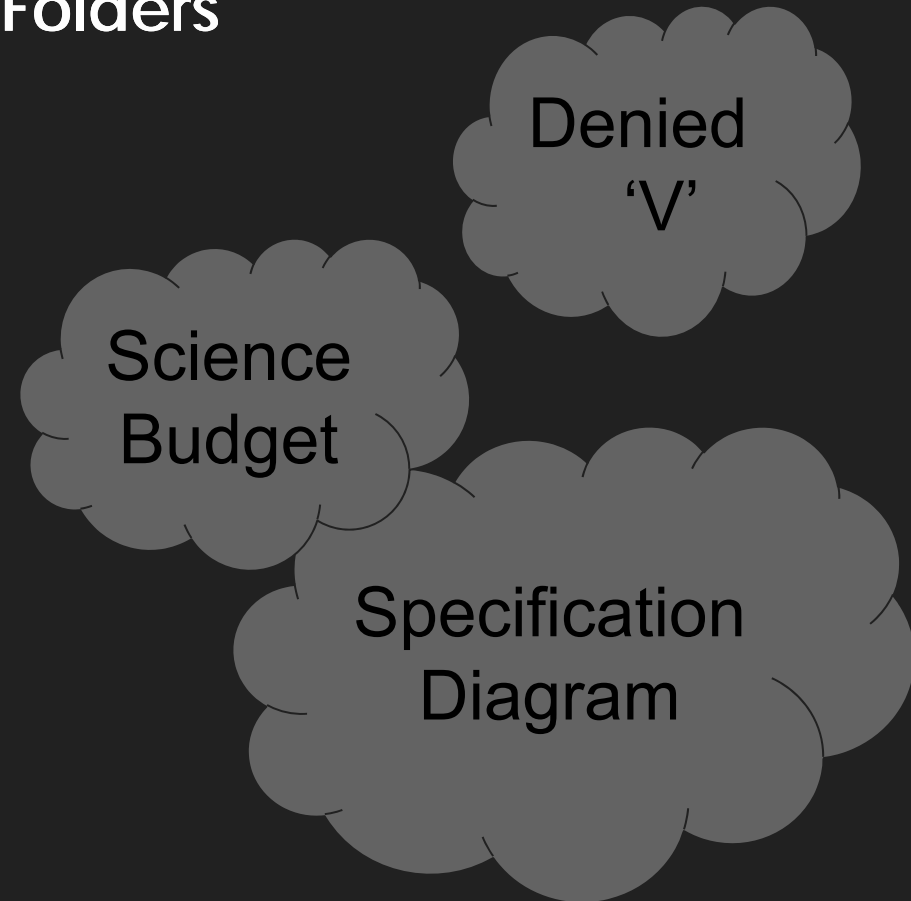
"Drill down" or zoom in and out of your collections.

# Text Comparison between Folders

**Significant Terms** are based on full text.

They are significant within overall scope of query.

**Significant Terms** can be used to distinguish neighboring folders or documents.

Denied 'V'

Science Budget

Specification Diagram

| parentURI: Descending ⇳ Q | Top 2 unusual terms in fulltext ⇳ Q | Count ⇳ |
|---|---|---|
| /Archive/ciber/RG 267 - Records of the Supreme Court of the United States/Orders and Journals/www.supremecourtus.gov/orders/courtorders/ | denied | 606 |
| /Archive/ciber/RG 267 - Records of the Supreme Court of the United States/Orders and Journals/www.supremecourtus.gov/orders/courtorders/ | v | 670 |
| /Archive/ciber/RG 359 - Records of the Office of Science and Technology/Office of Science and Technology Website/www.ostp.gov/pdf/ | science | 305 |
| /Archive/ciber/RG 359 - Records of the Office of Science and Technology/Office of Science and Technology Website/www.ostp.gov/pdf/ | budget | 288 |
| /Archive/ciber/RG 167 - Records of the National Institute of Standards and Technology/Visualization of Structural Steel Product Models, Construction Sites and Equipment, and the Virtual Cybernetic Building Testbed/cic.nist.gov/vrml/cis/lpm6/structural_frame_schema/lexical/ | specification | 314 |
| /Archive/ciber/RG 167 - Records of the National Institute of Standards and Technology/Visualization of Structural Steel Product Models, Construction Sites and Equipment, and the Virtual Cybernetic Building Testbed/cic.nist.gov/vrml/cis/lpm6/structural_frame_schema/lexical/ | diagram | 284 |

# DRAS-TIC
Institutional R&D Partners
Use cases for Parallel Compute
Fedora Sprinters

# Brown Dog

Try it on your Scientific Data
Become an Early Adopter of the API
Contribute Extractors & Converters

# UMD iSchool
Partner with the DCIC on Projects
Digital Curation Certificate Program
Computational Archival Science

# JOIN FORCES

# Thank you ISGC 2017 Team:
## Ludek, Simon, Stella, Vicky,
### *and many other staff*