# From Seed to Harvest: Web Archiving Program Considerations for SUL

## Nicholas Taylor
## @nullhandle

# hello, my name is Nicholas…

# Library of Congress Web Archiving



The Library of Congress >> More Online Collections

## Library of Congress Web Archives *Minerva*   BROWSE | SEARCH | TECHNICAL INFORMATION

LC Web Archives

### Web Archives Available:
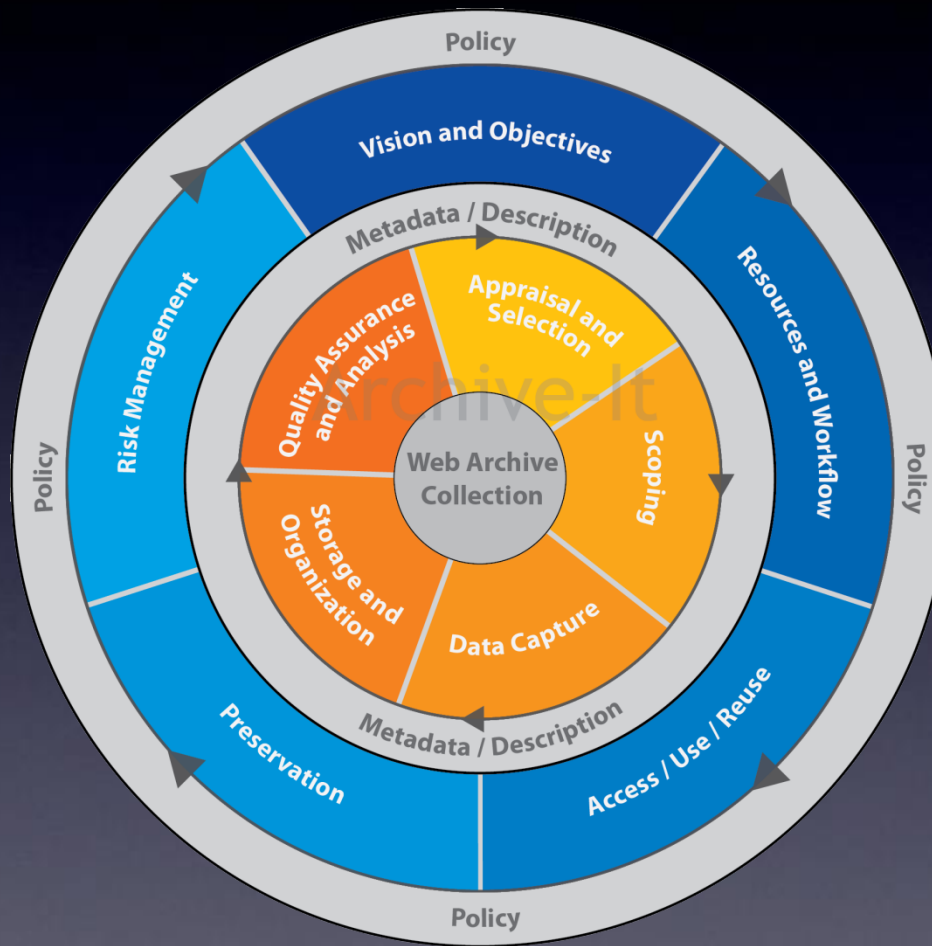
- Crisis in Darfur, Sudan, Web Archive, 2006
- Indian General Elections 2009 Web Archive
- Indonesian General Elections 2009 Web Archive
- Iraq War 2003 Web Archive
- Law Library Legal Blawgs Web Archive
- Library of Congress Manuscript Division Archive of Organizational Web Sites
- Papal Transition 2005 Web Archive
- September 11, 2001 Web Archive
- Single Sites Web Archive
- United States 107th Congress Web Archive
- United States 108th Congress Web Archive
- United States Election 2000 Web Archive
- United States Election 2002 Web Archive
- United States Election 2004 Web Archive
- United States Election 2006 Web Archive
- United States Election 2008 Web Archive
- Visual Image Web Sites Archive

The Library of Congress Web Archives (LCWA) is composed of collections of archived web sites selected by subject specialists to represent web-based information on a designated topic. It is part of a continuing effort by the Library to evaluate, select, collect, catalog, provide access to, and preserve digital materials for future generations of researchers. The early development project for Web archives was called MINERVA.

LC Web Archives

The Library of Congress >> More Online Collections
August 5, 2011

Contact Us
Library of Congress: "MINERVA"

# Web Archiving Life Cycle Model

# Web Archiving Life Cycle Model

**Program Elements**

- Vision and Objectives
- Resources and Workflow
- Access / Use / Reuse
- Preservation
- Risk Management

**Workflow Elements**

- Appraisal and Selection
- Scoping
- Data Capture
- Storage and Organization
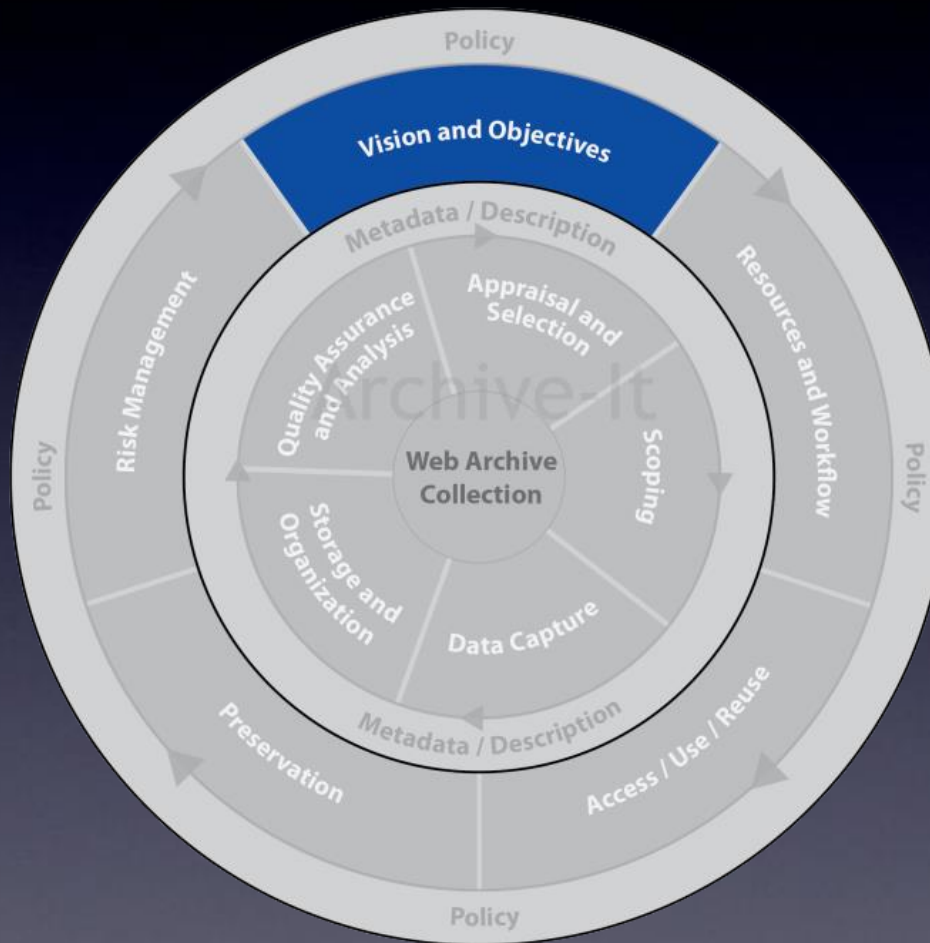- Quality Assurance and Analysis

Web Archiving

# PROGRAM ELEMENTS

# Vision and Objectives

# web archiving program vision

# SUL mission

*"The Stanford University Libraries (SUL) is more than a cluster of libraries; it connects people with information by* ***providing diverse resources and services*** *to the academic community."*

SUL: "About The Stanford University Libraries"

*"Stanford University Libraries...develops and implements resources and services...that* ***support research and instruction****."*

SUL: "SULAIR Brief Guide"

SUL: "Stanford University Libraries on Vimeo"

# DLSS mission

*"DLSS is the **information technology production** arm of the Stanford Libraries; it serves as the **digitization, digital preservation and access systems** provider for SUL; and it is the **research and development** unit for new technologies, standards and methodologies related to library systems."*

SUL: "SULAIR Digital Library Systems and Services (DLSS)"



SUL: "New Images of Rare Books and Digitization Devices"

# proposed program mission

*"The web archiving program will provide capabilities for the acquisition, preservation, and dissemination of resources that are increasingly and, often, exclusively accessible via the web that are necessary to support University research, instruction, and other purposes."*

# objectives

- build infrastructure
- develop expertise
- create research collections
- archive records and deprecated content
- mirror government documents

# Resources and Workflow

# cost modeling

# staffing

- service manager
- crawl engineer
- curators
- system administrators
- software engineers
- technical services
- legal counsel

# infrastructure



"Google Storage Server" by Flickr user Kazuya (Kaz) Yokohama under CC BY-NC-ND 2.0

# readily workflow-able

- collection management
- site nomination
- permissions tracking
- crawl scheduling
- data capture
- quality assurance



"Web Curator Tool User Manual Version 1.5.2"

# workflow challenges

- test crawling
- automated QA
- AIP/DIP generation
- SDR ingest
- indexing
- enabling access
- tools testing

# Access / Use / Reuse

# access policy

- dark archive
- data redistribution
- embargo
- onsite/offsite replay
- takedown requests

# browse and API: Wayback



Internet Archive: "Wayback Machine"

UK Web Archive: "Wayback Machine"

# many Wayback Machines

# discovery: Memento

# discovery: SearchWorks

# full-text search: Solr

# Preservation

# bit preservation

# preservation engineering

# Risk Management

# Risk Management

- "appified" web
- copyright
- ephemeral web
- financial sustainability
- fostering use

# Policy

# copyright

- § 108 (library exceptions)
- fair use
- notification vs. permission
- opt-out / takedown
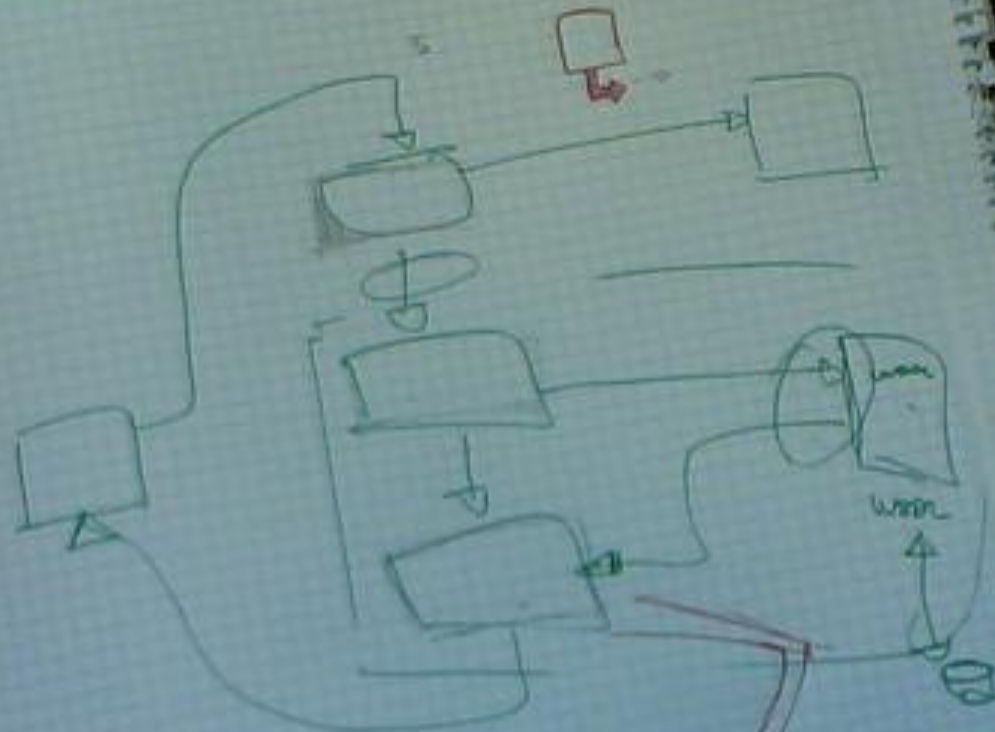- robots.txt
- third-party sites
- exceptions?

# collection development


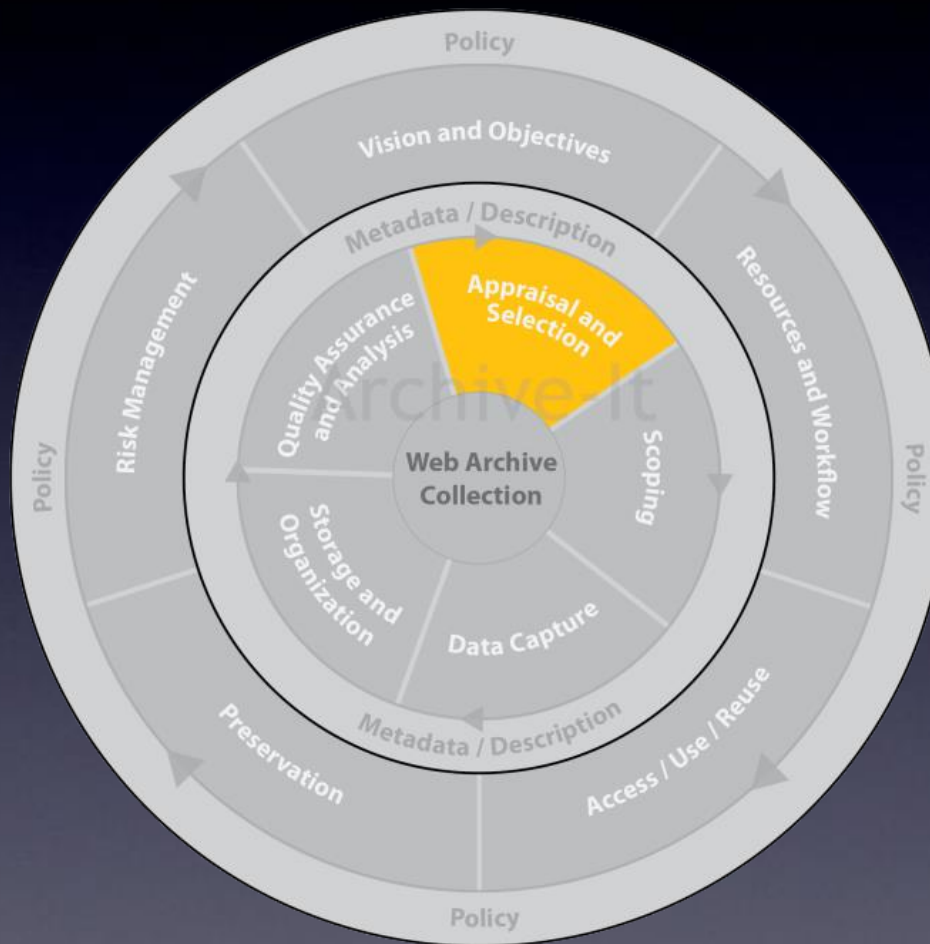
"leaf-cutter ants" by Flickr user Vilseskogen under CC BY-NC-SA 2.0

Web Archiving

# WORKFLOW ELEMENTS

# Appraisal and Selection

# informing selection

- value
- risk
- size
- extent to which archived

# TwitterVane



UK Web Archive: "TwitterVane"

# Wikipedia Live Monitor

**Wikipedia Live Monitor**

**Articles Edited Right Now:**

**Merging Right Now:**

None yet.
**Monitoring…**

**Article Clusters Repeatedly Edited In Short Intervals:**

None yet.
**Monitoring…**

**Breaking News Candidate Article Clusters Right Now:**

None yet.
**Monitoring…**

Proudly Developed In Europe By **Thomas Steiner** (tomac@google.com)
**It's Open Source** and Apache 2.0 Licensed (Fork This Project On GitHub)

Thomas Steiner: "Wikipedia Live Monitor"

# Wikipedia articles

# UNT Nomination Tool



University of North Texas Libraries: "Nomination Tool"

# Scoping

# the purpose of scoping

# Data Capture

# Heritrix

# other data capture tools

```
src/wget "http://www.archiveteam.org/" --mirror --warc-file="at"
```

Archive Team: "Wget with WARC output"

gwu libraries - social feed manager    about    search

## users by item count

| user | count | timeline (2012-01-01-2012-10-17) |
|------|-------|----------------------------------|
| ABC | 3,481 | |
| cnnpolitics | 3,345 | |
| BarackObama | 3,337 | |
| politicalticker | 3,301 | |
| CBSnews | 3,273 | |
| msnbc | 3,269 | |
| cnn | 3,262 | |
| ABCpolitics | 3,259 | |
| foxnews | 3,251 | |
| NBCpolitics | 3,237 | |
| Darrellssa | 3,233 | |
| NBCnightlynews | 3,231 | |
| NBCFirstRead | 3,231 | |
| SenSanders | 3,217 | |
| SpeakerBoehner | 3,211 | |
| foxnewspolitics | 3,208 | |
| jasoninthehouse | 3,200 | |

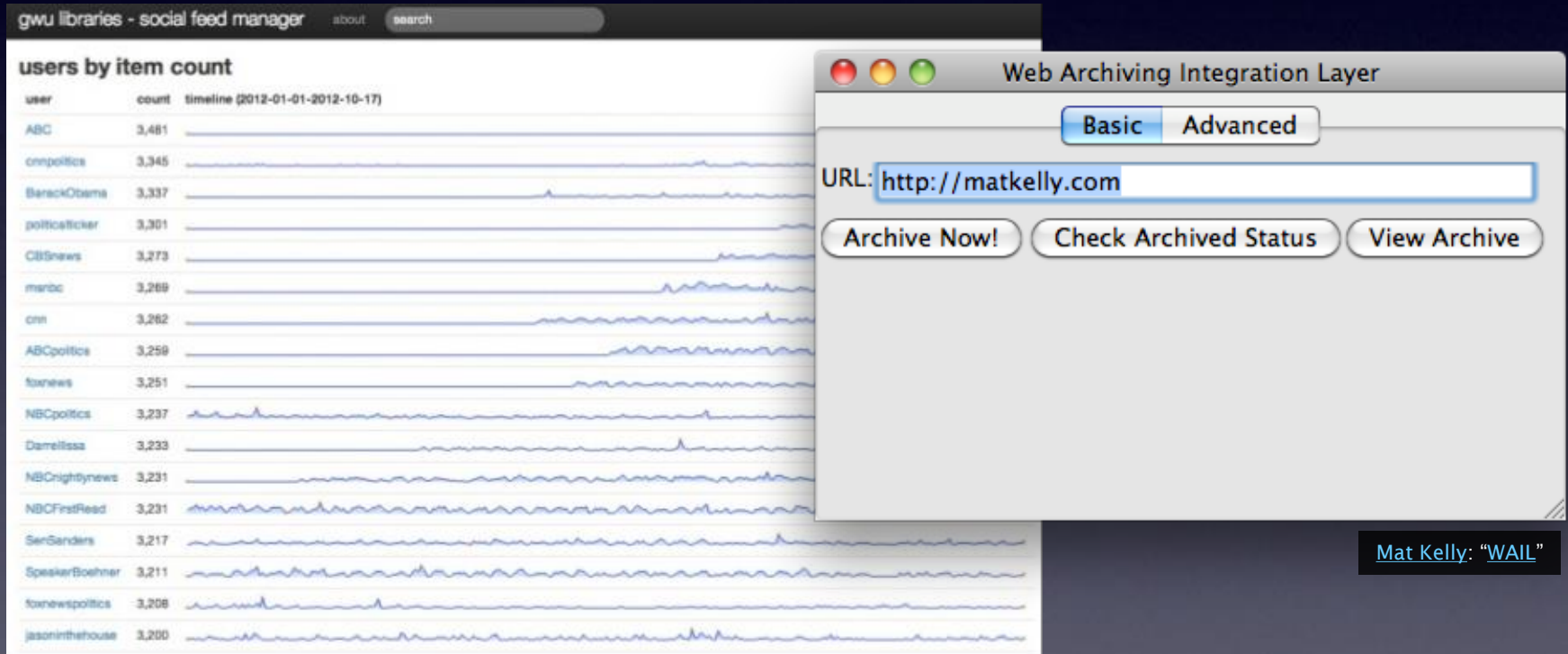Web Archiving Integration Layer

Basic    Advanced

URL: http://matkelly.com

Archive Now!    Check Archived Status    View Archive

Mat Kelly: "WAIL"

Dan Chudnov and Laura Wrubel: "social feed manager"

# the elusive web
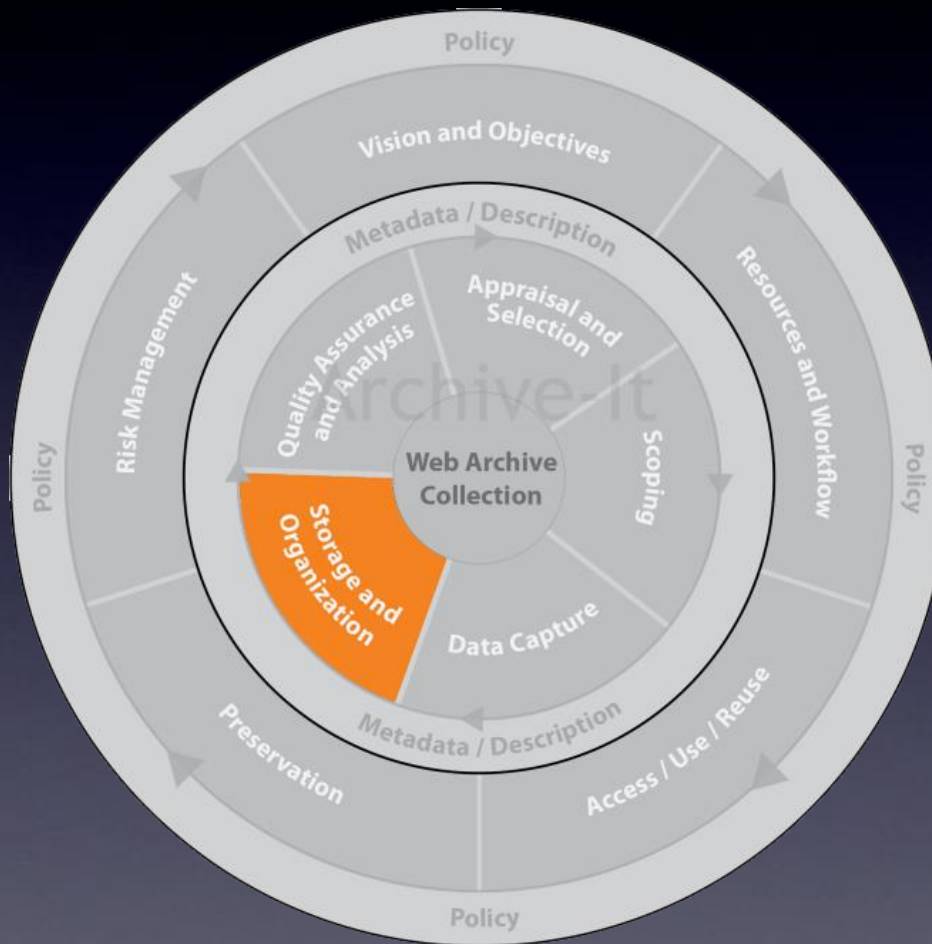
# scale



"chutes and ladders" by Flickr user reallyboring under CC BY-NC-SA 2.0

# Storage and Organization

# packages and their contents

# Quality Assurance and Analysis

# QA before, after, during

# Metadata / Description
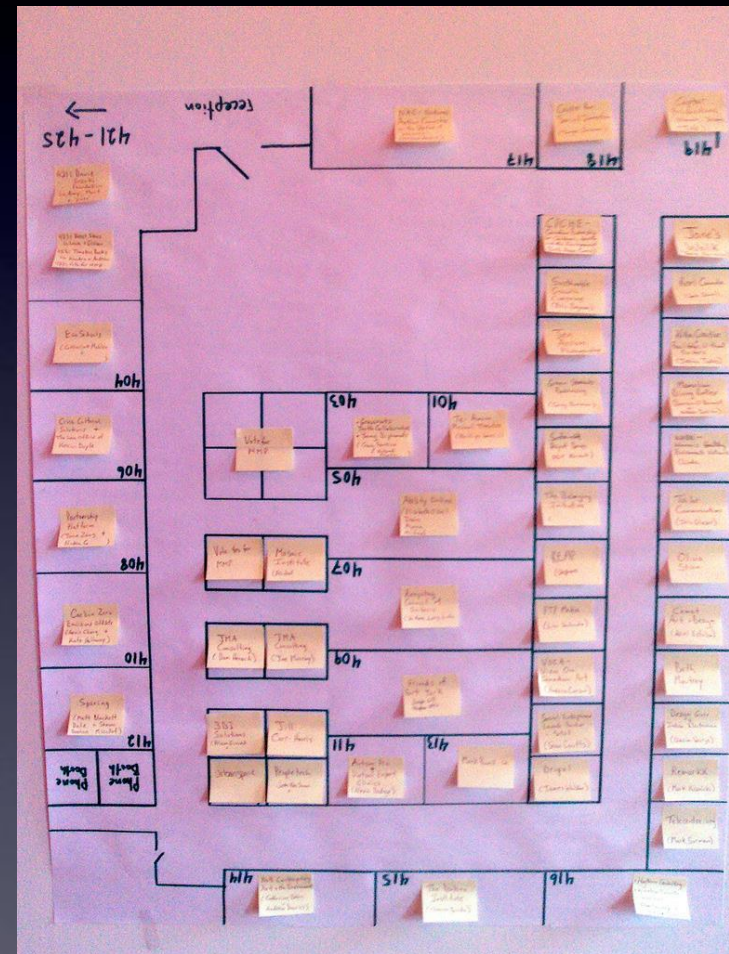
# Metadata / Description

Considerations

# BEYOND THE MODEL

# other program requirements

- marketing/outreach
- performance metrics
- service level definitions
- service roadmap
- training
- user documentation
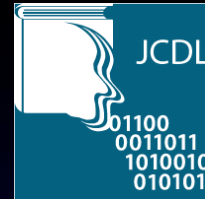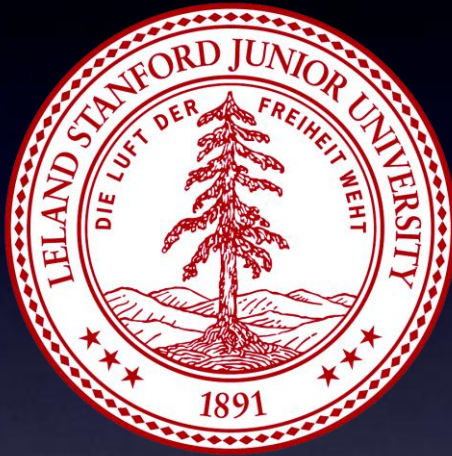


"Sticky notes" by Flickr user Kris Krug under CC BY-SA 2.0

# incorporating existing projects

- plan capacity
- normalize data
- ingest into SDR
- seek permissions
- process
- catalog
- enable access

# community engagement

# the web changes

Nicholas Taylor
@nullhandle