# Two Projects, One Challenge:
# Common research data issues in MIREX and HTRC

Presented by
J. Stephen Downie
University of Illinois at Urbana-Champaign

# Acknowledgements

- Most of today's slides are directly drawn (aka copied) from the slides presented at the HTRC UnCamp events in Bloomington, Indiana and Champaign, Illinois.

- Today's talk summarizes four days of excellent presentations and demonstrations!

- We thank the HTRC team and the UnCamp presenters for the use of their very informative slides.

# Agenda

- Introduce MIREX

- Introducing the HathiTrust

- Non-Consumptive Research
    - MIREX and NEMA
    - HTRC Software Architecture

- Workset Creation for Scholarly Analysis

- Next Steps

# In the Beginning…

- Work began on MIREX in Bloomington, Indiana in 2001
  - The "Indiana Manifesto"
- 2001-2003
  - Fact-finding meetings, planning meetings, and workshops funded by Mellon and NSF
  - Large-scale funding from NSF and Mellon 2003
- *Audio Description Contest* at ISMIR 2004
- MIREX first run at ISMIR 2005 in London
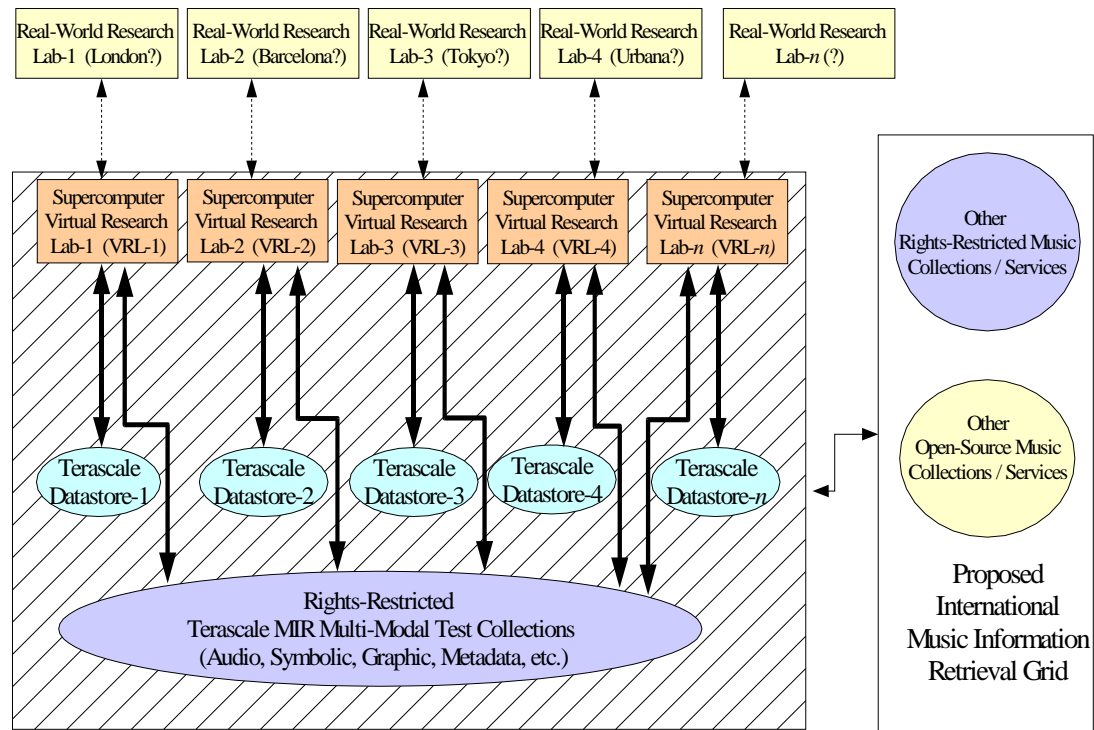
# MIREX Model

- Based upon the TREC approach:
  - Standardized queries/tasks
  - Standardized collections
  - Standardized evaluations of results
- Not like TREC with regard to distributing data collections to participants
  - Music copyright issues, ground-truth issues, overfitting issues

# IMIRSEL:  First Principles

1.  Security for the music materials
2.  Accessibility for international, domestic and internal researchers
3.  Sufficient computing and storage infrastructure for the computationally- and data-intensive MIR/MDL techniques examined

# IMIRSEL Model



Real-World Research Lab-1 (London?)
Real-World Research Lab-2 (Barcelona?)
Real-World Research Lab-3 (Tokyo?)
Real-World Research Lab-4 (Urbana?)
Real-World Research Lab-$n$ (?)

Supercomputer Virtual Research Lab-1 (VRL-1)
Supercomputer Virtual Research Lab-2 (VRL-2)
Supercomputer Virtual Research Lab-3 (VRL-3)
Supercomputer Virtual Research Lab-4 (VRL-4)
Supercomputer Virtual Research Lab-$n$ (VRL-$n$)

Other Rights-Restricted Music Collections / Services

Terascale Datastore-1
Terascale Datastore-2
Terascale Datastore-3
Terascale Datastore-4
Terascale Datastore-$n$

Other Open-Source Music Collections / Services

Rights-Restricted Terascale MIR Multi-Modal Test Collections (Audio, Symbolic, Graphic, Metadata, etc.)

Proposed International Music Information Retrieval Grid

Legend:

Super-Bandwith I/O Channel

NCSA *Music Data* **Secure** Zone

Command/Control/Derived Data traffic via Internet

Connection to International MIR Grid

# MIREX Overview

- Began as MIREX in 2005
- Tasks defined by community debate
- Data sets collected and/or donated
- Participants submit code to IMIRSEL
- Code rarely works first try ☺
- Huge labour consumption getting programmes to work
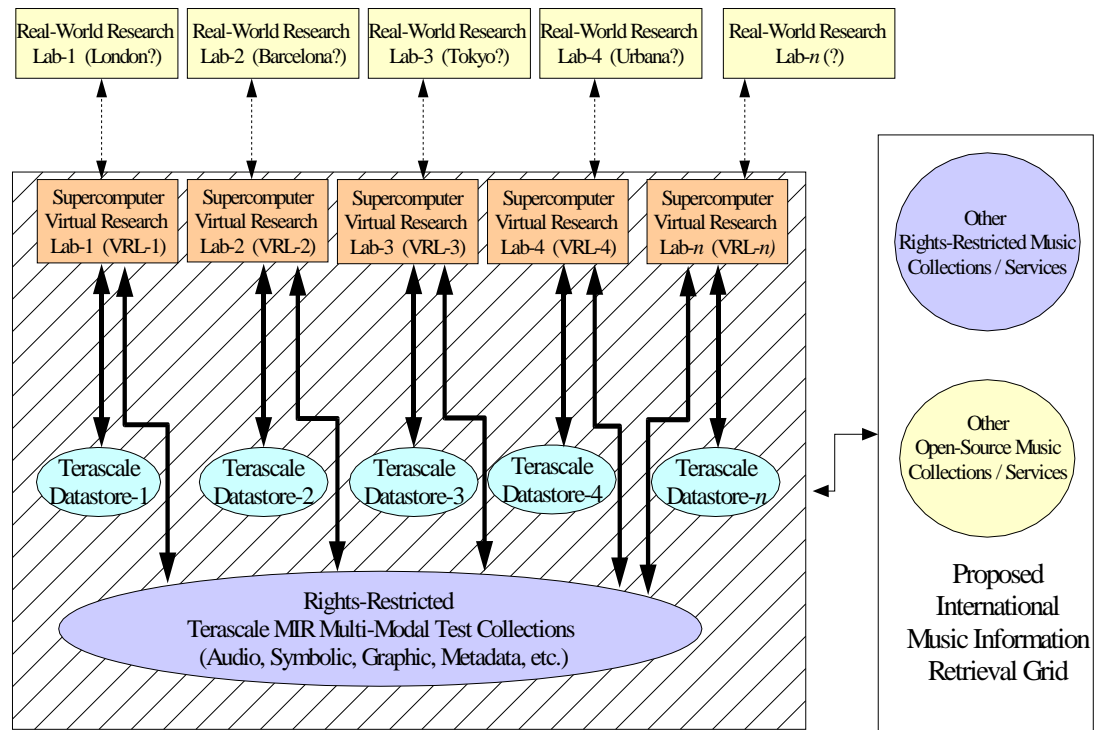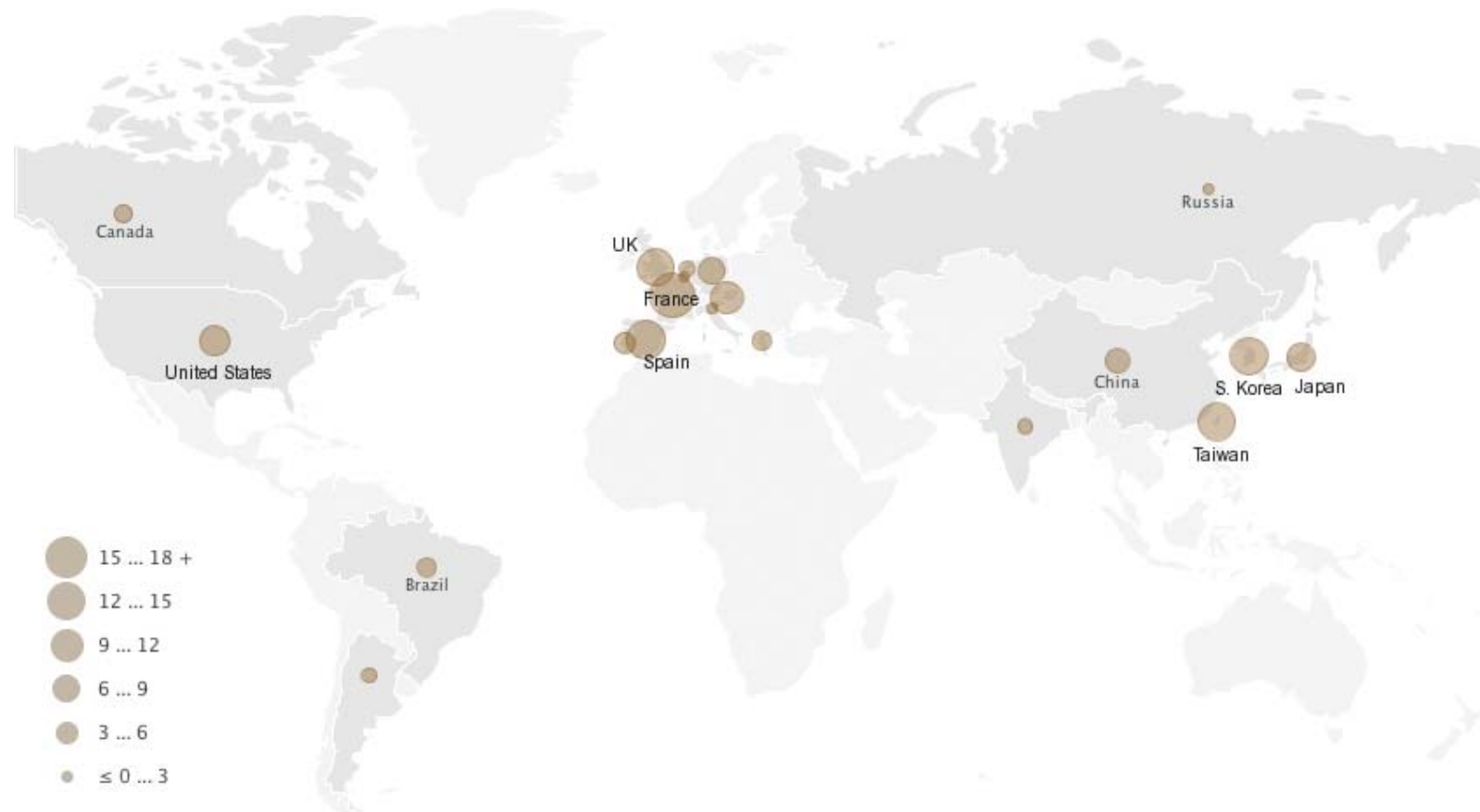- Meet at ISMIR to discuss results

# MIREX Model

- Based upon the TREC approach:
  - Standardized queries/tasks
  - Standardized collections
  - Standardized evaluations of results
- Not like TREC with regard to distributing data collections to participants
  - Music copyright issues, ground-truth issues, overfitting issues

# IMIRSEL:  First Principles

1. Security for the music materials
2. Accessibility for international, domestic and internal researchers
3. Sufficient computing and storage infrastructure for the computationally- and data-intensive MIR/MDL techniques examined

# IMIRSEL Model

# MIREX Participation



MIREX 2012: 109 participants from 20 countries

# MIREX Participation

| | 2005 | 2006 | 2007 | 2008 | 2009 | 2010 | 2011 | 2012 | 2013 |
|---|---|---|---|---|---|---|---|---|---|
| **Datasets** | 10 | 13 | 12 | 18 | 26 | 31 | 32 | 35 | 37 |
| **Individuals** | 82 | 50 | 73 | 84 | 138 | 152 | 156 | 109 | 116 |
| **Runs** | 86 | 92 | 122 | 169 | 289 | 337 | 312 | 302 | 328 |

# MIREX 2013

- 116 researchers
- More than 29 countries
- 37 datasets
- 24 tasks
- 328 completed runs

# MIREX 2013 TASKS

| | |
|---|---|
| Audio Artist Identification | Audio Onset Detection |
| Audio Beat Tracking | Audio Tag Classification |
| Audio Chord Detection | Audio Tempo Extraction |
| Audio Classical Composer ID | Multiple F0 Estimation |
| Audio Cover Song Identification | Multiple F0 Note Detection |
| Audio Drum Detection | Query-by-Singing/Humming |
| Audio Genre Classification | Query-by-Tapping |
| Audio Key Finding | Score Following |
| Audio Melody Extraction | Structural Segmentation |
| Audio Mood Classification | Symbolic Genre Classification |
| Audio Music Similarity | Symbolic Key Finding |
| Discovery of Repeated Themes & Sections | Symbolic Melodic Similarity |

# MIREX Participation

| | 2005 | 2006 | 2007 | 2008 | 2009 | 2010 | 2011 | 2012 | 2013 |
|---|---|---|---|---|---|---|---|---|---|
| **Datasets** | 10 | 13 | 12 | 18 | 26 | 31 | 32 | 35 | 37 |
| **Individuals** | 82 | 50 | 73 | 84 | 138 | 152 | 156 | 109 | 116 |
| **Runs** | 86 | 92 | 122 | 169 | 289 | 337 | 312 | 302 | 328 |

Total Runs: 2037!

# Introducing the HathiTrust

# Partnership

Allegheny College
Arizona State University
Baylor University
Boston College
Boston University
Brown University
California Digital Library
Colby College
Columbia University
Cornell University
Dartmouth College
Duke University
Emory University
Florida State University
Getty Research Institute
Harvard University Library
Indiana University
Johns Hopkins University
Lafayette College
Library of Congress
Massachusetts Institute of
    Technology
McGill University`
Michigan State University
New York Public Library
New York University
North Carolina Central
    University

North Carolina State
    University
Northwestern University
The Ohio State University
The Pennsylvania State
    University
Princeton University
Purdue University
Stanford University
Temple University
Texas A&M University
Tufts University
Universidad Complutense
    de Madrid
University of Alberta
University of British Columbia
University of Arizona
University of Calgary
University of California
    Berkeley
    Davis
    Irvine
    Los Angeles
    Merced
    Riverside
    San Diego
    San Francisco
    Santa Barbara
    Santa Cruz
The University of Chicago
University of Connecticut

University of Delaware
University of Florida
University of Houston
University of Illinois
University of Illinois at Chicago
The University of Iowa
University of Maryland
University of Massachusetts
University of Miami
University of Michigan
University of Minnesota
University of Missouri
University of Nebraska-Lincoln
The University of North
        Carolina at Chapel Hill
University of Notre Dame
University of Oklahoma
University of Pennsylvania
University of Pittsburgh
University of Queensland
University of Tennessee,Knoxvile
University of Utah
University of Virginia
University of Washington
University of Wisconsin-Madison
Utah State University
Wake Forest University
Washington University
Yale University Library

# Mission

To contribute to the common good by collecting, organizing, preserving, communicating, and sharing the record of human knowledge

# HathiTrust "Wow" Numbers

- 10,924,244 total volumes
- 5,719,252 book titles
- 285,776 serial titles
- 3,823,485,400 pages
- 490 terabytes
- 129 miles
- 8,876 tons
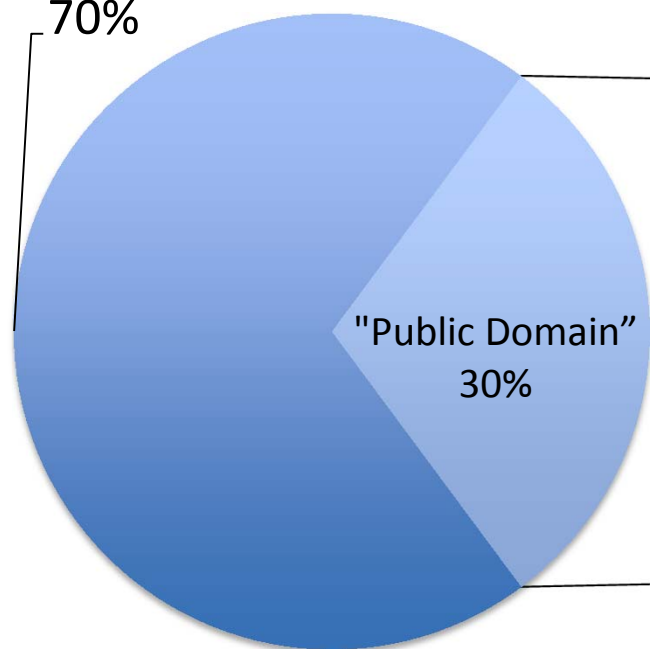- 3,565,657 volumes(~33% of total) in the public domain

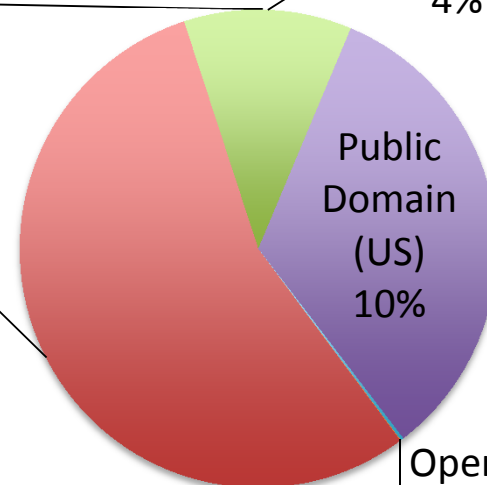# Content Sources

# Content Distribution



In-copyright or undetermined 70%

"Public Domain" 30%

U.S. Federal Government Documents (worldwide) 4%

Public Domain (US) 10%

Public Domain (worldwide) 15%

Open Access .1%

Creative Commons .01%

# Google PD Research Collection

- Public Domain Materials of the HatihTrust
  - 2,592,097 Volumes
  - Gigabytes
    - 2.3 TB in raw OCR'd text
    - 3.7 TB of managed OCR'd text
    - 1.85 TB solr Index
  - Monthly Updates
    - And irregular data 'take down' requests



■ Total volumes
■ Public Domain volumes

# Non-Consumptive
# Research Model

# Non-Consumptive Research Paradigm

- *No action or set of actions on part of users, either acting alone or in cooperation with other users over duration of one or multiple sessions can result in sufficient information gathered from collection of copyrighted works to reassemble pages from collection.*

- Definition disallows collusion between users, or accumulation of material over time. Differentiates human researcher from proxy which is not a user. Users are human beings.

# Non-Consumptive Research Paradigm

Bring the

COMPUTATION

to the

DATA!

Researcher 1    Researcher 2    Researcher 3

Wikis    Web Front Ends    NEMA Portal    Code Repositories    Mail Archives

# NEMA'S SEASR Framework

**High Level Services**

MIREX DIY

OMRAS2 TWM

jMIR

myExperiment

**Low Level Services**

Results Aggregators

Classification Modules
(ACE, M2K, OMRAS2)

Data Cleaning Tools
(MusicMetadataManager)

Data Exchange Tools
(ACE XML, M2K, RDF)

**Discovery & Sharing Services**

Greenstone

Maestro

Music Ontology

## OMEN

Web service calls (Only features returned)

Grid-based feature extraction tools
(M2K, jAudio, OMRAS@Home, etc.)

Grid-based feature extraction tools
(M2K, jAudio, OMRAS@Home, etc.)

Grid-based feature extraction tools
(M2K, jAudio, OMRAS@Home, etc.)

(Data passed internally)

Music database 1    Music database 2    Music database 3

# An External Classification Algorithm

# Networked Environment for Music Analysis
# NEMA DIY Interface

# Evaluation Reports

MIREX 2010: Audio Chord Description - MIREX09 Dataset

| Introduction | Summary | Detailed Evaluation Metrics | RRHS1 | RRHS2 | PVM1 | PP1 | MD1 | KO1 | EW1 | MK1 | OFG1 | EW4 |
| EW2 | EW3 | UUOS1 | MM1 | CWB1 | Comparative plots | Significance Tests | Raw data files |

## Introduction

### Description

| Field | Value |
|---|---|
| Task ID | 17 |
| Task Name | MIREX 2010: Audio Chord Description - MIREX09 Dataset |
| Task Description | Chord transcription task requiring participants to annotate and segment the chord events in the MIREX09chord transcription dataset. Please note that: <ul><li>Evaluations are performed at the triad level,</li><li>results for both pretrained algorithms and algorithms trained and tested under 3 fold cross-validation are reported here.</li><li>pretrained algorithms are likely to have been trained on the evaluation dataset hence they are expected to achieve higher results than algorithms evaluated on held out data.</li></ul> |
| Subject Metadata ID | 26 |
| Subject Metadata Name | Chord label sequence |
| Dataset ID | 33 |
| Dataset Name | MIREX09 Chord |
| Dataset Description | MIREX 2009 Chord transcription dataset composed of Christopher Harte's Beatles dataset (C4DM, Queen Mary's University of London) and Matthias Mauch's Queen and Zweieck dataset (C4DM, Queen Mary's University of London) |
| Date report generated | Aug 6, 2010 7:56:08 PM |

### Legend    [top]

| Submission code | Submission name | Abstract PDF | Contributors |
|---|---|---|---|
| CWB1 | ChordID | PDF | Taemin Cho, Ron Weiss, Juan Bello |
| EW1 | LabROSA Chord Train/Test 2010 | PDF | Daniel Ellis, Adrian Weller |

33

# HTRC Architecture

**Portal Access**

Blacklight

**Agent**

Job Submission | Collection building

Direct programmatic access (by programs running on HTRC machines)

**Security (OAuth2)**

Data API access interface | Solr Proxy

**Registry (WSO2)**

Algorithms | Meandre Workflows

Result Sets | Collections

**Audit**

Cassandra cluster volume store

Solr index

Compute resources

Storage resources

# HTRC Architecture

## Portal Access

Blacklight

## Agent

Job Submission

Collection building

S

## Registry (WSO2)

Algorithms

Meandr
Workflow

Result Sets

Collectio

Compute resources

## Portal Access

### HTRC Portal

Blacklight

## App SEAR

## App Blacklight



Flow Parameter · Universal Text Extractor · Search Text · HTRC Page Retriever · Text Cleaner · Sentence Detector · Sentence Tokenizer · Named Entity · Add Tuple Attribute · Add Tuple Attribute · Tuple Aggregator · Tuple To HTML · Stream Delimiter Filter · Push Text · Write To File

blacklight

# HTRC Architecture

## Portal Access

Blacklight

## Agent

| Job Submission | Collection building |

S

## Registry (WSO2)

| Algorithms | Meandre Workflows |
| Result Sets | Collections |

Cassandra cluster volume store

Solr index

Compute resources

Storage resources

## Agent

### HTRC Agent

| Job Submission | Collection building |

# HTRC Architecture

## Portal Access

Blacklight

## Agent

| Job Submission | Collection building |

## Registry (WSO2)

| Algorithms | Meandre Workflows |
| Result Sets | Collections |

cluster volume store

Solr index

Compute resources

Storage resources

# HTRC Registry

## Registry (WSO2)

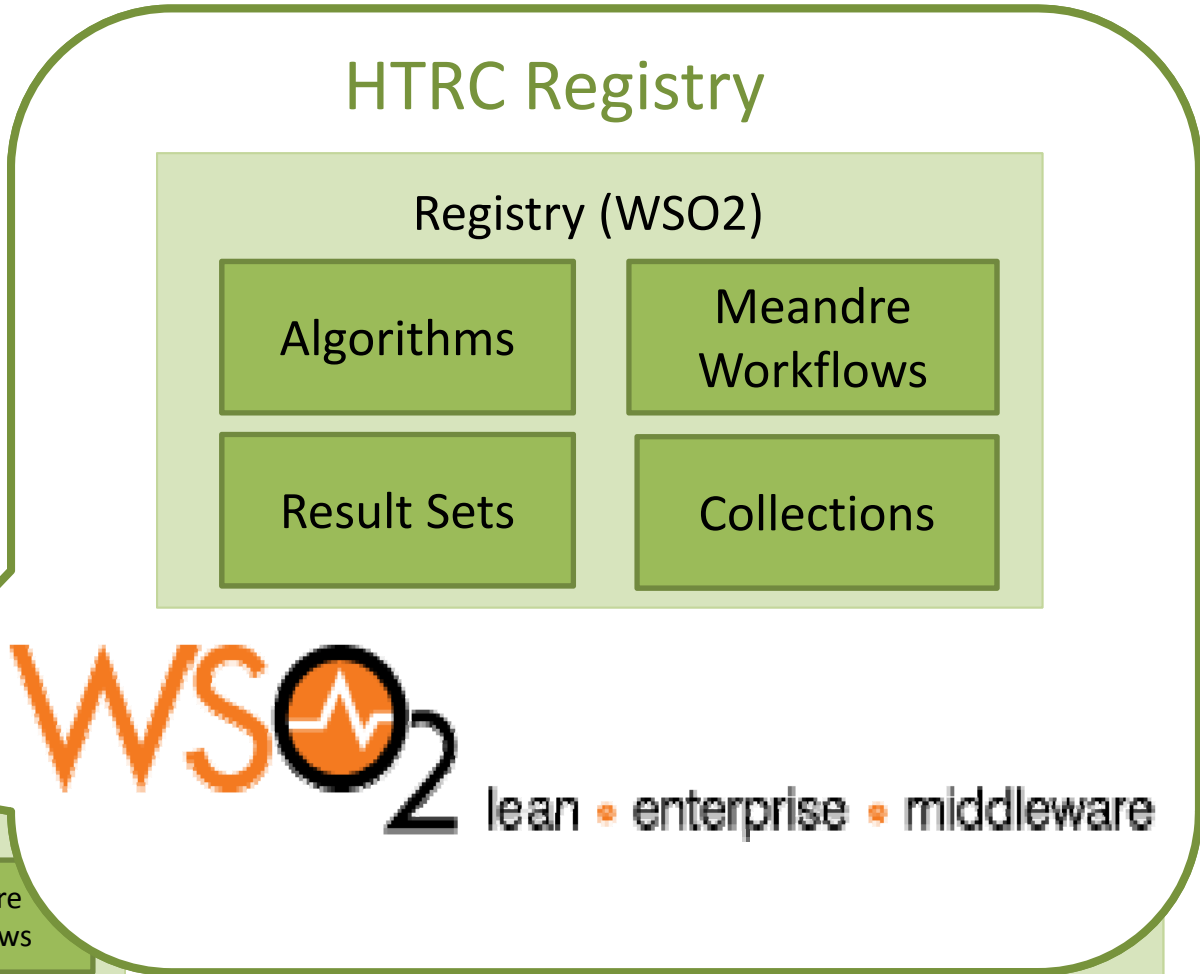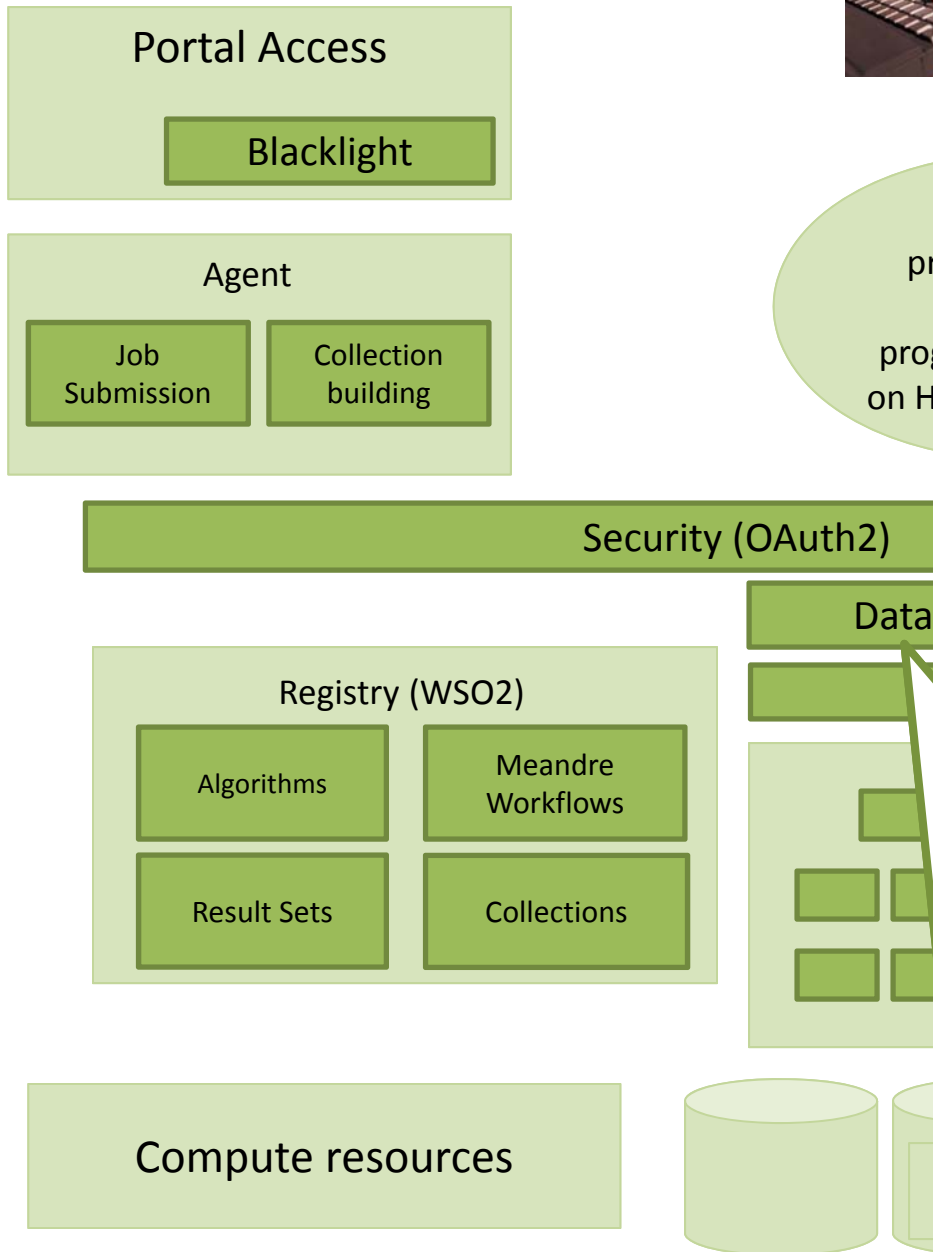| Algorithms | Meandre Workflows |
| Result Sets | Collections |

WSO2 lean • enterprise • middleware

# HTRC Architecture

**Portal Access**

Blacklight

**Agent**

| Job Submission | Collection building |

**Security (OAuth2)**

Data

**Registry (WSO2)**

| Algorithms | Meandre Workflows |
| Result Sets | Collections |

Compute resources

pr...
a...
prog...
on H...

# Secure Data API

- RESTful Web Service
  - Language agnostic
  - Clients don't have to deal with Cassandra
- Simple OAuth2 authentication
- HTTP over SSL
- Audits client access
- Protected behind firewall, accessible only to authorized IPs

HTRC

H7...

## Solr Proxy

Solr proxy

Solr service

RFS distributed file system

Direct
...ogrammatic
...access (by
...rams running
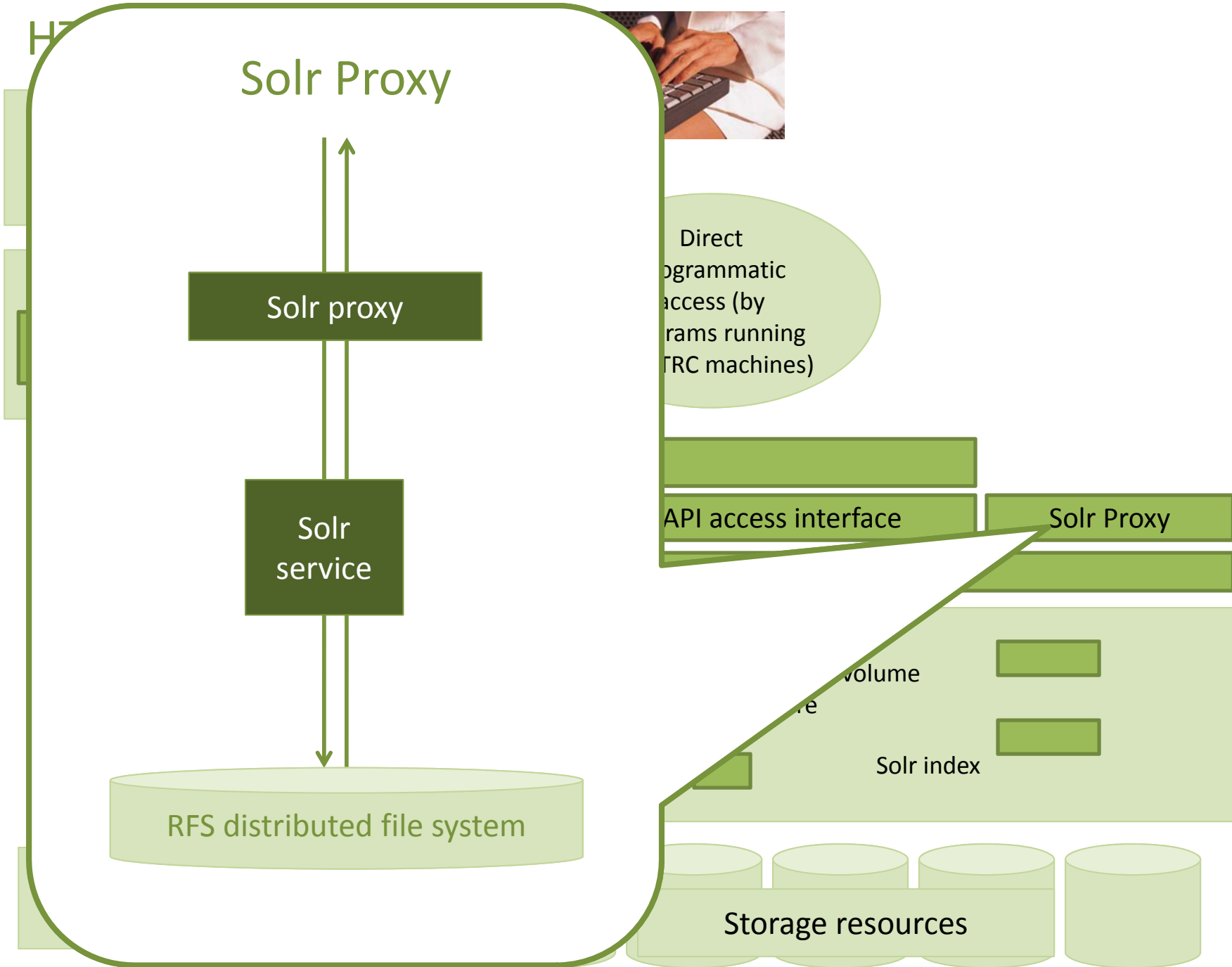...TRC machines)

API access interface

Solr Proxy

...volume
...

Solr index

Storage resources

Data Capsules VM
Cluster

HTRC Volume
Store and Index

Remote
Desktop
Or VNC

Provide secure
VM

Scholars

Submit secure
capsule
map/reduce Data
Capsule images to
FutureGrid.
Receive and
review results

FutureGrid
Computation
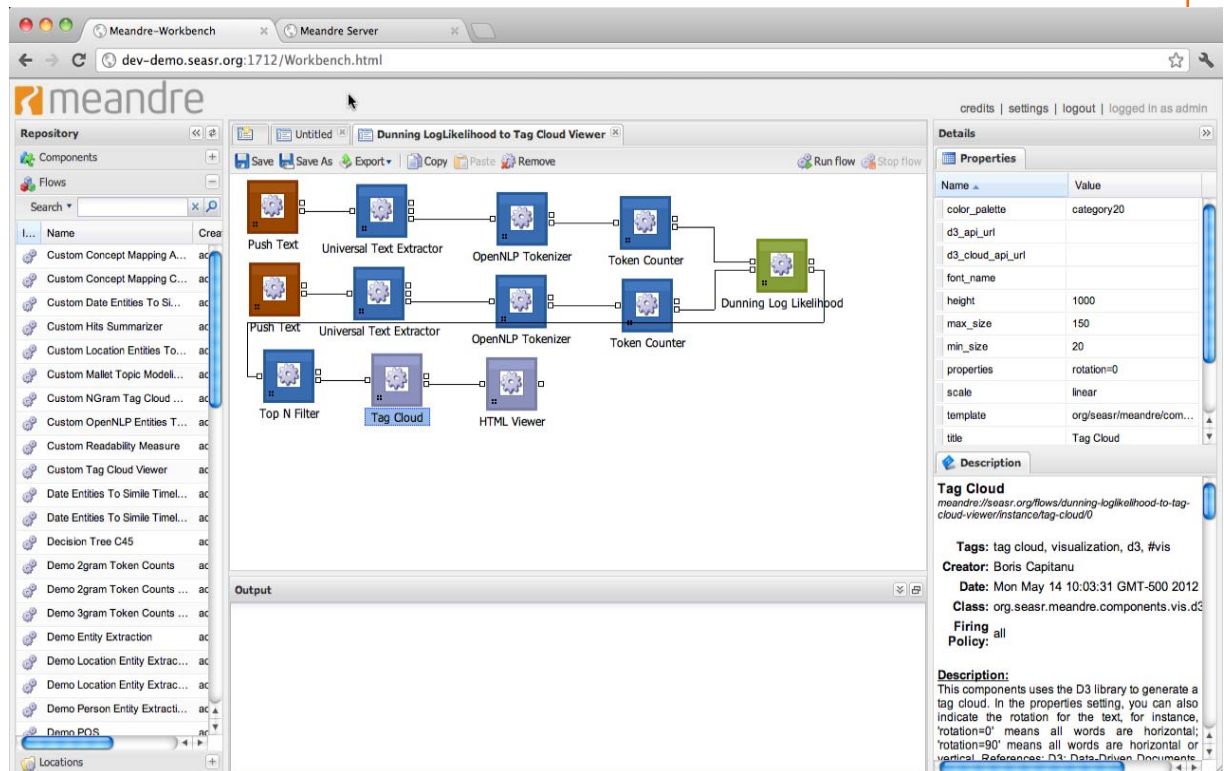Cloud

# Non-Consumptive Research-Secure Data Capsule

# Meandre: Workbench Existing Flow

- Web-based UI
- Components and flows are retrieved from server
- Additional locations of components and flows can be added to server
- Create flow using a graphical drag and drop interface
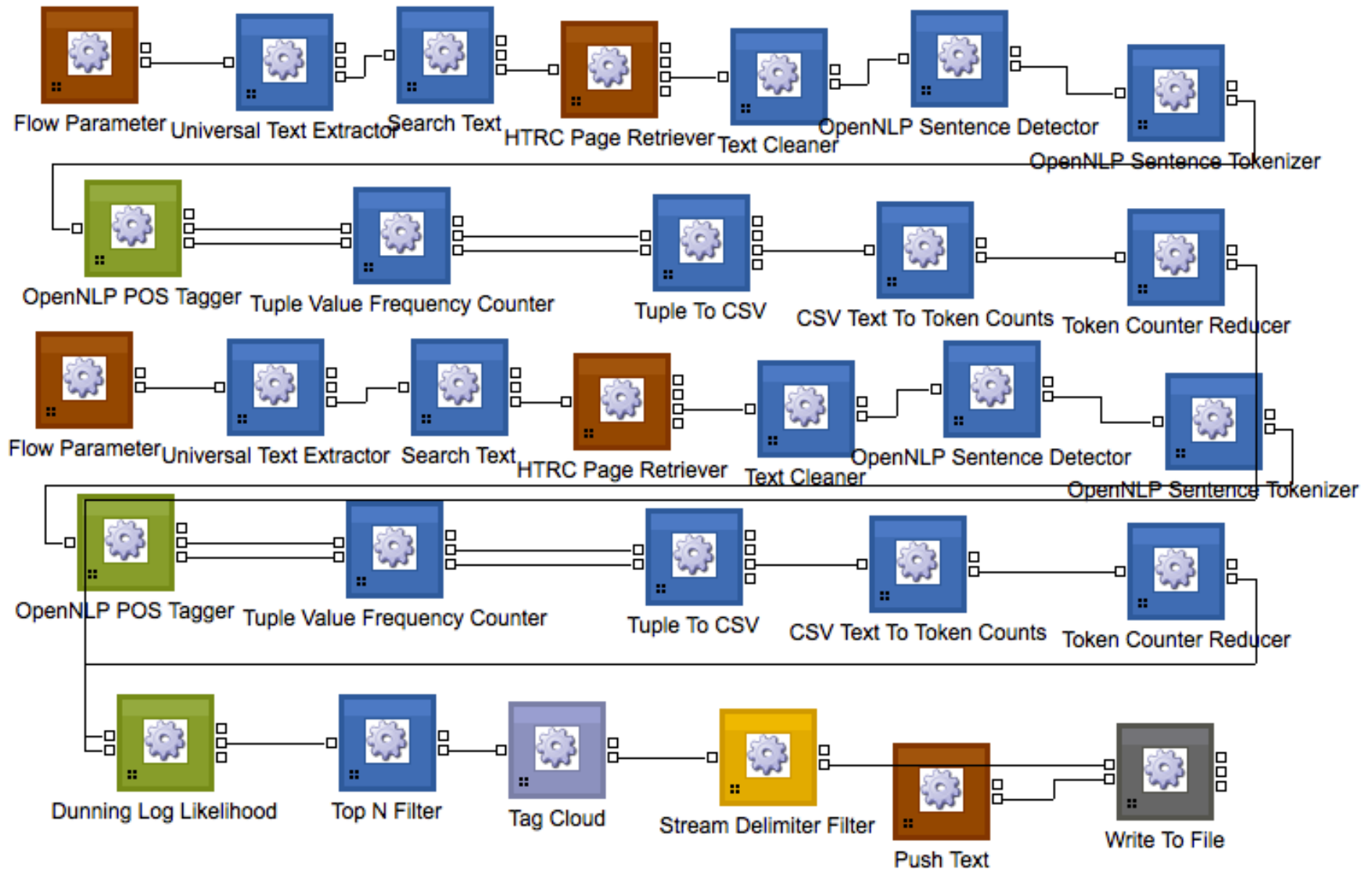- Change property values
- Execute the flow

# Meandre Flow

# Workset Creation for
# Scholarly Analysis (WCSA)

# WSCA Project Goals

- The **Workset Creation for Scholarly Analysis: Prototyping Project** (WCSA) seeks to address three sets of tightly intertwined research questions regarding:

1. **Enriching** the metadata describing the HathiTrust corpus through mining of the resources themselves and leveraging end-user annotations;

2. **Augmenting** string-based metadata with URIs to leverage external services and Linked Open Data to facilitate discovery and the process of organizing HathiTrust resources into collections and worksets; and,

3. **Formalizing** the notion of collections and worksets in the context of the HathiTrust Research Center.

# Motivation & Models

Collections, corpora, worksets, …:

- Aggregations of items brought together in some context:
  - Archival
  - Curatorial
  - Experimental
  - Referential
  - Thematic (for research)



Carl Spitzweg. 1850
*The Bookworm* (*Der Bücherwurm*)

# Grand Motivation

- *The ability to slice through a massive corpus constructed from many different library collections, and out of that to construct the precise workset required for a particular scholarly investigation, is an example of the "game changing" potential of the HathiTrust...*

# What is a Workset?

1.  A workset is an aggregation of materials brought together for the purpose of analysis.
2.  Worksets are conceptual and must be expressible in a variety of ways
    *   Need to allow creation outside of HathiTrust
    *   Need to facilitate inclusion of resources beyond HathiTrust
    *   Need to facilitate the inclusion of resources at many different levels of granularity beyond the book
3.  Worksets encapsulate the specific materials that underwent analysis.
    *   Need to capture provenance information
    *   Possible recording of parameters
4.  Worksets should be able to spawn descendants but otherwise immutable

# Dimensions of Workset Creation (Illustrative)

My work-set should contain (inspired by 2012 UnCamp):

- Volumes pertaining to Japan / in Japanese
- All volumes relevant to the study of Francis Bacon
- Music scores or notation extracted from HT volumes
- Images of Victorian England extracted from HT vols.
- Volumes in HT similar to TCP-ECCO novels
- 19th c. English-language novels by female authors
- Representative sample (by pub date & genre) of French language items in HT

# MARC Metadata Shortcomings I

| MARC Field | Percent of records in OCLC having instance of this field |
|---|---|
| 245 Title Statement | > 99% |
| 260 Publication Distribution, etc. | 92% |
| 500 General Note | 41% |
| 650 Topical Term / 653 Index Term – Uncontrolled | 39% / 13% |
| 050 LC Classification No / 082 Dewey Classification No | 17% / 13% |
| 655 Index Term -- Genre Form | 12% |

**Table 2. Frequency of MARC fields in OCLC Records**

# MARC Metadata Shortcomings II

| MARC Field | Percent of British Novel MARC records having instance of this field |
|---|---|
| 650 Topical Term | 6% |
| 050 LC Classification No / 082 Dewey Classification No | 27% / 4% |
| 655 Index Term -- Genre Form | 5% |

**Table 3. Frequency of MARC fields used in 2,386 descriptions of 19th century British novels digitized from UIUC collections**

# Why Worksets?

- The result of a first-level, rough filter

- Better scale for intensive analytics

- Provides essential scope for certain analytics
  - Word frequency scope over Bacon's essays

- Some tools (are trained to) work best on a narrow, homogeneous work-set

- Eliminate noise that would otherwise arise by asking questions across whole of HT

# Research Questions (Illustrative only)

- Can we enrich the HathiTrust corpus metadata by distilling analytics over full text?

- Can we augment string-based metadata with URIs for recognized entities – e.g., names, subjects, publication location, etc. -- and by doing so can we leverage external services to facilitate discovery and clustering of resources?

- Can we leverage existing, well-defined external corpora to identify complementary subsets of HT volumes, and having done so can we demonstrate the ability to create and perform analytics over an integrated workset that includes resources external to HT?

# Key Workset Questions

- Can we formalize the notion of collections and worksets in the HTRC context?

- What are the necessary elements of a "collection"? What are the necessary elements of a "workset"?

- How can we balance rigor with extensibility and flexibility?

- What roles do "data", "metadata", "annotations", "tags", "feature sets", and so on, all play in the conception, creation, use and reuse of collections and worksets?

# Two Project Streams

- Workset formal structures and semantics
  - Work in conjunction with Center for Informatics Research in Science and Scholarship at the Graduate School of Library and Information Science

- WCSA Prototyping Projects
  - Four projects funded by the grant but conducted by community teams

# WCSA Timeline

- July 2013:  Project Start
- Q1:   User needs assessments / focus groups
- Q2:   HT Corpus characterization
        Request For Prototype Proposals
- Q3:   RFP Finalist Workshop (Chicago) February 20
        Prototype experiment funding awarded
- Q4-6:  Prototype experiments done
         Metadata workflow & work-set modeling
- Q7-8:  Planning for prototype to production
         Report out
- June 2015: Project ends

# Prototype Grants

As part of project, HTRC will make 4 sub-awards

- $40K awarded to each of 4 non-HTRC teams

- HTRC will collaborate with each team
  - Access to representative test data / metadata set
  - Collaborate on work with HT / HTRC APIs, etc.

| RFP & Sub-Award Schedule | |
|---|---|
| 2013-11-15 | RFP Available |
| 2013-12-16 | Letters of Intent Due (preferred) |
| 2014-01-15 | Final Proposals Due |
| 2014-02-20 | Finalist Meeting |
| ~ 2014-03-15 | Award Notification for projects running April-Dec, 2014 |

# WCSA Summary

- Worksets are fundamental to the scholarly computational analysis enterprise
- We need a better understanding of their:
  - Constituent parts
  - Creation
  - Manipulation
  - Use and reuse
- Prototypes to lead to deeper tool development and metadata enhancement

# Next Steps

# Personal Goals for HTRC

- Engage in more collaborative projects
- Expand to have truly international partnerships
- Make sure to move beyond text
- Make sure to move beyond humanities!

# Redux: Ongoing Challenges

*How do we actually unlock the potential of 3 billion pages of human knowledge?*

- Data quality issues
- Data structure challenges
- Metadata shortcomings
- Overcoming copyright barriers to research
  - Non-consumptive research
  - Computation to the data
- Moving beyond text
- Community building important and ongoing

# Questions? Comments? Suggestions?

*Special thanks to:*

*Jeremy York, Stacy Kowalczyk and Loretta Auvil*