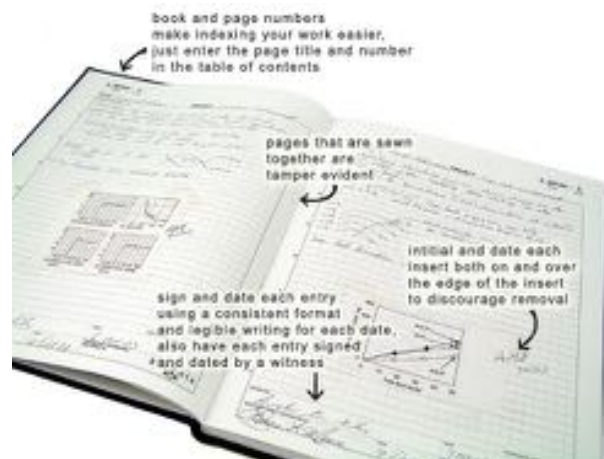




## Chain of Data Creation

1. Preparation
2. Creation of Metadata
3. Acquisition
4. Building a Permanent Record
5. Data Management
6. Storage
7. Data Sharing

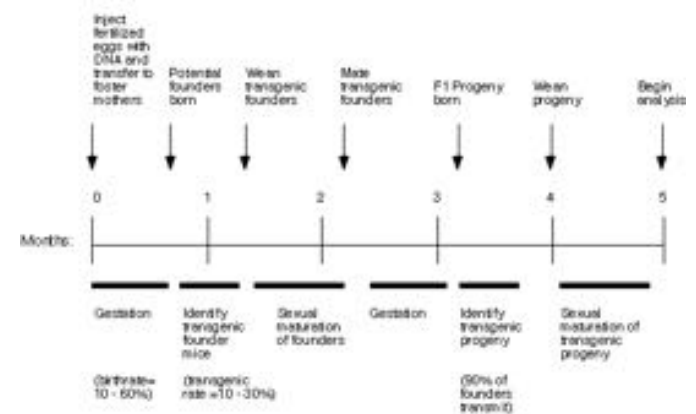
## Lab Notebook



- Record of hypotheses
- Record of Protocols
- Second brain

## Plan Your Experiment, Experiment With your Plan

### Timeline for Transgenic Mouse Analysis



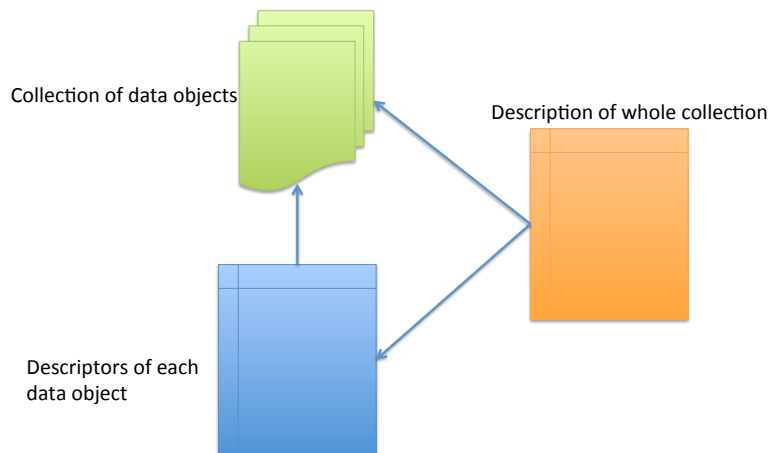
# Question to Ask About your Data Collection Activity

- What am I measuring?
- When am I measuring it?
- How am I measuring it?
- What are the tools I am using?
- What about the lab/field environment do I need to know?
- Is my protocol reproducible?

# Meta-Data



# What is Metadata?



# What is Metadata?

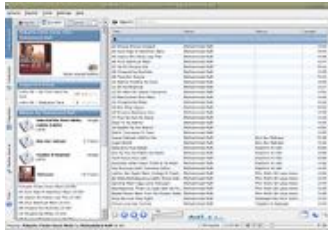
## Metadata is: Data 'reporting'

- **WHO** created the data?
- **WHAT** is the content of the data?
- **WHEN** were the data created?
- **WHERE** is it geographically?
- **HOW** were the data developed?
- **WHY** were the data developed?



## Metadata in Real Life

- Metadata is all around...



CC image by Nilsa on Flickr

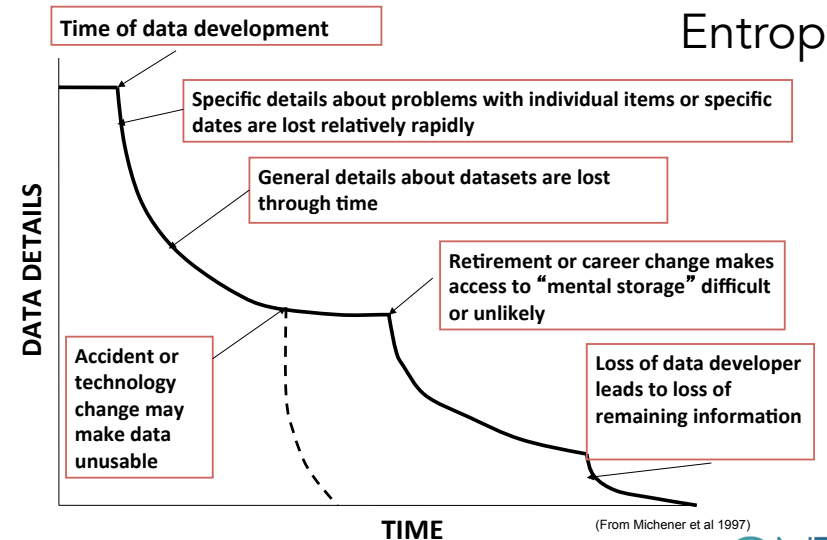
Nutrition Facts	
Serving Size 4 OZ. SERVING (112g)	
Servings Per Container VARIED	
Amount Per Serving	
Calories 170	Calories from Fat 70
% Daily Value*	
Total Fat 8g	12%
Saturated Fat 3g	15%
Cholesterol 65mg	22%
Sodium 70mg	3%
Total Carbohydrate 0g	0%
Dietary Fiber 0g	0%
Sugars 0g	
Protein 23g	
Vitamin A 0%	Vitamin C 0%
Calcium 0%	Iron 15%

CC image by USDA.gov on Flickr

**Author(s)** Boullosa, Carmen.  
**Title(s)** They're cows, we're pigs / by Carmen Boullosa  
**Place** New York : Grove Press, 1997.  
**Physical Descr** viii, 180 p ; 22 cm.  
**Subject(s)** Pirates Caribbean Area Fiction.  
**Format** Fiction

DataONE

## Information Entropy



DataONE

## Information Entropy

DATA DETAILS

**Sound information management, including metadata development, can arrest the loss of dataset detail.**



TIME

DataONE

## Data Management via Metadata



DataONE

# 'We Kill People Based on Metadata'

David Cole

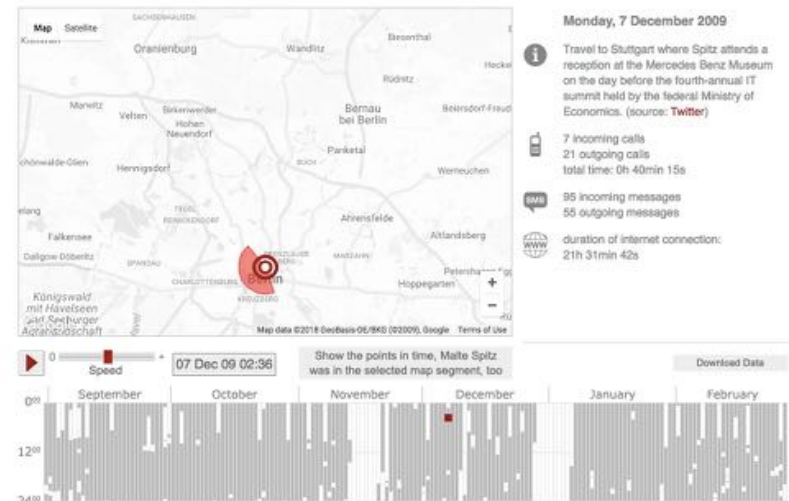


Rick Bowmer/AP Photo

The National Security Agency's \$1.5 billion data storage facility in Bluffdale, Utah, June 2013

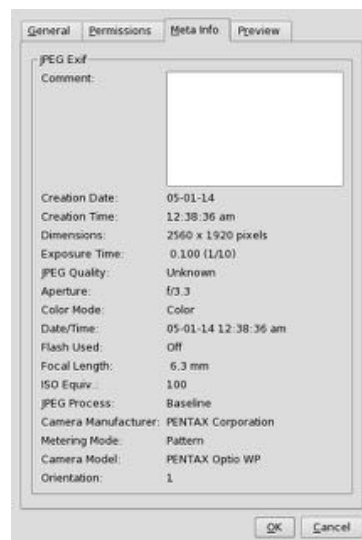
Supporters of the National Security Agency inevitably defend its sweeping collection of phone and Internet records on the ground that it is only collecting so-called "metadata"—who you call, when you call, how long you talk. Since this does not

## Cellphone Metadata



<http://www.zeit.de/datenschutz/malte-spitz-data-retention>

## EXIF Metadata



## What Meta-Data Do You Need?

- **Descriptive metadata** describes a resource for purposes such as discovery and identification
- **Administrative metadata** provides information to help manage a resource, such as when and how it was created
- **Rights management metadata**, which deals with intellectual property rights
- **Preservation metadata**, which contains information needed to archive and preserve a resource

# Structured Metadata

## Dublin Core Example

Title="Metadata Demystified"  
Creator="Brand, Amy"  
Creator="Daly, Frank"  
Creator="Meyers, Barbara"  
Subject="metadata"  
Description="Presents an overview of metadata conventions in publishing."  
Publisher="NISO Press"  
Publisher="The Sheridan Press"  
Date="2003-07"  
Type="Text"  
Format="application/pdf"  
Identifier="http://www.niso.org/standards/resources/Metadata\_Demystified.pdf"  
Language="en"

Understanding Metadata: niso.org

# Structured Metadata

```
<mcp:parameterName>
  <mcp:DP_Term>
    TAGS → <mcp:term>
              <gco:CharacterString>tile_lower_right_y_LL</gco:CharacterString>
            </mcp:term>
          <mcp:type>
            <mcp:DP_TypeCode codeList="http://schemas.aodn.org.au/mcp-2.0/
            schema/resources/Codelist/gmxCodeLists.xml#DP_TypeCode"
            codeListValue="shortName">shortName</mcp:DP_TypeCode>
          </mcp:type>
        <mcp:usedInDataset>
          <gco:Boolean>true</gco:Boolean>
        </mcp:usedInDataset>
      </mcp:DP_Term>
    </mcp:parameterName>
```

KNB

## A Simple Typology of Data

- Numerical
  - Continuous
  - Integers
  - Ordinal
- Controlled Vocabulary
  - Certain defined words with defined meanings
  - Has a reference 'dictionary'
- Dates/Times
  - Many formats – POSIX
- Raw Text
- Other Media

## What is a Metadata Standard?

- A Standard provides a structure to describe data with:
  - Common terms to allow consistency between records
  - Common definitions for easier interpretation
  - Common language for ease of communication
  - Common structure to quickly locate information
- In search and retrieval, standards provide:
  - Documentation structure in a reliable and predictable format for computer interpretation
  - A uniform summary description of the dataset



CC Image by carlbeard on Flickr



# Metadata Standards

## Metadata Standards

### ABCD - Access to Biological Collection Data

A standard for the access to and exchange of primary biodiversity data, including specimens and observations.

### Darwin Core

A body of standards, including a glossary of terms (in other contexts these might be called properties, elements, fields, columns, attributes, or concepts) intended to facilitate the sharing of information about biological diversity by providing reference definitions, examples, and commentaries.

### EML - Ecological Metadata Language

Ecological Metadata Language (EML) is a metadata specification particularly developed for the ecology discipline.

### Genome Metadata

Descriptive data about single genomes within the Pathosystems Resource Integration Center.

### ISA-Tab

A general purpose framework with which to capture and communicate metadata for data files from 'omics-based' experiments employing combinations of technologies.

### MIBBI - Minimum Information for Biological and Biomedical Investigations

A common portal to a group of checklists of Minimum Information in nearly 40 biological disciplines.

### Observ-OM

Used to integrate and compare observation data across experimental projects, disease databases, and clinical biobanks.

### OME-XML - Open Microscopy Environment XML

A metadata standard and data file format for biological light microscopy data.

<http://www.dcc.ac.uk/resources/subject-areas/biology>



Britishlibrary.co.uk

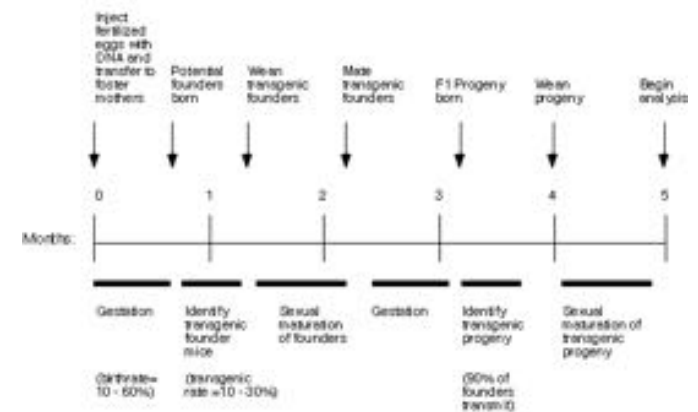
## Case Study 1



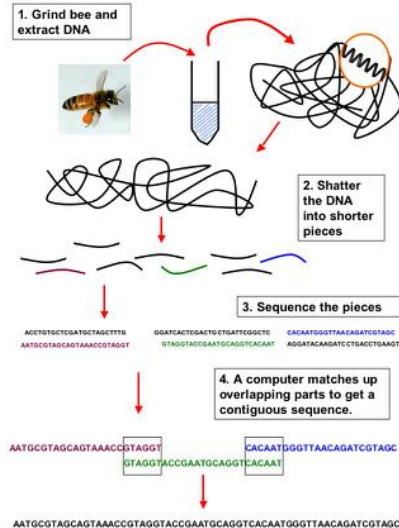
Digital Image © 2002 University of Wisconsin-Madison Libraries

## Case Study 2

### Timeline for Transgenic Mouse Analysis



## Case Study 3



beespotter.org

## Data Collection



## Creating a Good Data Gathering Sheet

- How easy is it to read?
- Are column and row definitions clear?
- Is there metadata?
- How similar is it to your digital data entry form?
- Can you use it at 4am?

## After the Collection...

- Preserve original data
- Created digital archive of raw data
- Implement robust storage strategy
- Quality Control

## Scanning

KEEN GOM BENTHIC MONITORING DATA BENTHIC UNIFORM POINT CONTACT DATA SHEET Scanned by: \_\_\_\_\_

Site: SW Appleton Transect: N. Raposa Date: 2/4/14 Observer: JD

Depth (m) at 0m: 0.2m Depth (m) at 10m: 2.3 Velocity (m/s): \_\_\_\_\_

Inshore	Substrate	Offshore	Substrate
HJ B	B	HJ VEL EC	B
HJ	SS	HJ SL & DIAT	GR
HJ	SL	HJ CF	B
SL VEL	B	SL SL HJ	B
SL HJ BACF	BL	HJ ATM	BL
HJ CF	BL	HJ ATM	BL
HJ EL	BL	HJ DIAT	BL
HJ CF	BL	HJ SL CF	BL
HJ CF	BL	HJ EC	BL
HJ ATM	BL	SL HJ ANSP	BL
HJ CF	B	HJ SL ANSP	B
HJ CF	B	HJ	SS
HJ CF	B	HJ ANSP	B
HJ CF	B	HJ HIRU	BL
HJ SL CF	B	DEVI HJ BUTU ATM	BL

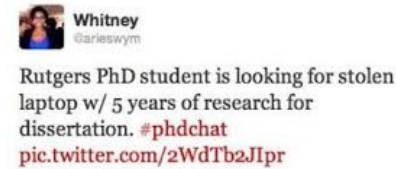
**Green Algae**

- CP Codium fragile (fragile)
- FG Filamentous Green (fine hairs)
- GV Grassy Ulva (slab & enteromorpha)
- TV Tubular Ulva (slab & enteromorpha)

**Red Algae**

- BOHA Boreosiphonia hamifera (hooks)
- CRSP Ceramium spp. (circular, vorticifer)
- CHCR Chondrus crispus
- CYPO Cylindrocapsa (bottle)
- DUCO Dumortiera contorta (flat cylinders)
- EUCR Euthalia crassa (flat branching blades, leafy)
- HJ Heterosiphonia japonica (bright, tuft, solid mass)
- MAST Mastocarpus stellatus (thamnidial blades, pebbles)
- PAPA Palmaria palmata (hand-like blades)
- PHRU Phycodryas rubens (dark leaves)
- POLS Polysiphonia sp. (variable tufts)
- PORD Pordetia rotunda (regular thick branches)
- PORE Porphyra spp. (thin sheet)
- PTSE Ptilota serrata (branches w/ branchlets)
- RAT Red Alga Turf (thin UMD mat, many spp.)
- SPHE Sphaerostichum repens (thick, tuft, cylindrical)
- CO Corallina officinalis (erect coralline)
- URS Unidentified Red Blade
- UPR Unidentified Filamentous Red
- SPRSP Sphaerostichum repens (thick, tuft, cylindrical)
- HLR Heterosiphonia rubra (not called oval)
- CLSP Cladophora spp. (smooth and thick)
- LESP Leptocarpus spp. (smooth, very thin, no band)
- PHSP Phyllocladus spp. (rough surface, white center)
- LGL Lithothamnion glaciale (bottle)
- UP Unidentified Filamentous (bottle)

## Storage: Physical



DO NOT LET THIS BE YOU

## Storage: Physical



## Storage: The Cloud





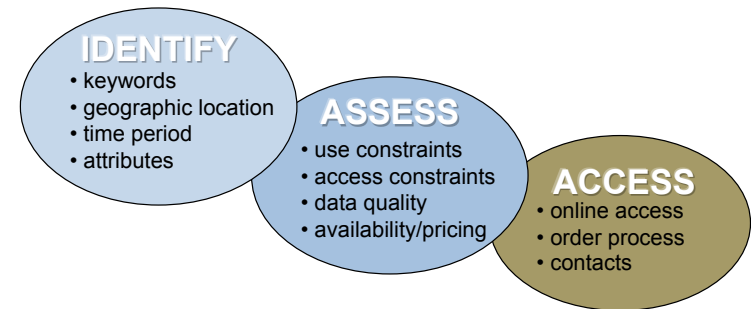
## Data Sharing



blog.veritythink.com

## Distribution: Data Discovery

- The descriptive content of the metadata file can be used to identify, assess, and access available data resources.



## Things to Consider when Data Sharing

1. Is what you did understandable?
2. How do you want your work credited?
3. Will your data sharing service be around in 50 years?

## Why Share Data?

- One scientist can only do so much
  - More data = more Power
- Science must be reproducible
- Who paid for this data collection?

## Examples

- <https://www.dataone.org/>
- <http://blast.ncbi.nlm.nih.gov/Blast.cgi>
- <http://datadryad.org/>
- <http://www.oceandataportal.org/>

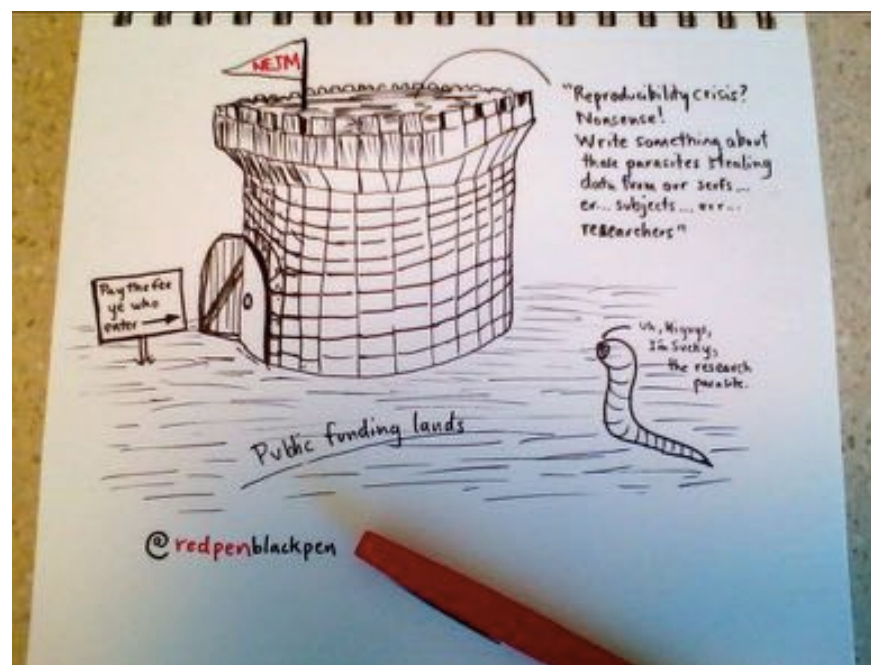
## Backlash?



The NEW ENGLAND

"A second concern held by some is that a new class of research person will emerge — people who had nothing to do with the design and execution of the study but use another group's data for their own ends, possibly stealing from the research productivity planned by the data gatherers, or even use the data to try to disprove what the original investigators had posited. There is concern among some front-line researchers that the system will be taken over by what some researchers have characterized as 'research parasites.'

quality information carefully reexamined for the possibility that new nuggets of useful data are lying there, previously unseen? The potential for leveraging existing results for even more benefit pays appropriate increased tribute to the patients who put themselves at risk to generate the data. The moral imperative to honor their collective sacrifice is the trump card that takes this trick.



## Should all Data be Open?



## Let's Look at Data & Metadata Examples

<https://biol355.github.io/datasets.html>