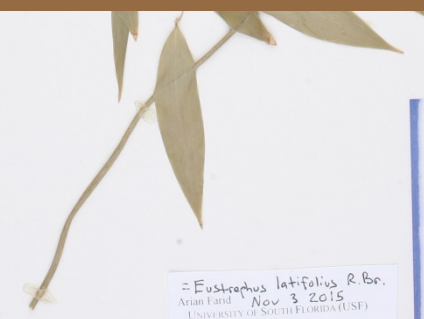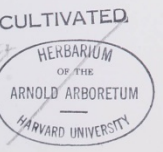# Incorporating collector behavior into large-scale range models for digital biodiversity data

**Kelley D. Erickson**, Stephen J. Murphy and Adam B. Smith

MISSOURI BOTANICAL GARDEN

INSTITUTE of Museum and Library SERVICES

# Rapidly expanding access to an enormous amount of digital biodiversity data



Get data    Share    Tools    Inside GBIF

GBIF | Global Biodiversity Information Facility

Free and open access to biodiversity data

**iNaturalist**    Explore  Community ⌄  More ⌄

Observations

Species            Location

The World    **21,179,569** OBSERVATIONS    **215,638** SPECIES    **80,829** IDENTIFIERS

Map  Grid  List  Places of Interest

Cliff Swallow
(Petrochelidon)
Whatcom Cou
• Jun 2, 2019

Unknown
Unknown • Ju

Mississippi
Cittina missis
Tallahassee, FL

Green Iguan
(Iguana iguana)
Provincia de Al

Comocladia platyphylla observed in Rafael Fre

| Occurrence records | Datasets | Publishing institutions | Peer-reviewed papers using data |
|---|---|---|---|
| 1,304,475,217 | 44,934 | 1,409 | 3,697 |

News

Angola becomes the newest member of the GBIF network
20 May 2019

Data use

On the evolution of food customs
4 June 2019

2019 Ebbe Nielsen Challenge

News

2019 GBIF Ebbe Nielsen Challenge seeks open-data innovations for biodiversity
Deadline: 1 August 2019

News

Data mobilization and capacity building essential to address global biodiversity crisis
6 May 2019

**Tropicos®**

Home  Names  Specimens  References  Projects  Images  More⌄  Tools⌄

Tropicos® was originally created for internal research but has since been n
available to the world's scientific community. All of the nomenclatural, bibl
and specimen data accumulated in MBG's electronic databases during the
years are publicly available here. This system has nearly 1.3 million scienti
and over 4.4 million specimen records.

Quick Name Search          Search    Search E

Common Name

**iDigBio**    About iDigBio | Research
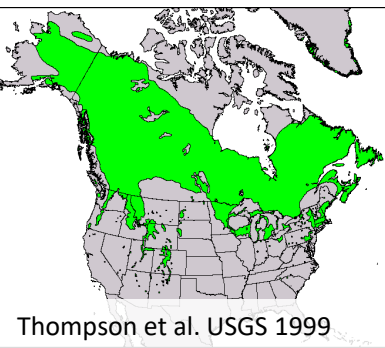Integrated Digitized Biocollections

Google Cus

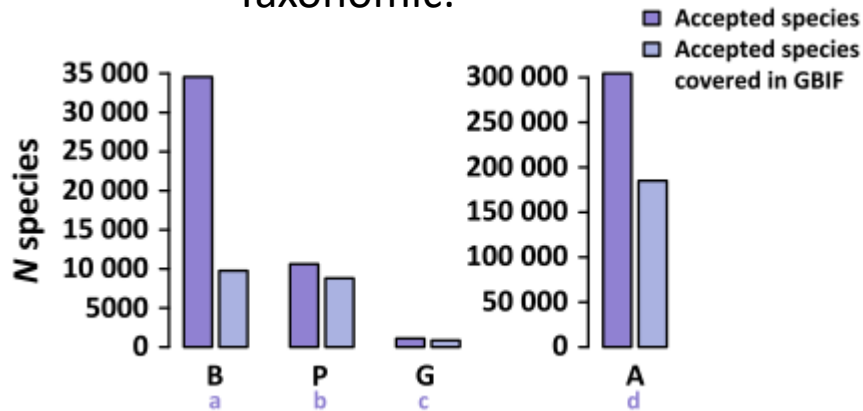Making data and images of millions of biological specimens available on the web

119,163,881
Specimen Records

30,380,997
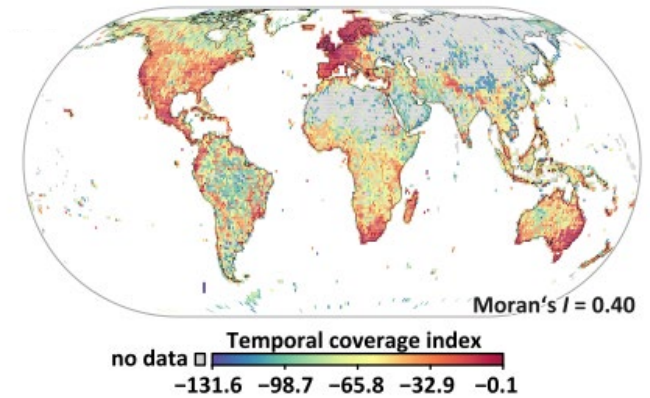Media Records

1,614
Recordsets

Search the Portal

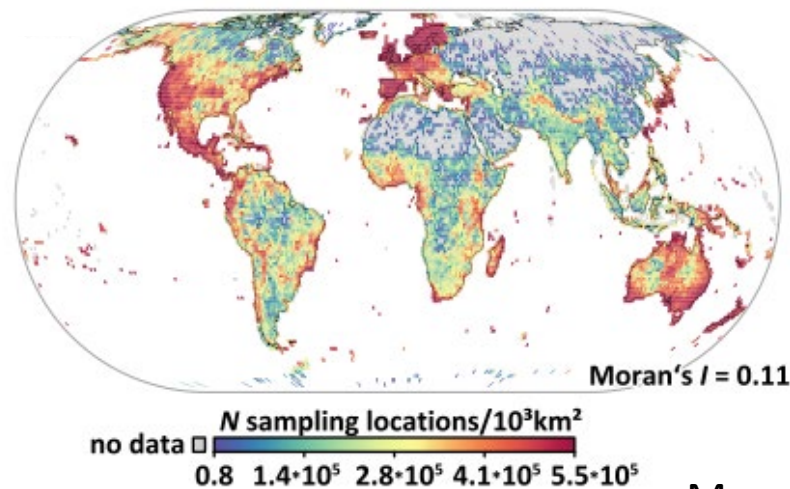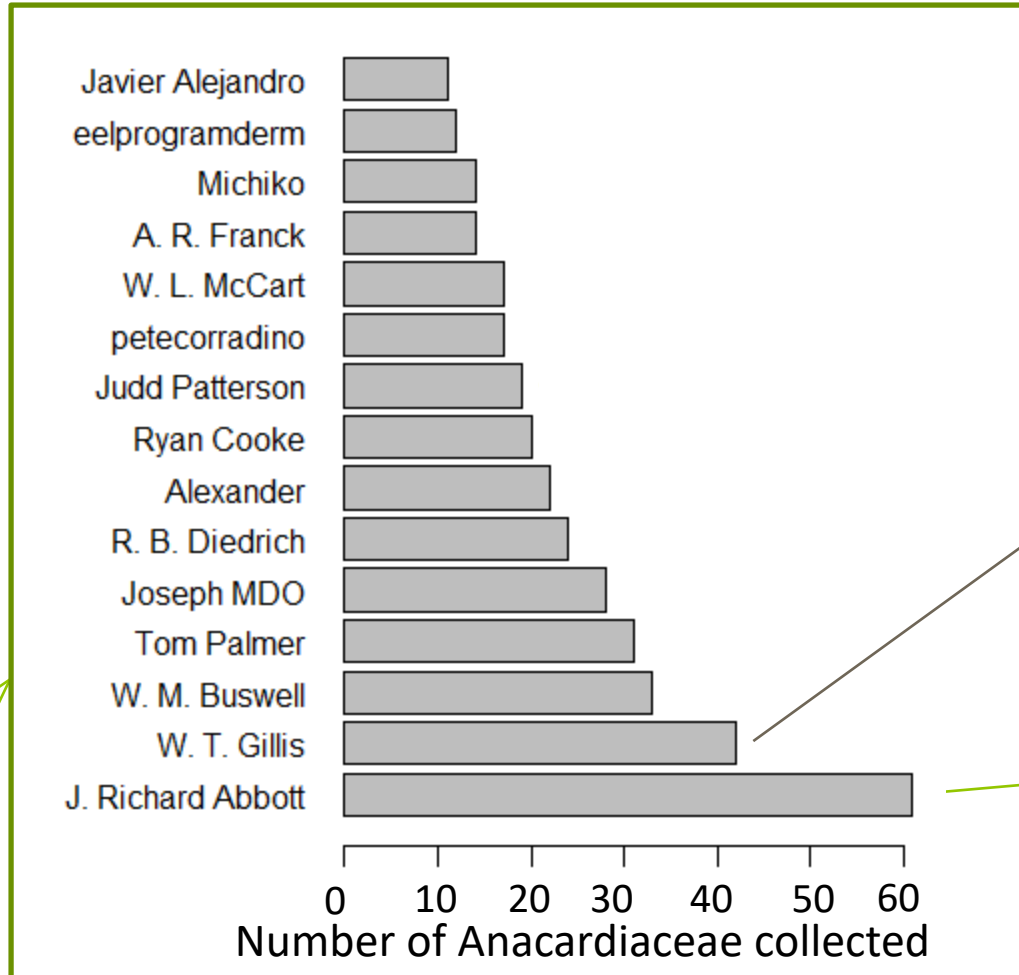# Gaps and biases in occurrence data for species ranges


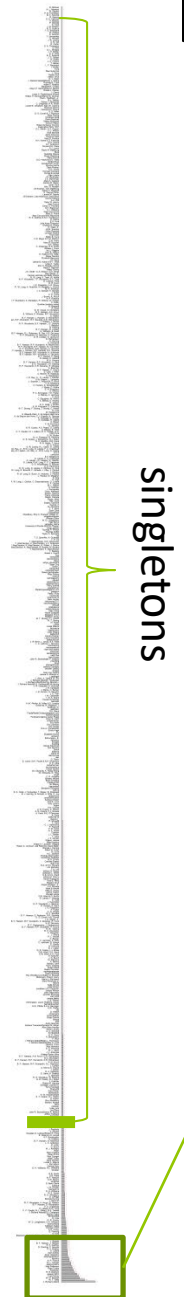Thompson et al. USGS 1999

Taxonomic:



Temporal:



Geographical:



Meyer *et al. Ecology Letters* 2016
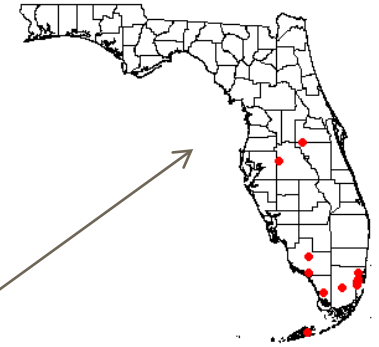
# Most collectors collect only once
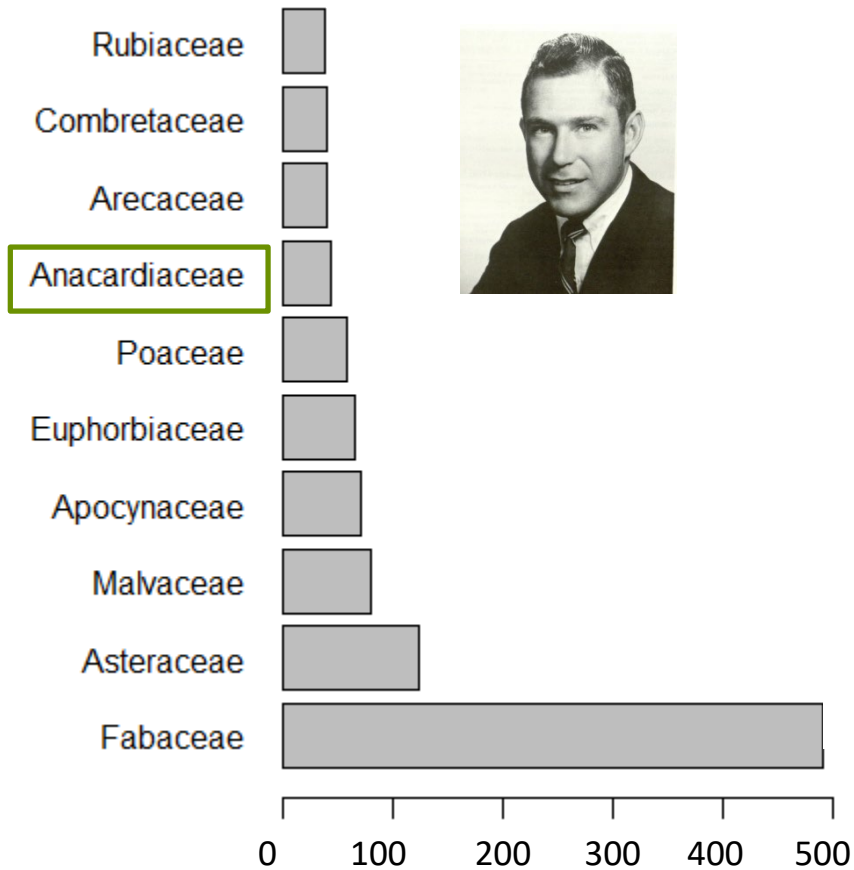


singletons

Number of Anacardiaceae collected
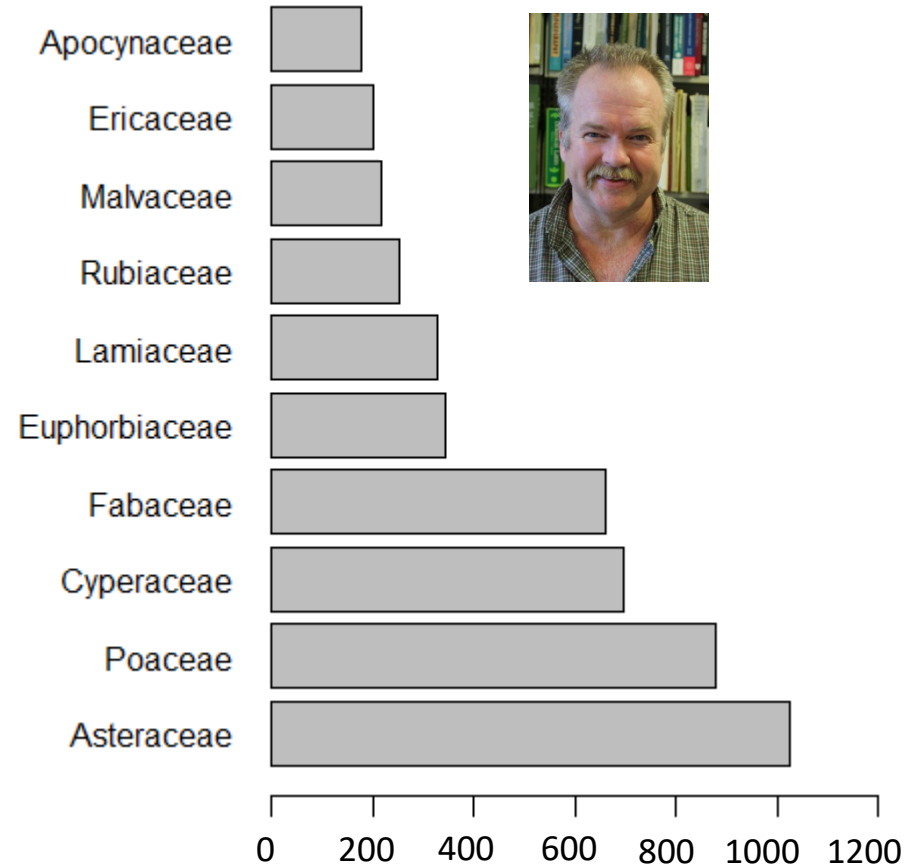
Even those that collect a lot don't collect everywhere:

# Taxonomic focus varies among collectors: Number of Records by Family

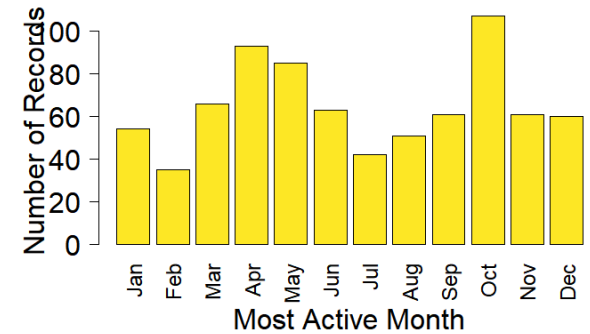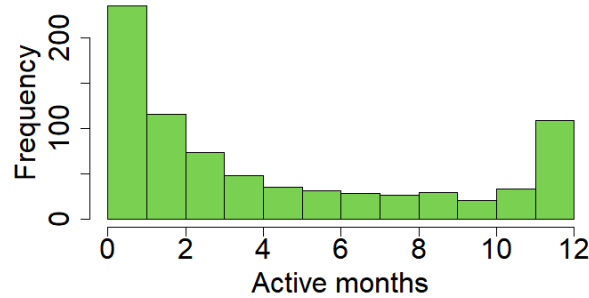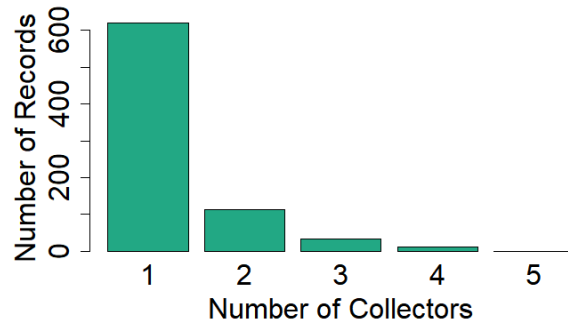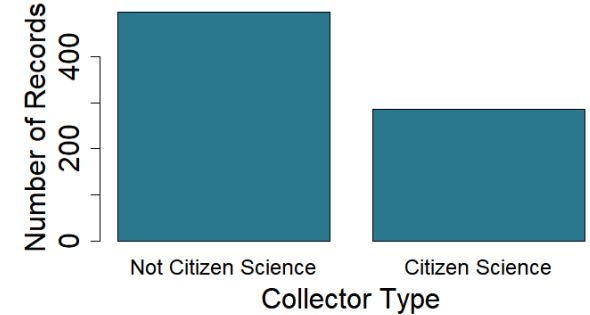# Other ways collectors differ from each other

# Other ways collectors differ from each other



Differences among collectors give rise to observed biases in occurrence data

# Occupancy-detection models

$Z_A = 1$

**Ecological Process**

$$Z_i \sim Bernoulli\ (\psi_i)$$

$Z_B = 1$

$Z_C = 1$

$Z_D = 0$

$Z_E = 1$

# Occupancy-detection models



Site / Survey

*Latent State*

$Z_A = 1$

$Z_B = 1$

$Z_C = 1$

$Z_D = 0$

$Z_E = 1$

Ecological Process

$$Z_i \sim Bernoulli\ (\psi_i)$$

MacKenzie et al. *Ecology* 2002

# Occupancy-detection models



MacKenzie et al. *Ecology* 2002

# Occupancy-detection models



MacKenzie et al. *Ecology* 2002

# Occupancy-detection models

**Observed States**

**Latent State**

| Site | 1 | 2 | 3 | 4 | |
|------|---|---|---|---|---|
| A | 1 | 1 | | 1 | $Z_A = 1$ |
| B | 0 | | 1 | 0 | $Z_B = 1$ |
| C | 0 | 0 | 1 | | $Z_C = 1$ |
| D | 0 | | 0 | 0 | $Z_D = 0$ |
| E | 0 | 0 | 0 | 0 | $Z_E = 1$ |

**Survey**

**Ecological Process**

$$Z_{i,t} \sim Bernoulli\ (\psi_i)$$

**Observation Process**

$$Y_{i,t} \mid Z_i \sim Bernoulli\ (p_i \cdot Z_i)$$

$$\text{logit}(p_i) = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

Detectability can depend on characteristics of collectors

MacKenzie et al. *Ecology* 2002

# Translating occupancy-detection framework to collections data context

- How do we define a **site**?

  o a county or other defined locality
  o buffer around a given coordinate

# Translating occupancy-detection framework to collections data context

- How do we define a **site**?

  

  - a county or other defined locality
  - buffer around a given coordinate

- What is the time scale of the **sampling period?**

  

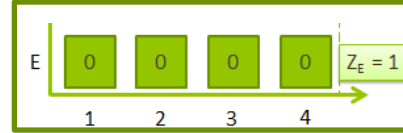  - ***closure assumption:*** *occupancy status of site does not change between sampling periods*
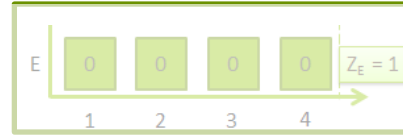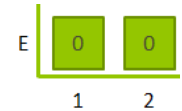
# Translating occupancy-detection framework to collections data context

- How do we define a **site**?

  - a county or other defined locality
  - buffer around a given coordinate

- What is the time scale of the **sampling period?**

  - *closure assumption: occupancy status of site does not change between sampling periods*

- What counts as a single **survey**?

  - all collections by a **single individual** (or collection group) within a **site** during a single **sampling period**

# Translating occupancy-detection framework to collections data context

- How do we define a **site**?

  

  - a county or other defined locality
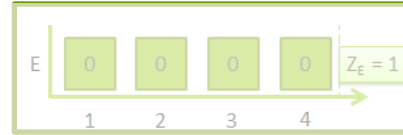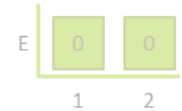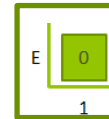  - buffer around a given coordinate

- What is the time scale of the **sampling period?**

  

  - *closure assumption: occupancy status of site does not change between sampling periods*
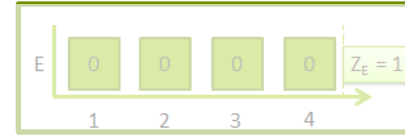
- What counts as a single **survey**?

  

  - all collections by a **single individual** (or collection group) within a **site** during a single **sampling period**

- What is a **non-detection event?**

  

  - **Surveys** by collectors that do not include a single record of the focal species
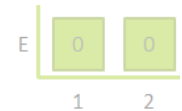
# Translating occupancy-detection framework to collections data context

- How do we define a **site**?

  

  - a county or other defined locality
  - buffer around a given coordinate

- What is the time scale of the **sampling period?**

  

  - ***closure assumption:*** *occupancy status of site does not change between sampling periods*

- What counts as a single **survey**?

  

  - all collections by a **single individual** (or collection group) within a **site** during a single **sampling period**
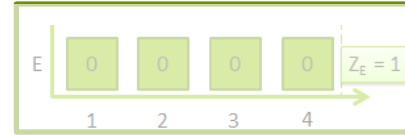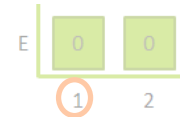
- What is a **non-detection event?**

  

  - **Surveys** by collectors that do not include a single record of the focal species
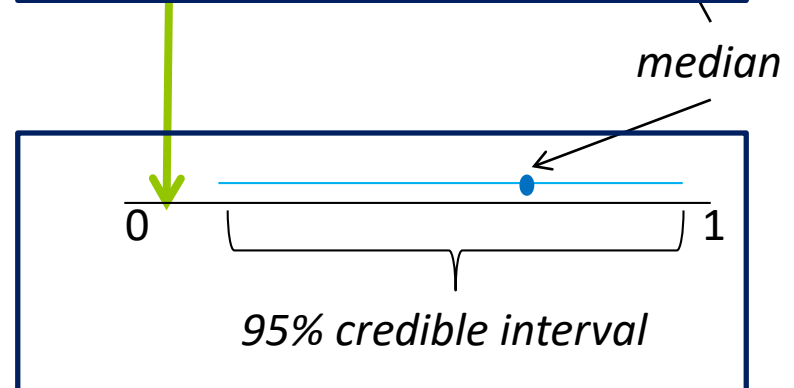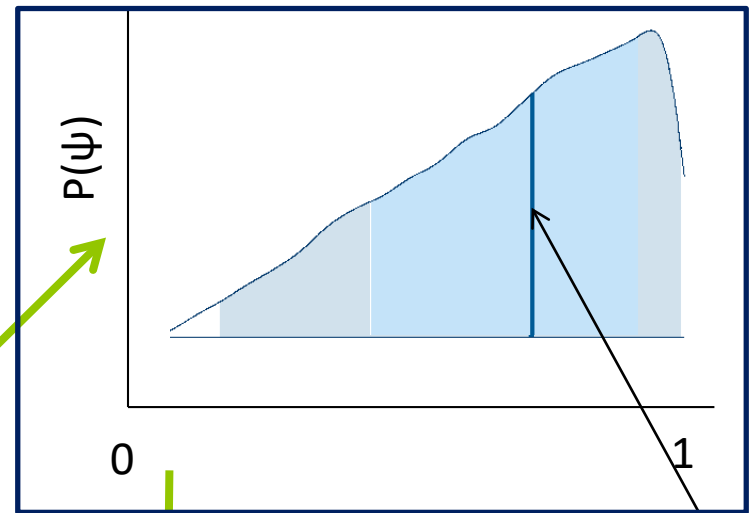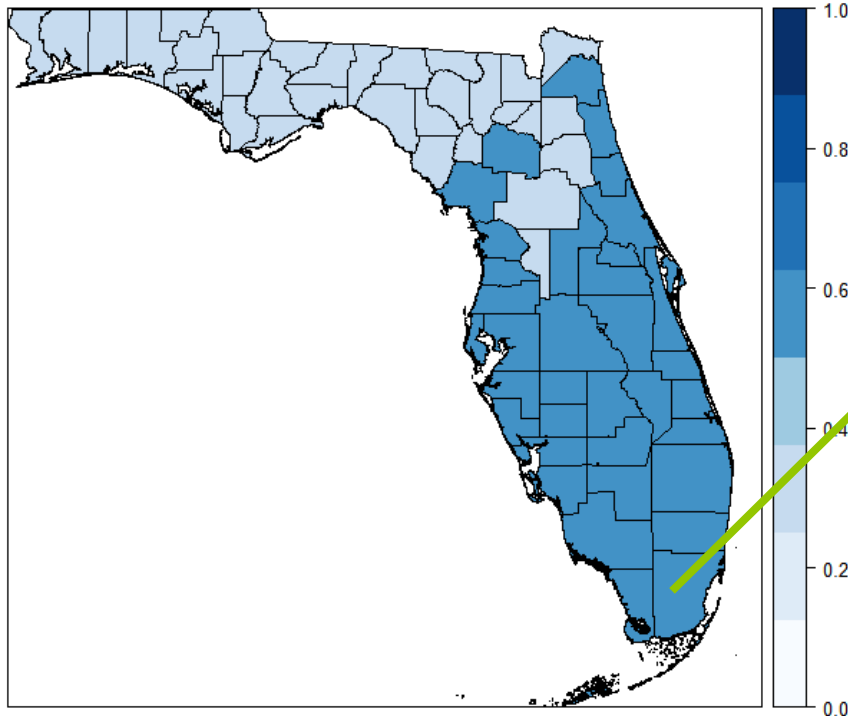
**Answers to these questions are study specific!**

# Case study: Distribution of *Schinus terebinthifolia* in Florida

# Simple occupancy-detection model: Occupancy

$$\text{logit}(p_i) = \varepsilon_i \quad \longleftarrow \quad \textsf{random effect of county}$$

Median of posterior distribution of ψ

Posterior Distribution of Occupancy for Miami-Dade County
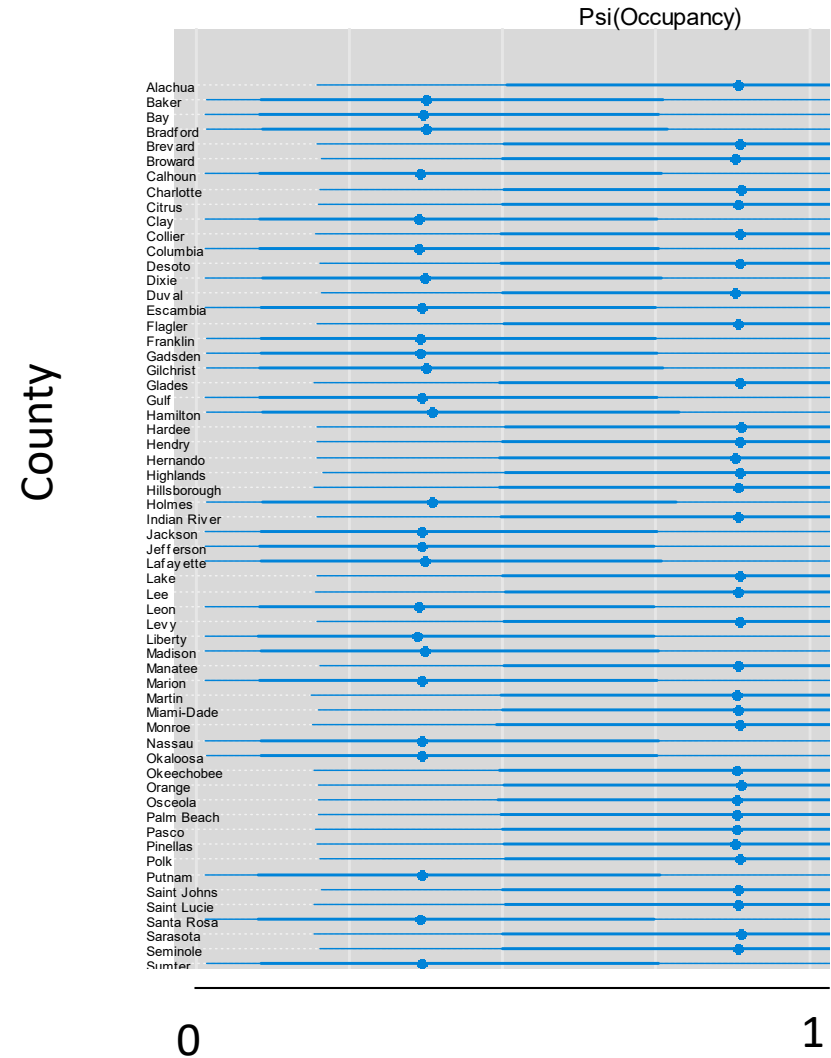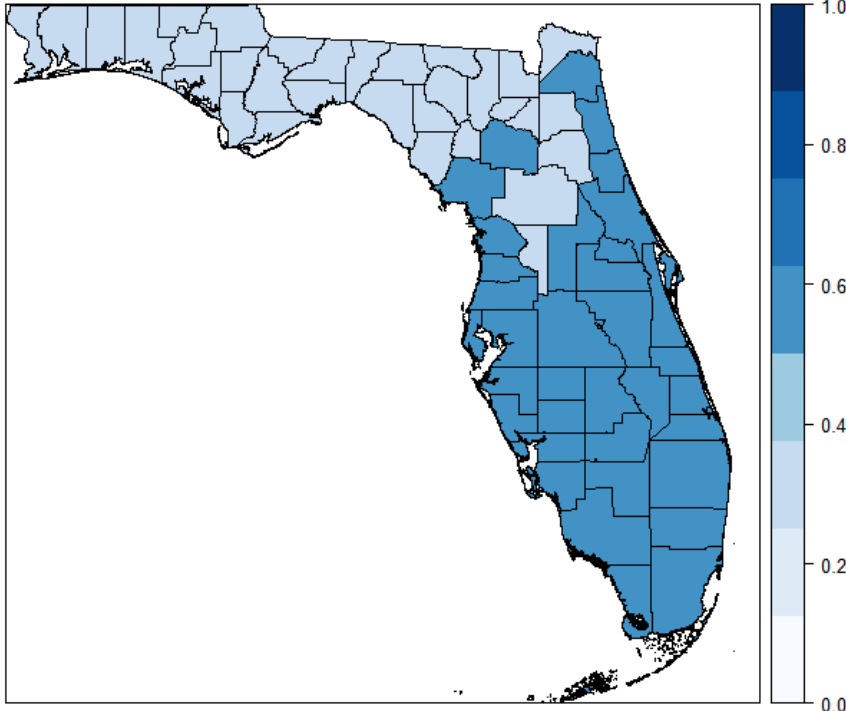


median

95% credible interval

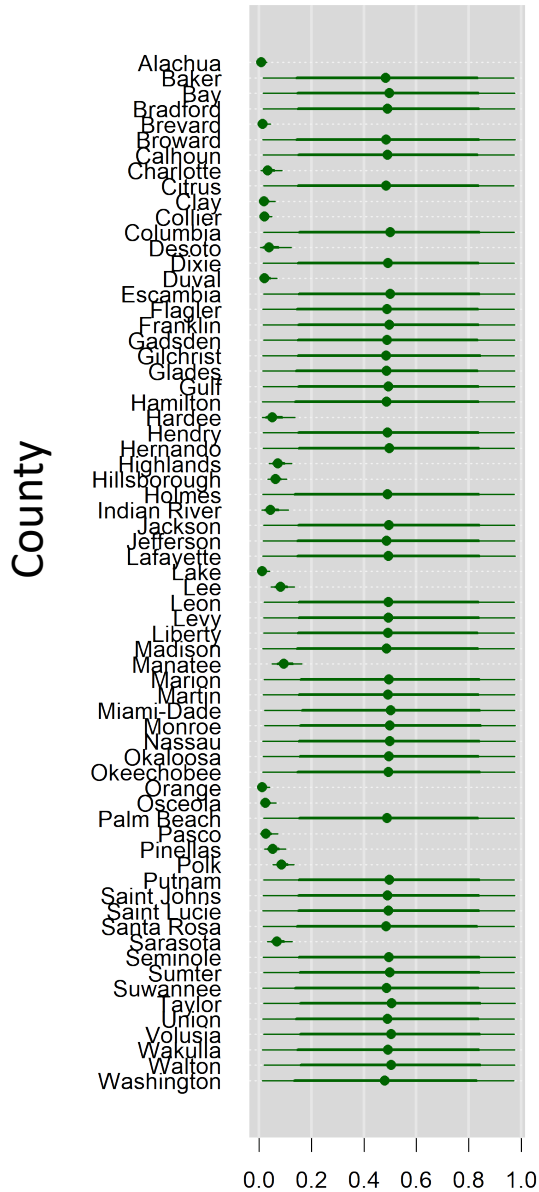Large uncertainty in estimates of occupancy

# Large uncertainty in estimates of occupancy

$$\text{logit}(p_i) = \varepsilon_i \quad \longleftarrow \quad \text{random effect of county}$$

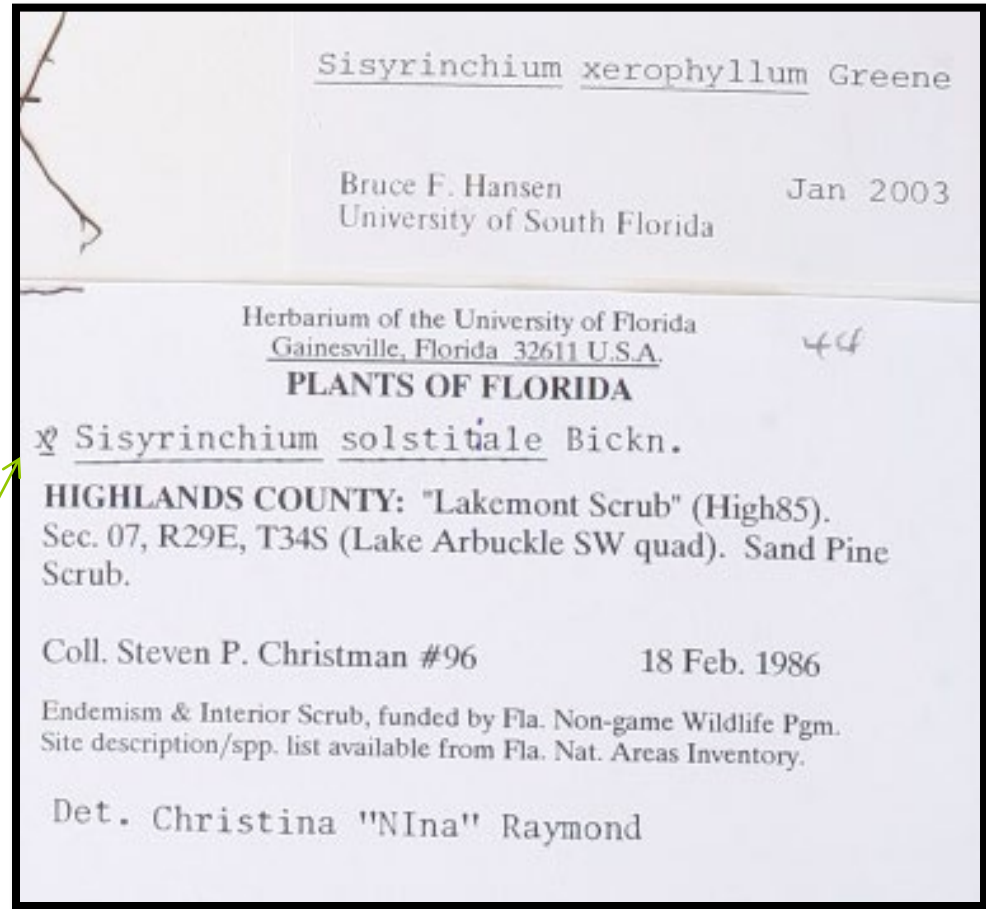Median of posterior distribution of ψ
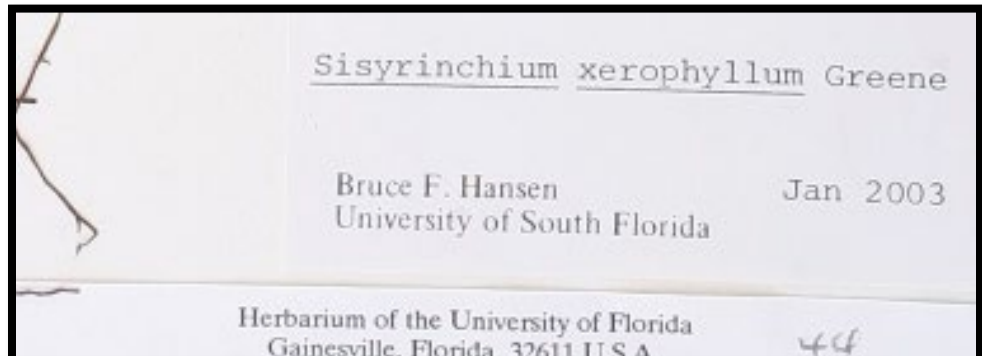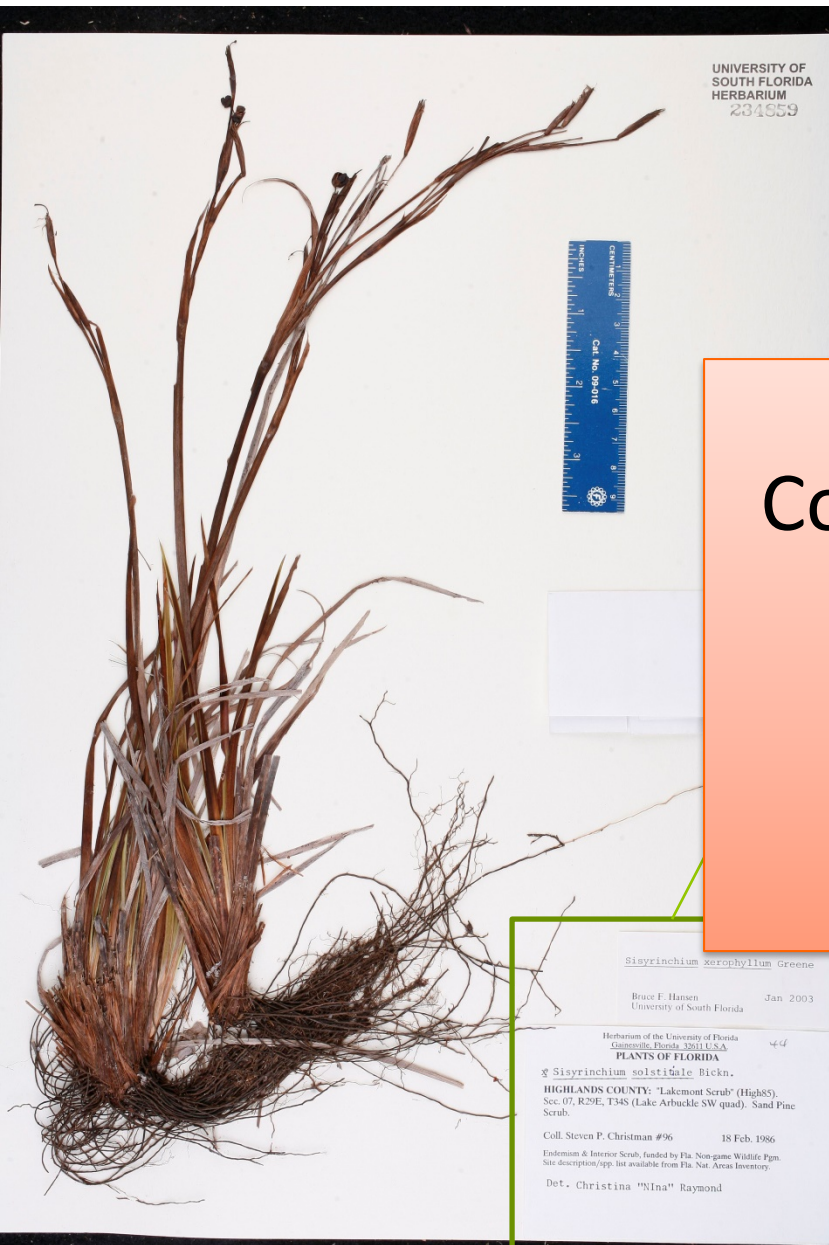
# Simple occupancy-detection model: Detectability



Large uncertainty in detectability when ignore collector behavior

# Collector data is messy



Sisyrinchium xerophyllum Greene

Bruce F. Hansen          Jan 2003
University of South Florida

Herbarium of the University of Florida
Gainesville, Florida 32611 U.S.A.          44
**PLANTS OF FLORIDA**

☿ Sisyrinchium solstitiale Bickn.

**HIGHLANDS COUNTY:** "Lakemont Scrub" (High85).
Sec. 07, R29E, T34S (Lake Arbuckle SW quad). Sand Pine
Scrub.

Coll. Steven P. Christman #96          18 Feb. 1986

Endemism & Interior Scrub, funded by Fla. Non-game Wildlife Pgm.
Site description/spp. list available from Fla. Nat. Areas Inventory.
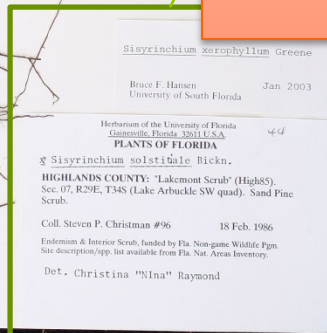
Det. Christina "NIna" Raymond

L. E. Arnold & Erdman West
Erdman West & Lillian Arnold
L. Arnold & Erdman West
L. E. Arnold, Erdman West
Arnold & West
L. E. Arnold; Erdman West
E. West, L. Arnold
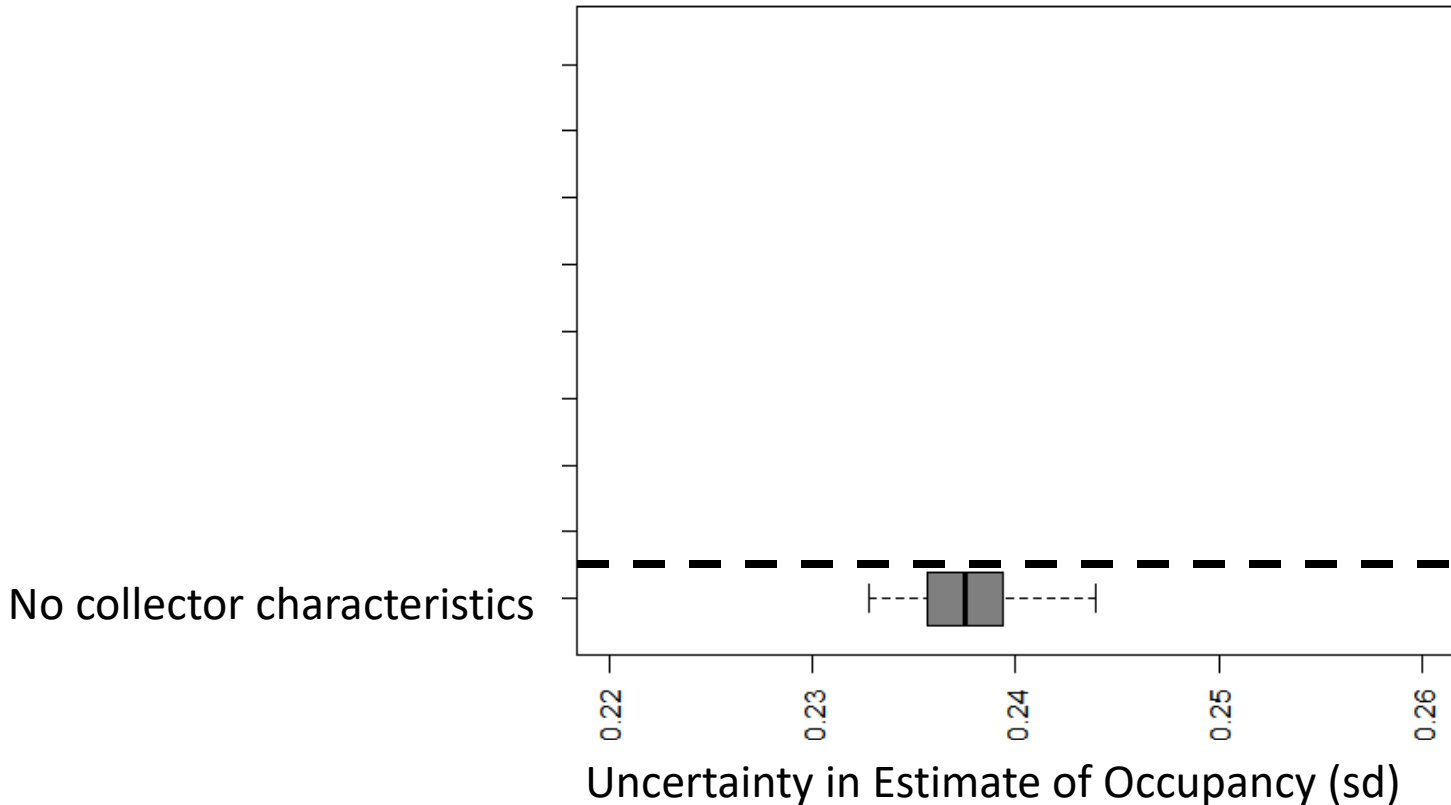
# Collector data is messy

Constructing new R package
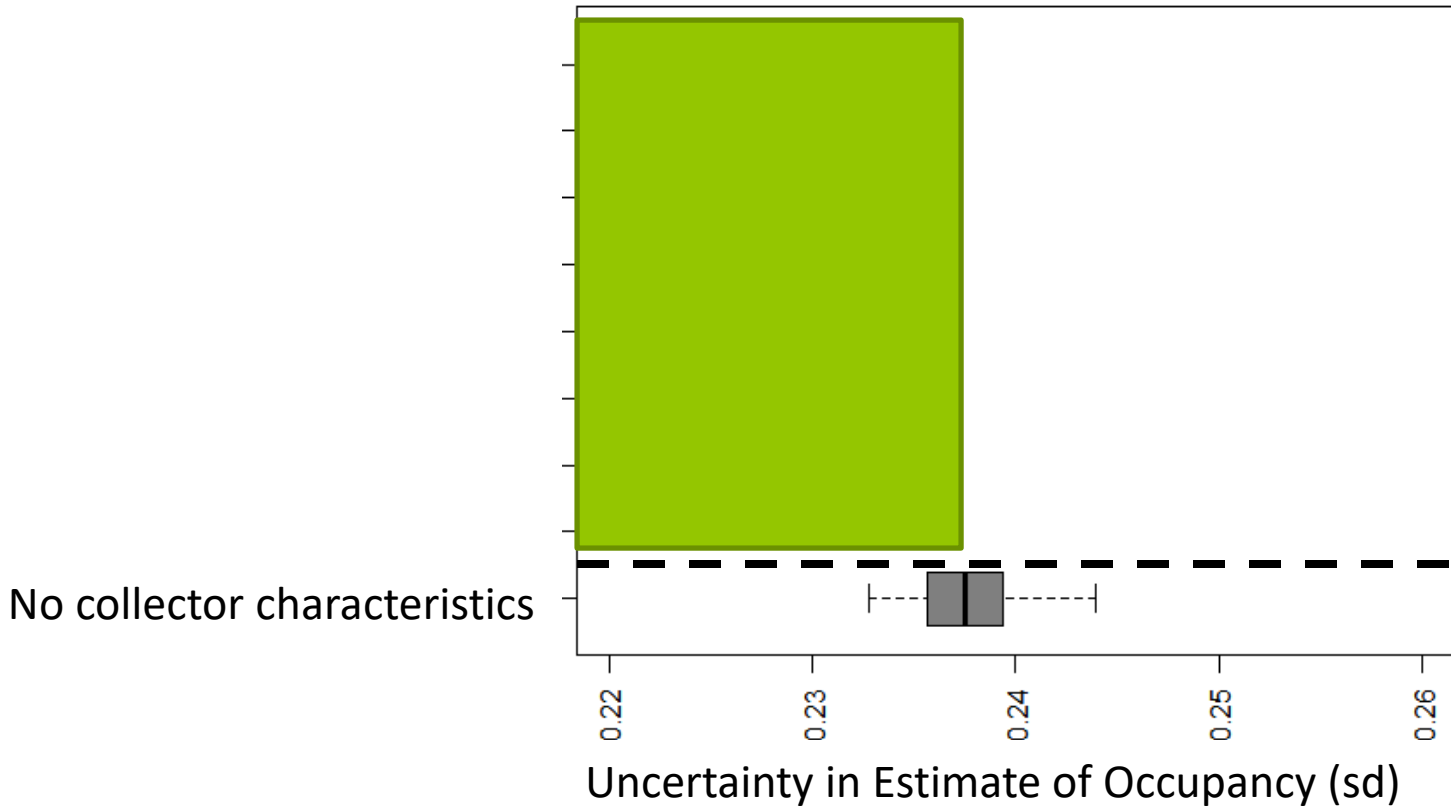collectR!
-automate cleaning of
specimen data

L. E. Arnold & Erdman West
Erdman West & Lillian Arnold
L. Arnold & Erdman West
L. E. Arnold, Erdman West
Arnold & West
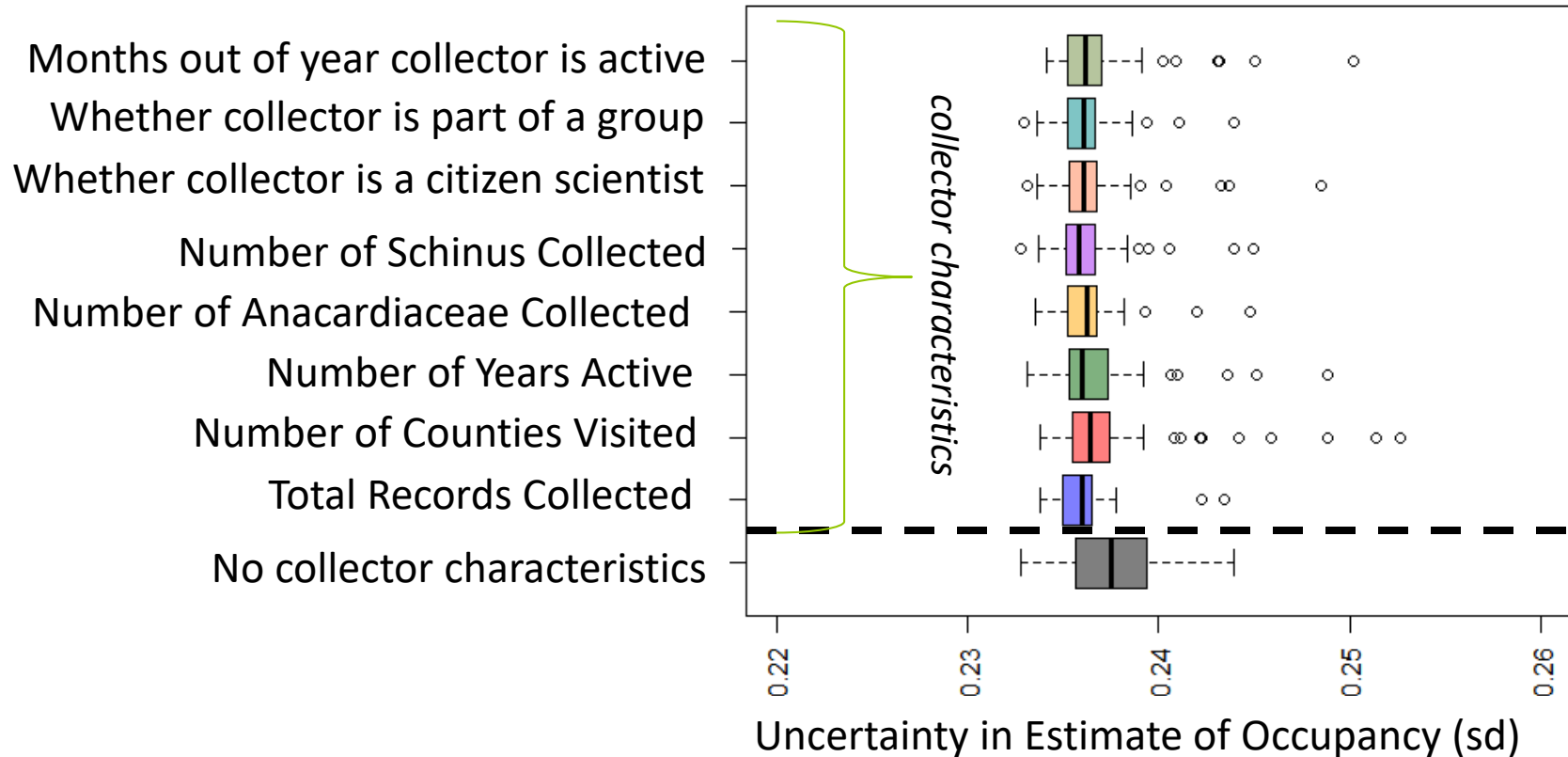L. E. Arnold; Erdman West
E. West, L. Arnold

# Uncertainty in estimate of occupancy for model without collector behavior

# Goal: Shrink uncertainty



No collector characteristics
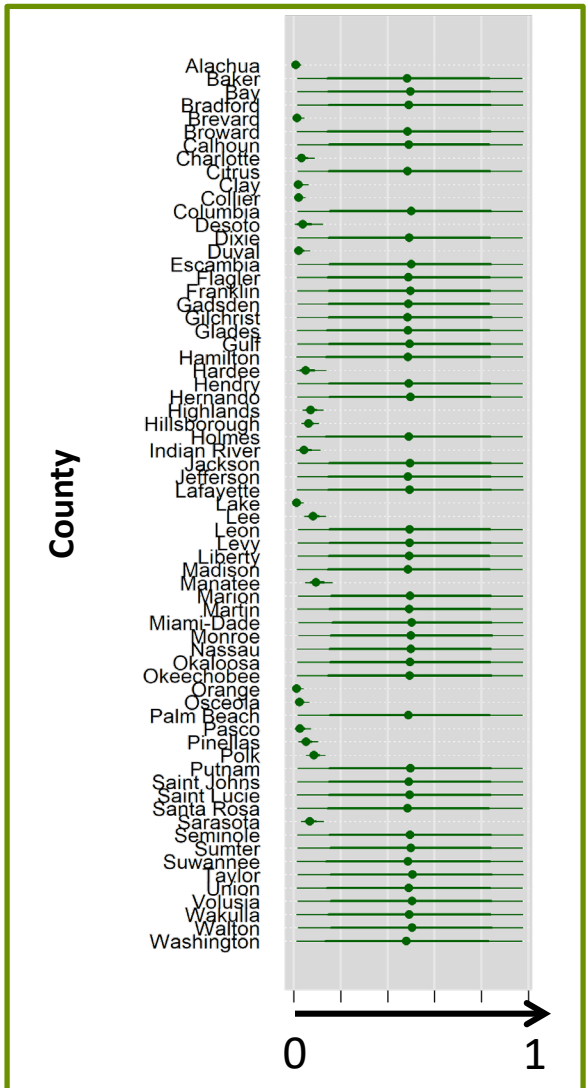
Uncertainty in Estimate of Occupancy (sd)

# Models that incorporate collector covariates shrink uncertainty in posterior

# Including collector behavior decreases uncertainty

No collector behavior:



Including collector behavior:

*Detectability = f(# of records)*

# Takeaway

- Accounting for collector behavior improves models

- Standardizing collector name entry important

- Developing new R package -> collectR

Broad-scale efforts to standardize and clean collector covariates are a worthy investment to improve the efficacy of digital biodiversity data for modeling species' ranges.