Data Creation

# Chain of Data Creation

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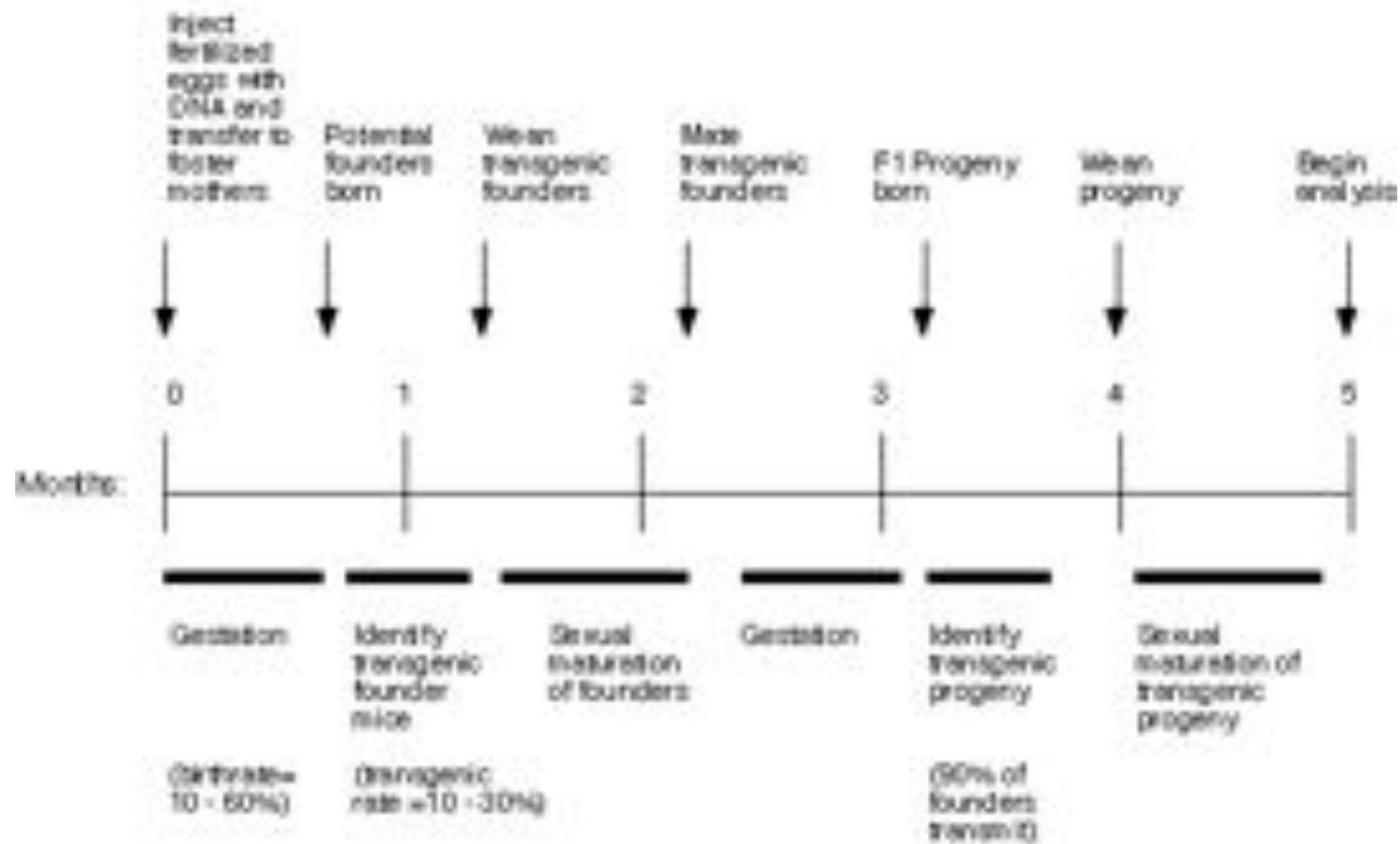7. Data Sharing

# Lab Notebook

book and page numbers
make indexing your work easier,
just enter the page title and number
in the table of contents

pages that are sewn
together are
tamper evident

initial and date each
insert both on and over
the edge of the insert
to discourage removal

sign and date each entry
using a consistent format
and legible writing for each date,
also have each entry signed
and dated by a witness

- Record of hypotheses
- Record of Protocols
- Second brain

# Plan Your Experiment, Experiment With your Plan

## Timeline for Transgenic Mouse Analysis
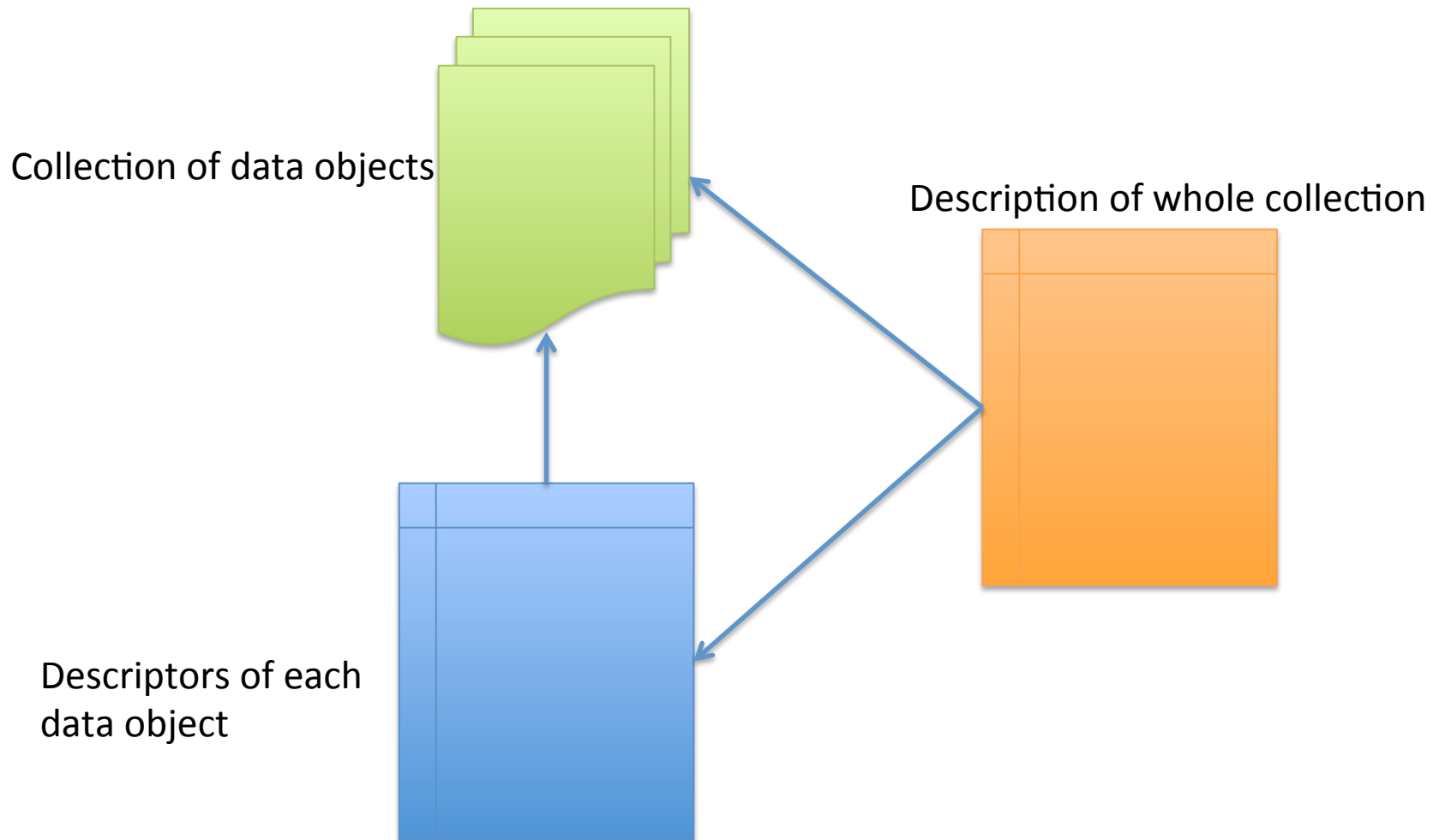


medicine.umich.edu

# 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?

Collection of data objects

Description of whole collection

Descriptors of each
data object

# 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

Understanding Metadata: niso.org

# Structured Metadata

<u>Dublin Core Example</u>

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"

# Structured Metadata

**TAGS** →

```
<eml>
  <access
    authSystem="ldap://ldap.ecoinformatics.org:389/dc=ecoinformatics,dc=org"
    order="allowFirst">
    <allow>
      <principal>uid=alice,o=NASA,dc=ecoinformatics,dc=org</principal>
      <permission>read</permission>
      <permission>write</permission>
    <allow>
  </access>
  <dataset>
  ...
  ...
  <dataTable id="entity123">
  ...
    <physical>
    ...
      <distribution>
      ...
        <access id="access123"
        authSystem="ldap://ldap.ecoinformatics.org:389/dc=ecoinformatics,dc=org"
        order="allowFirst">
          <deny>
            <principal>uid=alice,o=NASA,dc=ecoinformatics,dc=org</principal>
            <permission>write</permission>
          </deny>
        </access>
      </distribution>
    </physical>
  </dataTable>
  <dataTable id="entity234">
    ...
    <physical>
    ...
      <distribution>
        ...
        <access>
          <references>access123</references>
```
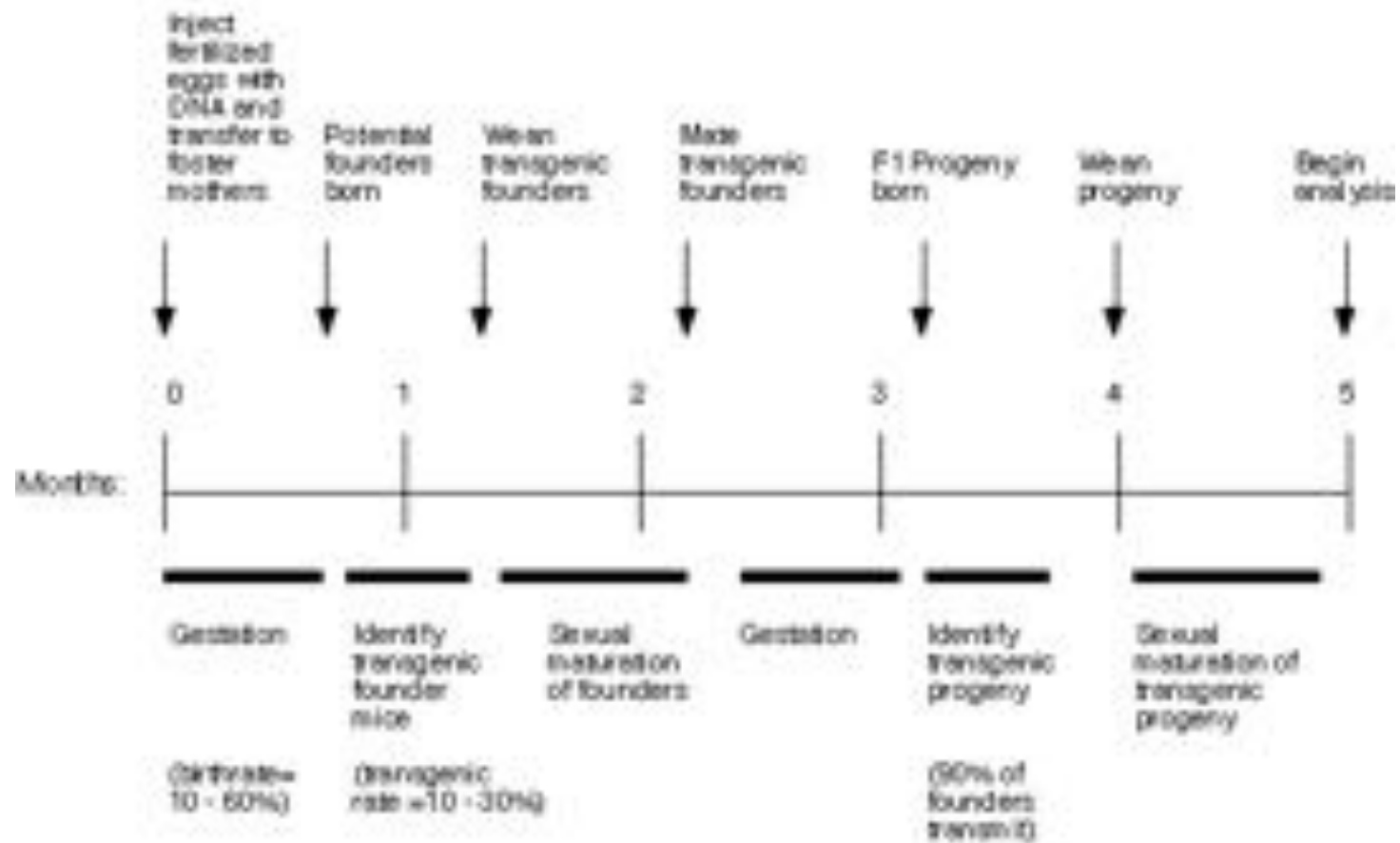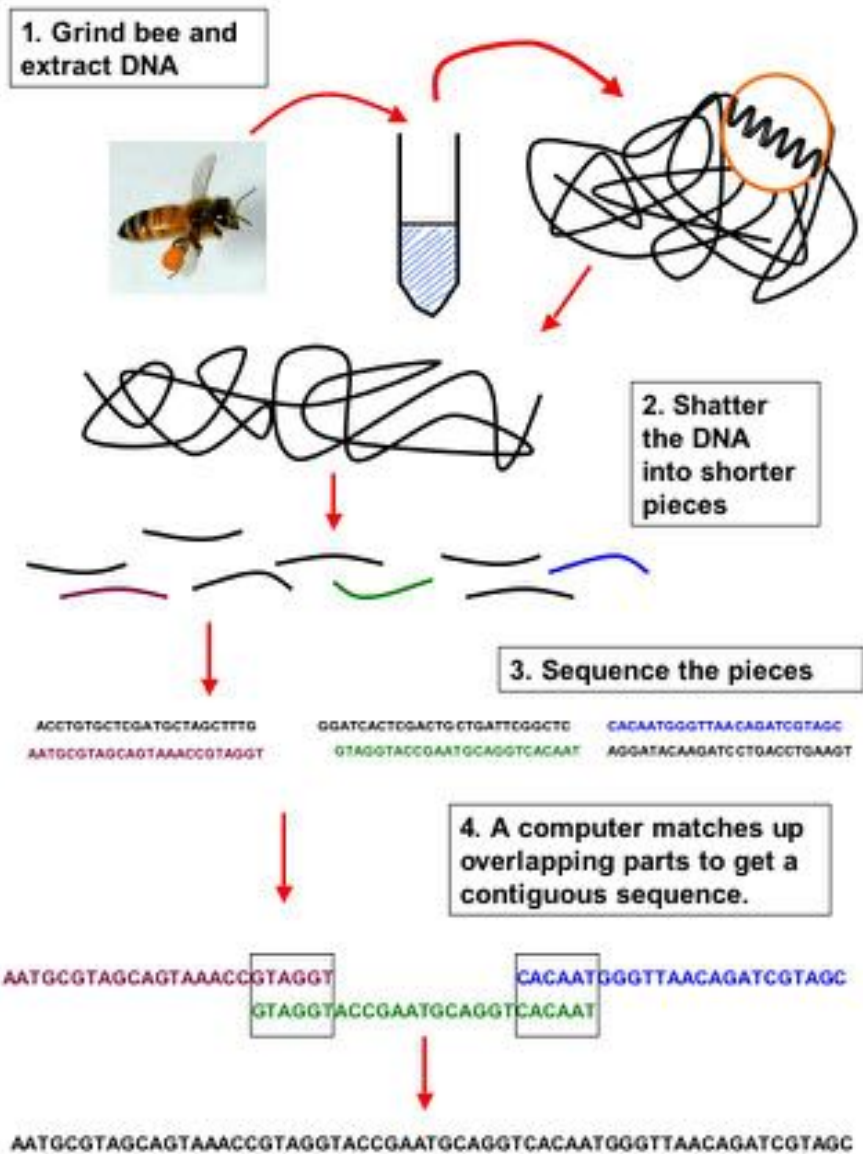
# Case Study 1

# Case Study 2



Timeline for Transgenic Mouse Analysis

# Case Study 3



1. Grind bee and extract DNA

2. Shatter the DNA into shorter pieces

3. Sequence the pieces

ACCTGTGCTCGATGCTAGCTTTG
AATGCGTAGCAGTAAACCGTAGGT

GGATCACTCGACTGCTGATTCGGCTC
GTAGGTACCGAATGCAGGTCACAAT

CACAATGGGTTAACAGATCGTAGC
AGGATACAAGATCCTGACCTGAAGT

4. A computer matches up overlapping parts to get a contiguous sequence.

AATGCGTAGCAGTAAACCGTAGGT
           GTAGGTACCGAATGCAGGTCACAAT

CACAATGGGTTAACAGATCGTAGC

AATGCGTAGCAGTAAACCGTAGGTACCGAATGCAGGTCACAATGGGTTAACAGATCGTAGC

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 (next time)

# Scanning

# EXCEL TIME!

Entry

Fills

Basic Functions

Functions for Error Checking

Controlled Vocabularies

# Storage: Physical



**Whitney**
@arieswym

Rutgers PhD student is looking for stolen laptop w/ 5 years of research for dissertation. #phdchat
pic.twitter.com/2WdTb2JIpr

DO NOT LET THIS BE YOU

# Storage: Physical

# Storage: The Cloud

# Data Sharing

# 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/](https://www.dataone.org/)

- [http://blast.ncbi.nlm.nih.gov/Blast.cgi](http://blast.ncbi.nlm.nih.gov/Blast.cgi)

- [http://datadryad.org/](http://datadryad.org/)

- [http://www.oceandataportal.org/](http://www.oceandataportal.org/)
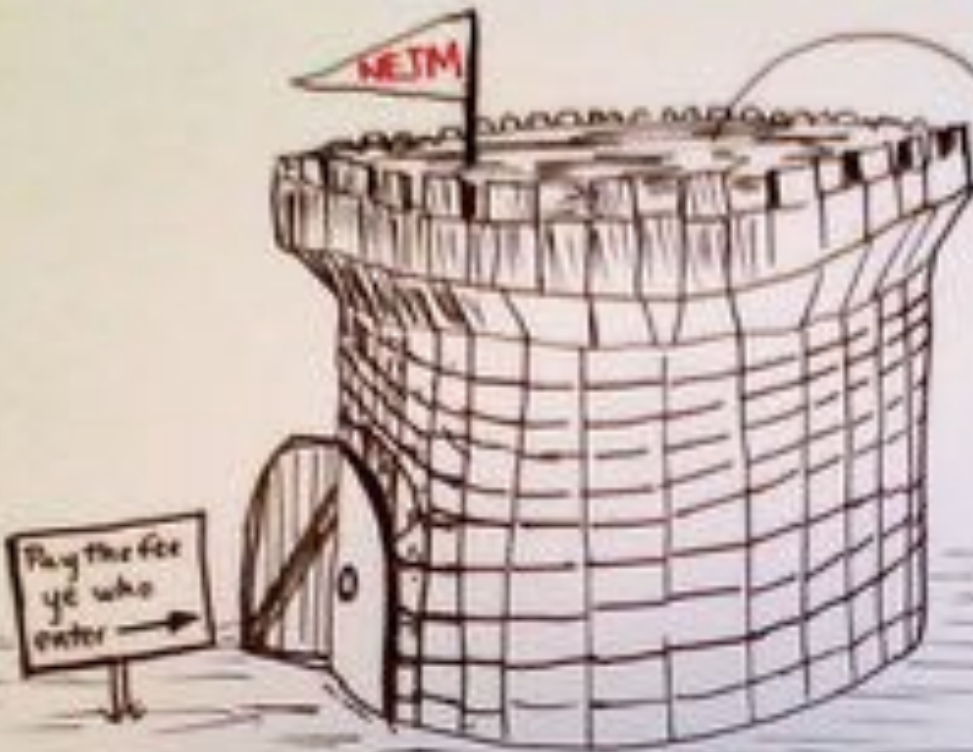
# Backlash?

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?