



June 3-5, 2015

Steve Ruggles

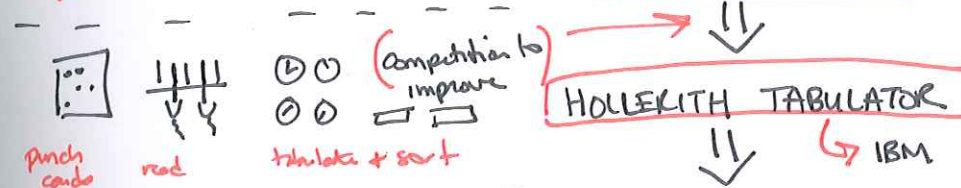
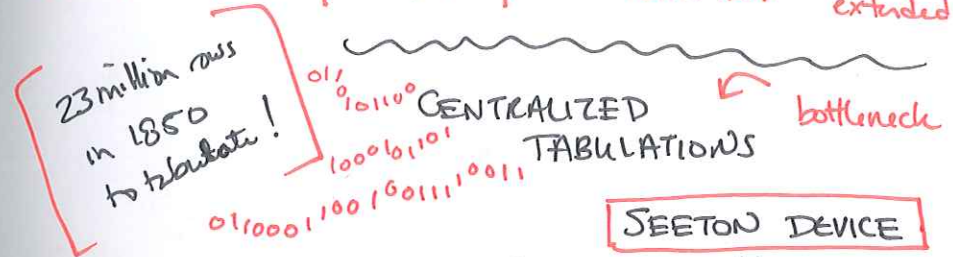
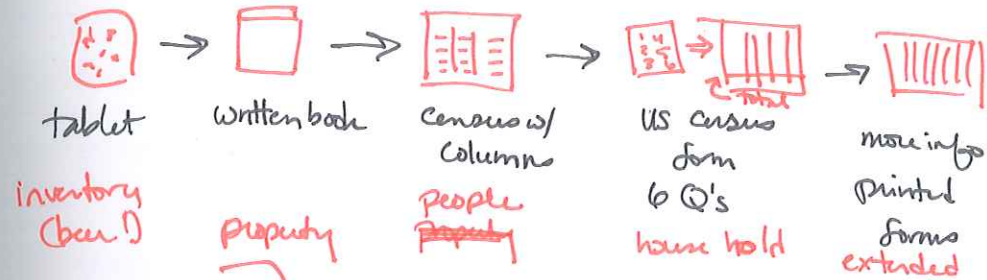
The History of Data:

Technological Change and The Census
1790-2020

Minnesota
Population
Center

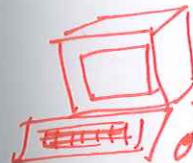


DATA CAPTURE



CENSUS BUREAU

directly or indirectly sponsored
two biggest computer companies
of early 20th century



UNIVAC
bought by
Remington
Rand

1960 mail out census form



photo ~~scanning~~ reading
for collection



* purchased

1967 invention of GIS

1970 mailed out
AND mailed back



used data of #2 pencils

Compared w/ 1960's data due to compatibility

1980 } faster computers ...
1990 }

2000 data capture outsourced to Lockheed Martin

lots of scanning, lots of errors

2010 "was a colossal failure"



ended up w/ Lockheed again due to

handheld device? failures by other contractors, Harris

CANCELED internet response option

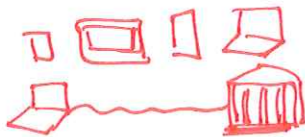
↑ CANADA did this no problem

2020? Not looking good...

bring your own device

internet response option

neither looking good



US CENSUS TIMELINE

Recover?
Compare old data

1960-1990



Build longitudinal study
from 1940's on...

CLIP



link to other
existing microdata

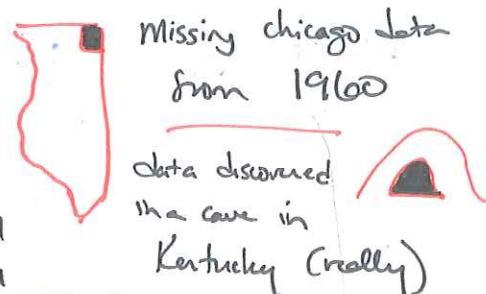
truly an extraordinary
moment to be a
data nerd!

CENSUS BUREAU

→ have huge innovations
in computing in
20th century

→ changed recently with
outsourcing


→ still ahead in data
(depth, breadth, free)



data discovered
in a cave in
Kentucky (really)






SESSION A5 - RDM SERVICES

 Jungwon Yang
University of Michigan

"Stakeholder Analysis for RDM services for public policy researchers"

SERVICES @ UMICH

- teach each other
- data ed working group
- program-oriented RDS (not project)

TENURE   stakeholders  faculty's main interest

publish grants

84 Faculty
41 research
37 clinical

Journals

- doing a lot with data
 - stats
 - micro / macro
 - code books

Grants

- applied for / 48 have federal grants (NSF + NIH)
- IRB big component
 - sensitive data

<CLOSUP>

- survey group on policy
- had own data curation specialist!

ISSUES

- reliable / trusted repository
- preparing data for sharing
- anonymizing sensitive data

WANT TO SHARE DATA
(persuade from APSA)

- guide for anonymizing sensitive data
- data education for students
- work on data repository

FUTURE PLANS

Maya Ishida
U. Manitoba



Sarah Williams
UIC

Developing Research Data Services Vision(s):
An analysis of ARL Libraries



← Lots of Libraries developing services

72% had RDS web page
48% other related pages

① Driver docs (fed req's + memos)

② ARL Libraries web pages + strategic doc

③ Semistructured interviews w/ admin

PRELIMINARY FINDINGS

24% DMP assistance

- Consults / instructions
- Data deposit / repos
- Storage
- Sharing / re-use
- Best practices / info dissemination
- Archiving / preservation
- Data processing / Analysis
- 3% Metadata

ALL SERVICES

<Strategic planning>

□ 56% mention data in library-wide strategic plan

□ 7% have data-focused vision

INTERVIEWS 

COLLABORATION AND ENGAGEMENT

library can't be only place supporting data partner not just service

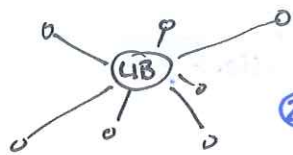
ENABLING RESEARCH ACROSS LIFECYCLE

STRONG EXPERTISE

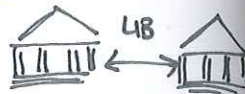
integrate and bring others on board

VISION THEMES

① library should be HUB of RDM network



VISION SUGGESTIONS



② provide support across disciplines

U Minnesota | Amy West
Amy Neesen

A Coordinated, Decentralized
Approach to RDM services:
from education to everyday

we don't have a choice but
to do things as cooperatively
as we can

UMinn

- varying levels of expertise
- culture of decentralization
- dept/working group silos



WHAT WORKS

- o at least one person who facilitates
- o Communication {share content!}
- o freedom to experiment

Never an end to relationship
building



maintain communication &
workflows between groups

.....

Quantitative
Qualitative
Simulation
Other

rough breakdown of
data on campus

Use this to tailor
teaching by discipline

DISCIPLINARY VARIATIONS

- * Acknowledge diversity
- * Focus on users
- * Develop general content
customize as needed
- * Iterative simultaneous efforts to train
 - library staff
 - researchers

<STAFF TRAINING>

just start!

like an
interview
training

Slightly different training
than researchers

RESEARCHER TRAINING

← leverage RCR requirement

734 attendees since 2009!

- ① wave
come one come all!
- ② targeting different audiences
students, faculty
bring other librarians on board

SPECIALIZED

- ▷ Civil engineering course
↓ repurposed
- ▷ life science cabin
- ▷ Social science workshops
- ▷ grad students
- ▷ DMP consultants



It takes time ...

... and you're never done

SESSION BS - Building on Common Ground

Lizzy Rolando - Georgia Tech

Kelly Chatain - University of Michigan

Bethany Anderson - UIUC



Research Center

Kelly

Survey Group hired to do data archiving/records management

20-30 studies per year
~140 employees

Records Management

- Assess The situation
 - regulatory / compliance issues
- Records Survey
- Research / Business Task Analysis

Strengths and Challenges

Retention Schedule

- granular
- accessible
- annually reviewed
- actionable

engage researchers upstream in data lifecycle



Targeted Guidelines

- * Email
- * File naming / digitization
- * Archive checklist
- * Project closure

Training

Annual updates

OUTREACH

Records Day

"From The Archives" articles



encouraging good behavior
improve documentation
and more positive changes!

SESSION BS

it started with a thesis...



that had a dataset...



this was a snes

Lizzy R.

Intersection of



Archives (Wendy)



Data (Lizzy)

ARCHIVES would maintain in current form
would want to fix in DATA SERVICES



limited technical resources

Borrow processing from ARCHIVES
learn about digital sustainability from RDS

how do we handle records?
how is it different b/c different services?



Can we build

1

System for Archives & Data?

THEORY

DATA

+ prioritize use over preservation
+ context valuable w/ respect to original material

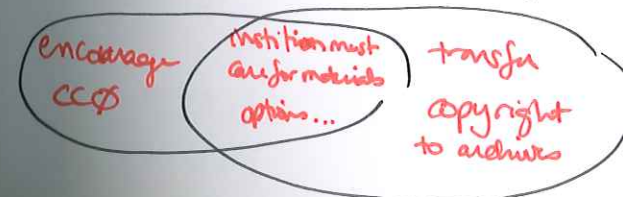
+ prioritize preservation over use

+ context valuable on its own

< different priorities >



need to determine retention and deaccessioning



messy stuff here...



Borrow workflows from archives



Trying to share tools as well...



ARCHIVES

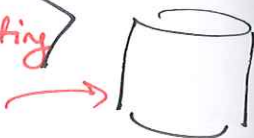


DATA SERVICES

<Betnamy> Appraising and Preserving Data in Context

<Appraisal means identifying and selecting>

Many reasons to appraise



How to appraise

- mission alignment
- institutional compliance
- scientific/historical value
- uniqueness
- potential for redistribution
- non-replicable
- economic factors
- contextual documentation and metadata



expensive to keep

maintain everything

Storage is going to be a problem w/out careful selection

Haas, Samuels, Toppel, Simmons
Appraising the Records of Modern Science and Technology
1985



[What context is needed for data?]

methods, how data was acquired, etc.

PRESERVE
MEMORY

better than data

C3: Data Sharing Behavior and Policy

"The road to data sharing is paved w/ good intention"



Laurence Horton
(London School of Economics
Political Science)

Astrid Recker
(Leibniz Institute for the
Social Sciences)

More policies
from UK (20-30%)
than Germany & France

UK: strong centralized funding
Germany: federal structure strong
tradition of independence

UK EPSRC

diving this trend



puts requirement
on institutions
to support data

< RUSSELL GROUP >

(high research institutions)

not huge difference b/t
in or out of group

250 - ~1000 - 6000
WORD# average

- o short strong values
weak on details
- o longer on requirements

two general Policy types

- what to keep
various specificity
- data retention time
when given, 10yrs for
10 years
- ethical use/reuse
- access
- availability
(open access)
- cost of DPM (open access)

different work for this
type conditions

[POLICY CONTENTS]

- define data (i value)
- identify university roles/
responsibilities
- DMP requirement
83% required!
- who it covers
(Staff v. students)
- ownership
27% state owned
- external funding reference

□ subject to review
(review after X years)

table available via DEC
data available via UKDA

(seeing data under IP policy
in a lot of "new" universities)

In Australia, you
can copyright a phone-
book. Wow.

QUESTIONS

{ Army Pienta }
ICPSR

Looking at ICPSR, do some social
and behavioral sciences disciplines share
more than others?

SURVEY

- o NSF & NIH PIs
- o social / behavioral science data
- o range of topics on data sharing
 - behavior
 - perceptions
 - barriers
 - publications



Data Sharing

12% self archived / IR
45% informal sharing
website
by request
43% no sharing

879 respondents in
this study

See similar results
in observed sharing
and actual sharing

- poly sci / econ more likely to
archive than
- psychology less
likely to share +
informally archive
- psychology / health,
informal preferred
- health sciences less likely to
archive
- poly sci / econ more likely to share informally
+ archive

remove
demographics

Confidentiality
Lack of time
documentation
won't understand

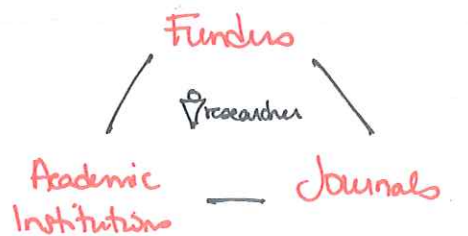
Barriers to Sharing

different main
barriers b/t
disciplines

target sources
based on this

Sharing behavior/
news different
by discipline

Alexia Katsanidou
GESIS - Leibniz Institute
for the Social Sciences



Big Journal
Journal Data Policies

Political science

higher impact factor
more likely to have policy

BUT researchers don't always comply...

ask authors why
shared data?
(same journals as
previous study)



questionnaire

initial results

- researchers less likely to share if
 - others may find something new in data
 - or might be scooped
 - have to use standard (raw) data
- more likely to share if
 - data will enhance research
 - give back to community
 - suitable repositories available
 - attitude of data sharing

ASSIST Day 2



DIVIDE



Curtiss Cobb-Facebook-

Measuring the Digital Divide:
Using Existing Data Sources and New
Data Collection to Understand Between
Country Differences

The internet is all about services
to people and communities, allowing
them to achieve their potential

⇒ extend internet
access to developing
world



of the world's 7 billion
people

only ~2.7* have access



*mckinsey

Key Facebook principles
internet access



The Plan: Lower cost x less data
mobile-based internet

Barriers of



- access
- education (literacy)
- gender divide



Unconnected demographics

- * poor
- * rural
- * illiterate
- * age



How do we measure progress?

need to have #'s for background
but hard to do in many countries

run into a lot of
social issues here

Numbers
unknown
Forecasts are
black boxes

How do we use
3rd party data?

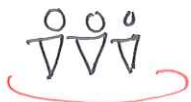
- price/availability
- source
- timeliness
- relevance
- just check

DATA
EVALUATION

Millions of Facebook
users have no idea
that they're using
the internet

CONTEXT AND
CULTURE MATTER

D1: Data Professionals



Adetoun Oyedele - Data Professionals' training challenges in dynamic work environments
University of Ibadan



How do we train ourselves to get ahead?

Interviewed 10 data professionals } different positions

< Self-training prominent >

CHALLENGES

- o Funding for training
- o infrastructure
- o bosses not understanding need for training
- o Understanding technicalities
- o Knowing of training opportunities
- o lobbying for time for training
- o Synchronizing time for training online

SOLUTIONS

- + Sponsorship
- + Social media
- + Showcase results
- + reading grey lit
- + Learn lobbying
- + grant skills
- + work around obstacles
- + politicizing issues
- + Sacrifice



Interview results from 10 professionals in Nigeria

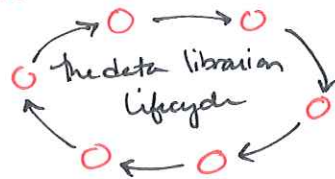
WHAT MORE?

- Support from above
- Professional org support
- trainers should train others

- Share via social media
- institutional mentoring
- invest in ourselves

use Constructive analogies

Michelle Edwards [CISER]



Stolen Borrowed from Research Data Lifecycle

THE CYCLE:

concept new source collection community info processing needs assessment Archiving write it down Distribution how it location Discovery marketing Analysis

Respond to change => repurpose

- meet w/ local librarians
- learn what's been done
- create new information/ processes
- raise awareness & establish hours
- integrate into community, campus & library
- is it working?

M. Edwards thinking about this w/ respect to recent job move



Guelph



Corneil

Line Bouchard • PURDUE

Comparing policies for open data from publicly available accessible institutional sources

Had to negotiate agreements b/t provider & Purdue
concerned about data reuse & access



EXAMINE OTHER COUNTRY'S POLICIES

(data protection directive)

CAM² project

760,000 webcam streams



Keep data on Amazon servers b/c it's BIG

(don't store stream, just data analysis)

often

PRIVACY contradicts REUSE



UK has most cameras

(estimates that you'll be photographed ~500x/day in London)

US has patchwork of regulations



(different for personal & industry)

click-time consent

regulations written before arrival of BIG DATA

→ policies not suited to BIG DATA



VARIETY OF POLICY PROVISIONS

ES: Restricted-Use Data Support in Academic Libraries



Jen Darragh - Johns Hopkins
 Ryan Womack - Rutgers
 Jen Doty - Emory
 Jamene Broder-Kieffer - U. Kansas
 Sarah Irwin - Penn State

Physical restriction

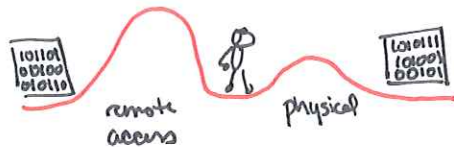
grad students don't have own space

Remote access

much more technological many issues here



library IT not here



How do researchers store sensitive data?



Jamene Broder-Kieffer

Survey of institution researchers producing PII (not clinical/medical)

Researchers

Storage

Developing services

△ Ryan - student need + extra year-end \$

Sensitivity

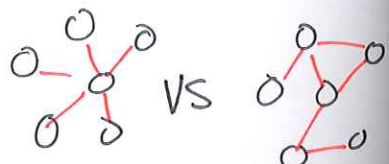
△ J. Darragh - limited resources in one area pushed services in other areas



Special room for limited access

△ Sarah - service limitations due to funding source

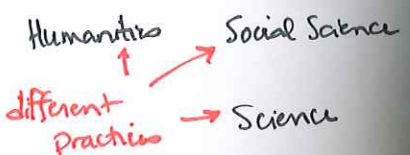
△ J. Doty - lots of anecdotal evidence (e.g. ICPSR application)



Centralized VS decentralized

overall institutional structure MATTERS

Contains identifiers or triggers uni classification policy SENSITIVE DATA



Institutional Networked storage not researcher top choice

Researchers use

- 1 local computer
- 2 hard disk
- 3 DVD / CD / USB
- 4 Cloud storage
- 5 institution storage
- 6 institution server
- 7 institution cloud
- 8 HPC cluster

Researchers want EASY

researchers using personal accounts here

PROBLEMS WITH NETWORKED STORAGE

- low capacity, high cost difficult to use
- insufficient support for sensitive data
- sync & collaborate difficult
- incompatible work conditions

Mismatch between needs and services



IF researchers need help w/ own data

THEN need help with sensitive data

J. Darragh: library will be your security (physical & electronic)

Ryan Womack: library workflows? policies for use of restricted space

Can only do so much when there's another person involved

Library can't have total control...

re: security

Sarah Irwin:

Takes time to develop services

working group looking @ cross-campus workflows



little monitoring of output from this process



what is library's role in this area?

J. Doty: sensitive data is part of RDM service, partner with IRB

training on mitigating disclosure risk

J. Darragh



- IRB
- IT
- Contracts people
- individual departments
- Campus working groups
- research office

stakeholders

LIABILITY

- issues with admin
- misunderstanding risks / services

+ keep services limited

+ schools know of these issues

\$ of Services
(currently + going forward)

□ J. Darragh \$10-15k + IT help
(room + tech setup) *

□ Ryan W <\$2k initial
~\$3k for new lock system +

Is this part of
regular data services?
Should it be?



J. Darragh

↑ Isn't at
Johns Hopkins
maybe should be

EDUCATION!

Need to be in addition to storage

@bjunul:

← Bo Wandschneider

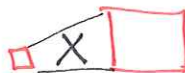
- User is right in wanting it easy
- institutions need to do better
- Cloud is doing it better than we probably can
- Partner w/ IT



LIBRARY

Solution that helps while
university solves (hopefully)
larger issues

*not scalable,
need more manpower



may end up out of
library

(+ again, not scalable)

BIG MESSAGE
ABOUT SCALE



↑ if can't provide service

be CONCEALGE

↑ stay just a bit ahead of
users...

F5: Using DMPs as a Research Tool for Improving data services in academic libraries

Amanda Whitman - Oregon State

Lizzy Bolando - Georgia Tech

Brian Westra - Oregon

bit.ly / dmpresearch
url for the project

SO MUCH DATA!

researchers don't know how to handle data



libraries can help!

but need to audit this as verifiable

inform library services



DMP reflects

- + knowledge
- + capabilities
- + practices
- + needs



NSF requirements different b/c divisions



general rubric



used to test rubric to make sure it's clear

How develop services / understand researcher needs?

- o surveys
- o interviews
- o review DMPs

this project

Need a tool

Using a rubric

GOAL to make rubric accessible to anyone

tested on real DMPs and entered assessment via Qualtrics

careful about inter-rater reliability and other errors

* plans from funded projects

* intra-class correlation



Lizzy R.

using dmps to look at institutional info and inform data services

Results of analysis

- o 8 plans no data
- o 5 plans from math - not producing data?
- 42 analyzed

Previous work

- o 40% mention IR
- o researchers share text

used plagiarism software

correct misinformation put more info online (instead of come to library) learned more than interviews or surveys

REDUX

researchers get sharing message

less good at archiving

BUT repository does both!

bad at documentation (not even metadata) sad

share data via journals difficult / more work here

- o good at describing data
- o bad at describing metadata
- o mixed re: formats
- o bad at data availability where
- o mixed at public availability
- o bad at security measures
- o N/A at PI / IP data
- o bad at reuse
- o bad at redistribution
- o bad at production of derivatives
- o mixed at archiving / preservation
- o unfortunate data sharing places journals?

Intervene @ school of math improve web presence and add boiler plates

FUTURE STEPS

communicate to campus

repo-technical reg's metadata and documentation

<lots of Qs / work here>

Brian Westra



NSF biology guidance

data AND metadata

appear together

assess
⇒ Post-award

now moving to
annual
report

BIO sharing

- personal website
- on request
- named data centers
- journals

What is data?

How do you
publish it?



still an
issue



CHEM Guidance language

- peer-reviewed journals
- personal website
- data repository

*this is me
shaking my fist at
my fellow chemists

↑ actually are
sharing "data" as
journal supplements

Named data centers

crystallographic
data big one

NEED MORE

DATACENTERS IN
CHEMISTRY



Lots more focus on by-request
<collective sigh>
actually in guidance language

Plenary 3

Johnson:

- △ Data analytics
- △ web development
- △ entrepreneurship
- △ interest in politics



Bill listing where tax money goes

"What we Pay For" website

TRANSPARENCY



People like this

FEAR OF OPEN DATA

- open data will make us look bad
- without context, data could be misinterpreted
- people could do anything with this data!
- will cost us millions!
- folks will have an unfair advantage

digital divide in city



ran for City Council

HE WON!

IT Subcommittee
Budget director
work w/ CIO



DATA VIZ CHALLENGE
via google

"it will take 4 years to do"

BENEFITS

- + public trust
- + shifts relationship
- + reduces requests
- + better data quality
- + empowers staff
- + empowers community to further city goals

seen examples of this already

POLICY
PORTAL
CULTURE

challenge

- legacy systems
- open data v. open process
- prioritizing data
- portal politics
- political process

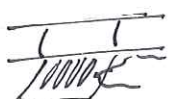
GOALS

empower leaders
gain buy-in

success stories
more pressure to be open

dataset can make

from main data



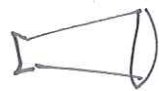
The Politics of Open Data

Andrew Johnson
- Minneapolis City Council -



last-minute attempt to make policy optional

POLITICS



outreach to stakeholders

↑ restaurants: food safety info going on line



NEXT STEPS

- compliance report
- revisit portal
- events for codes
- continue culture change
- revisit policy

add open policies into purchasing

16th city w/ open data policy

G2: Planning RDM Services

Kate McNeill
MIT Libraries

RDM ↔ OA

- Focus on open access to information
- Focus on information producers
- Compliance / funding reqs
- Publisher policies

Complexities

Overlap w/ other areas

Whole lifecycle

Complex collaborators

RDM

Copyright complications
end of lifecycle
more familiar

OA

DIFFERENCES

RDM Meets Open Access
(RDM = Open Access)



In MIT Libraries...

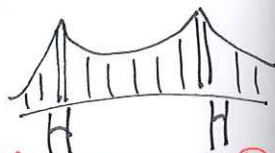
"data is specialized services"

Using the team approach



OA } office of scholarly publishing
copyright, & licensing

OA logo! where are RDM logs?



OA ↔ Data

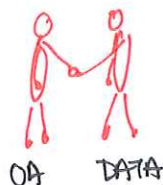
make services work together

Use bottom-up approach but make structural

Working jointly

- Communication
 - in library
 - on campus

- increased awareness
- synergies



OA DATA

JOINT GROUP

Coordinate but not provide services

[High-priority library project]

+ OSTP Memo
+ share personnel

meets bi-weekly

COLLABORATE

UNCERTAIN

who does what

with

authority

presenting

other

service

strategic planning?

level of collaboration?

Partnerships in a Data Management Village

Researchers get a lot of services along data journey



how do we reach researchers more?

need to add data steps to research process

IT TAKES A VILLAGE

Where are our village people?

lots of places in uni to connect

+++++

Connectivity & Awareness

more beyond network to partners

most / most engaged partners

- STRATEGIC FOCUS
 - policy

- COLLABORATION OPPORTUNITIES
 - community of practice
 - RDM library group
 - design thinking workshops

U MINN CAMPUS

Tom L.

[College of liberal arts + Libraries]



how long is the new forever?

Lisa J.



↑ re: data retention

- DATA VILLAGE

- does not happen
- overnight, work in progress!

has to prep data for sharing

- documentation templates
- de-identification
- DRUM repository
- curation model (distributed)
- enable / track reuse

think beyond each village

Carol Perry
University of Guelph

data management on a shoestring budget

00000
▽▽▽▽▽

} training from
various sources

5 members of
data + schol comm
team

REPOSITORY
PLATFORM
Database +
DataCite

SCALE

- focus local
- not huge branding
- pilot
- small data sets

have
manage
stuff?

COLLABORATE

- + library
- + office of research
- + computing & communication
services



DEVELOPED
SUPPORT
MATERIAL

website
modules
workshops



Planning for upcoming
DMP requirements

RDM Services

- Environmental scans
of services
 - models
 - lit review
 - brief faculty survey
 - interviews deans /
dept chairs
- DMP support
- Support for new funder
requirements
- Data consultation
- Data curation

OUTREACH

Outside funding for
repository
research data canada
seminar
advise another grant
agency on setting
up repo
consortial ~~memberships~~ ~~collab~~

creation of DMP
template
(Portage)