# Biodiversity & Big Data: Potential for AI

Annika Smith

iDigBio
Florida Museum of Natural History
University of Florida

# Biodiversity & Big Data: Potential for AI

- What are natural history collections?
- What is iDigBio?
- What kinds of biodiversity data are available?
- What kinds of questions do researchers ask with specimen data?
- How can specimen data be linked with other types of data?
- How has AI been used with specimen data, and what are potential future uses?
- What are current challenges to AI and specimen data?

# Museum Collections: The Library of Life

~1,600 natural history collections in the US

1–2 billion specimens in the US
3–4 billion specimens worldwide

# Systematics & Taxonomy



Linnea (twinflower)

Carl Linné, aka Carolus Linnaeus

# Museum Collections:  The Library of Life

Genetics
Genomics
Chemistry…

Species interactions
Phenology
Biogeography
More!

# Museum Collections:  The Library of Life

Most specimens locked away in cabinets, unavailable for general use.

# Museum Collections: The Library of Life

Most specimens locked away in cabinets, unavailable for general use.

**DIGITIZATION!!!!**

# Label Data from Herbarium Specimens

- Scientific name – including authority
- Date
- Collector
- Location – state, county, specific site, GPS coordinates
- Associated species
- Notes

# iDigBio:  www.idigbio.org



iDigBio
Integrated Digitized Biocollections

**National Coordinating Center**
**For Digitization of Biodiversity Collections**
Ingest, serve, integrate data:
Localities
Dates
Images

# Digitized Data & Biodiversity Research



www.idigbio.org

# Search Specimen Records: idigbio.org



www.idigbio.org

# Search Specimen Records: idigbio.org



Top 1 Taxa
- Acer rubrum
- other

Enter search criteria before using this map.

www.idigbio.org

# Search Specimen Records: idigbio.org

# Search Specimen Records: idigbio.org

## Specimen Record

Plantae > Tracheophyta > Magnoliopsida > Sapindales > Sapindaceae

### *Acer rubrum* L.

From Louisiana State University, Shirley C. Tucker Herbarium

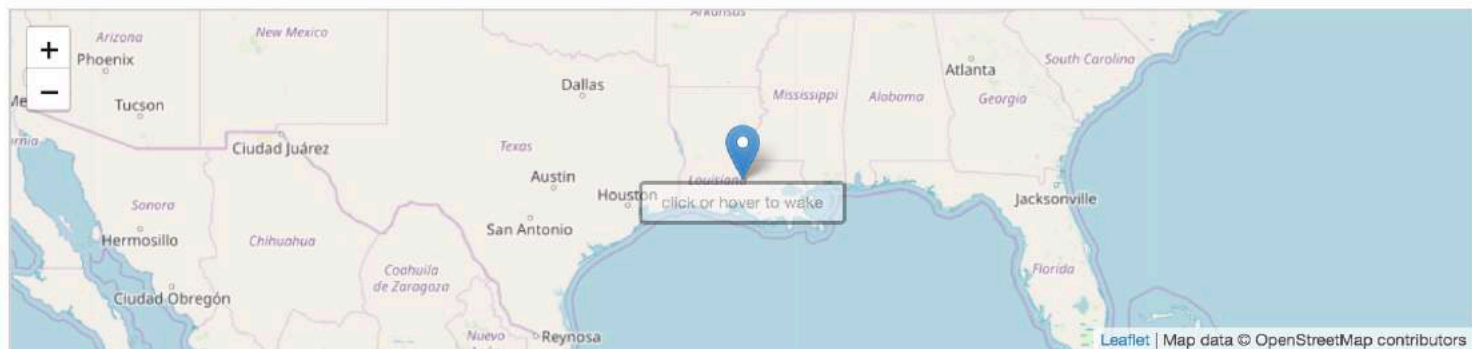| | |
| --- | --- |
| Continent | North America |
| Country | United States |
| State/Province | Louisiana |
| County/Parish | St. Martin |
| Locality | Atchafalaya National Wildlife Refuge: Sherburne Wildlife Management Area: Sse Of Krotz Springs. Collections Along Wooded Edge Of La Hwy 975, Ca. 8.5 Mi Nnw Of Junction La Hwy 975 And Interstate Hwy 10.; Atchafalaya National Wildlife Refuge |
| Latitude | 30.4581 |
| Longitude | -91.7302 |

| | |
| --- | --- |
| Institution Code | Lsu |
| Collection Code | Vascular Plants |
| Catalog Number | Lsu00132915 |
| Collected By | Marisa Conner |
| Date Collected | 2007-03-17 |



Leaflet | Map data © OpenStreetMap contributors

## Media

# Search Specimen Records: idigbio.org



## Media Record

Plantae > Tracheophyta > Magnoliopsida > Sapindales > Sapindaceae

### *Acer rubrum* L. view specimen record

From Louisiana State University, Shirley C. Tucker Herbarium

Media retrieved from:
http://images.cyberfloralouisiana.com/images/specimensheets/lsu/0/13/29/15/LSU00132915.JPG

Open in browser

Download File

# Search Specimen Records: idigbio.org



## Media Record

Plantae > Tracheophyta > Magnoliopsida > Sapindales > Sapindaceae

### *Acer rubrum*  L.  view specimen record

From Louisiana State University, Shirley C. Tucker Herbarium

Media retrieved from:
http://images.cyberfloralouisiana.com/images/specimensheets/lsu/0/13...

Open in browser

Download File

# Specimen Localities in iDigBio



Record Density

1
4
13
43
149
521
1,818
6,348
22,174
77,457
270,573

3000 km
2000 mi

Leaflet | Map data © OpenStreetMap

# Other Data Aggregators

# Components of iDigBio

G. Nelson, Director



Cyber-Infrastructure
(Fortes)

Digitization
(Riccardi)

Serving the
Research
Community
(Soltis)

Education &
Outreach
(MacFadden)

# Components of iDigBio

G. Nelson, Director



Cyber-Infrastructure (Fortes)

Serving the Research Community (Soltis)

Digitization (Riccardi)

Education & Outreach (MacFadden)

# Using Specimen Data for Research

## Big questions in biodiversity research:

– How many species are there?

– Why are species where they are?

– How does habitat change affect species?

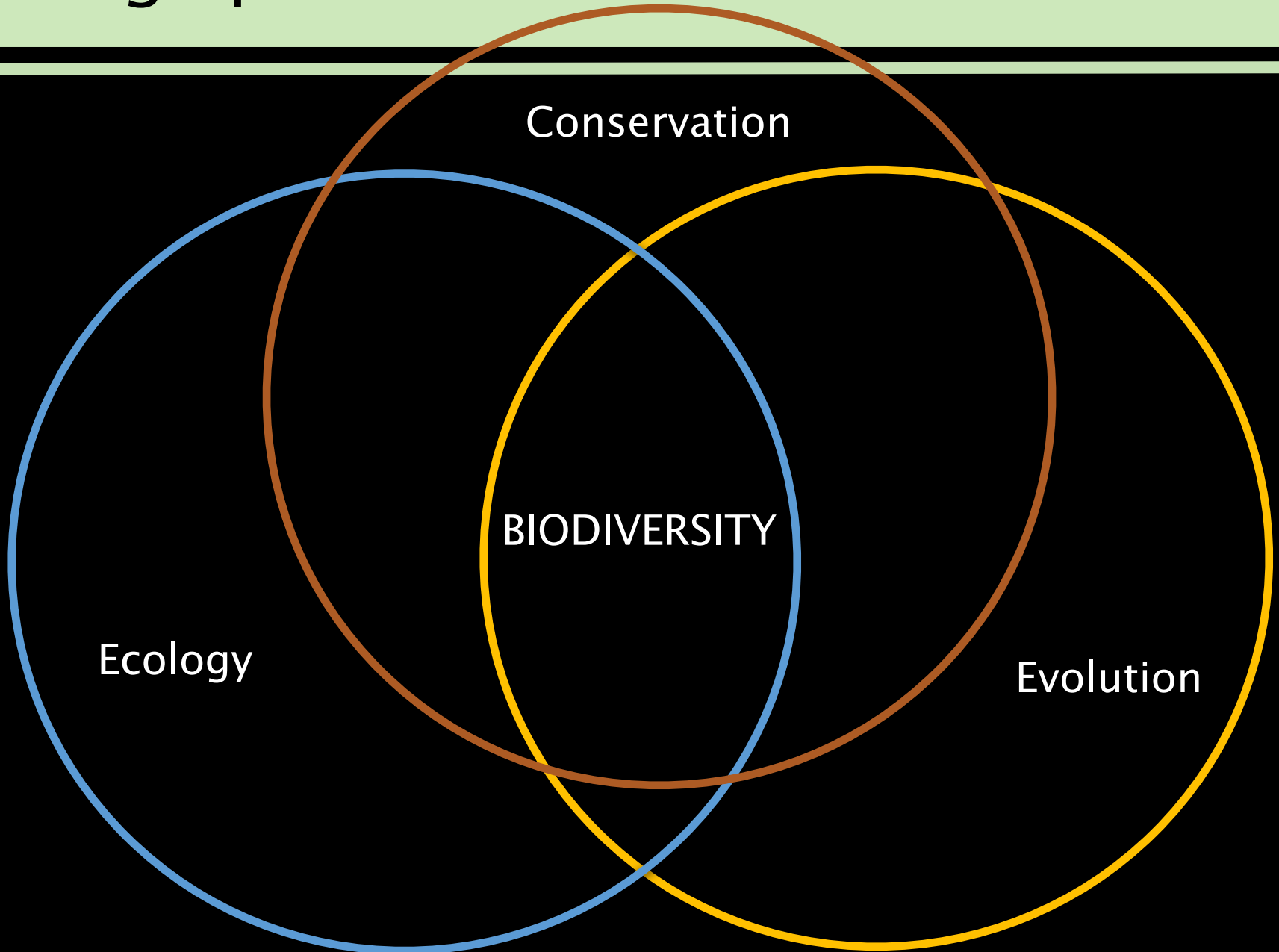– How does invasion by exotic species affect the stability of communities of species?

– How will a changing climate affect species— extinction rates, distribution, evolution?

# Using Specimen Data for Research

Conservation

Ecology

Evolution

BIODIVERSITY

# What is **<u>biodiversity</u>**?

"<u>Biodiversity</u> defies easy definition, but we value it nonetheless, much the same way we value <u>justice</u>, <u>freedom</u>, and <u>nature</u>, similarly difficult terms to define."

Naeem S et al. 2016, Proc. R. Soc. B

# Different Measures of Biodiversity

- Species richness— How many species are in the area?

- Richness of endemic species- How many endemic species are there?

- Functional diversity – How similar are the functional traits of species in one area to one another?

- Phylogenetic diversity- How closely related are the species in an area?

Naeem S et al. 2016, Proc. R. Soc. B

# Using AI to generate data for biodiversity research

- Can AI be used to identify plants?

- Can AI be used to label phenological stages of plants?

# Using AI to generate data for biodiversity research

- Can AI be used to identify plants?

- Can AI be used to label phenological stages of plants?
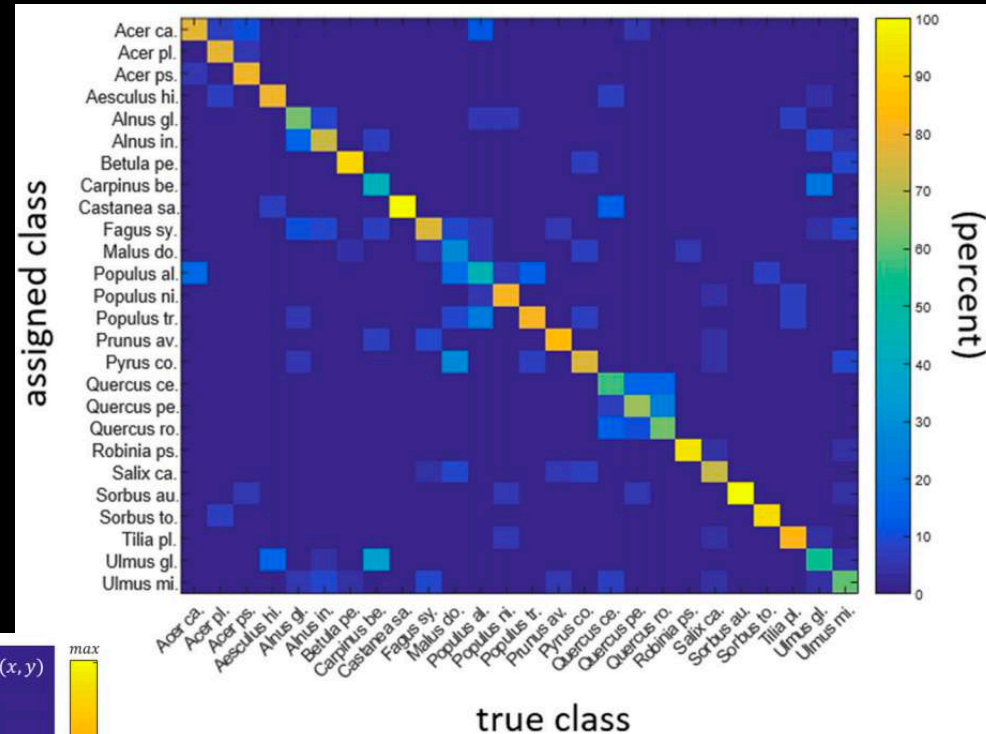
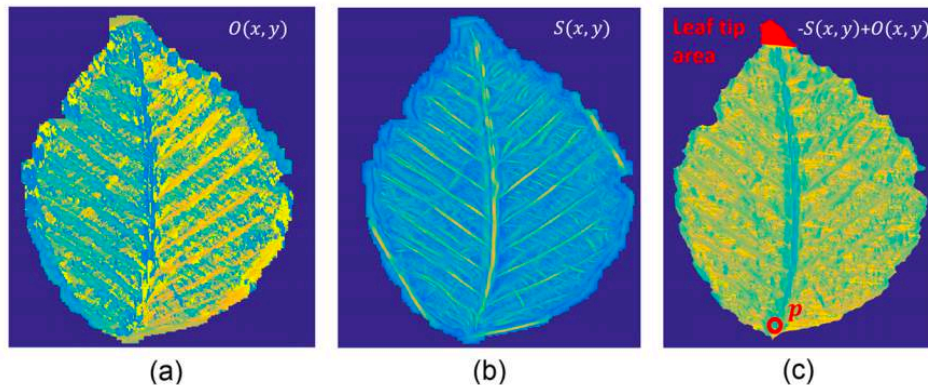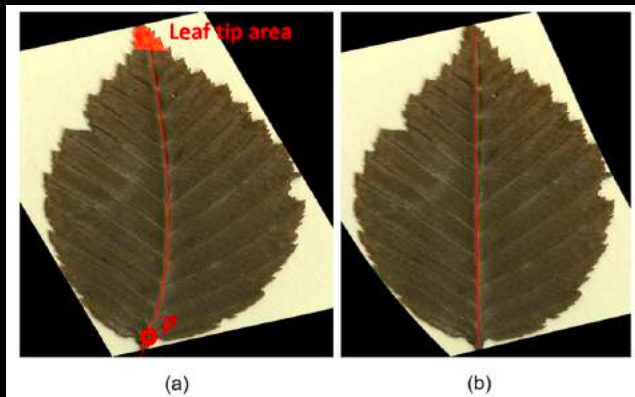<span style="color:yellow">YES!</span>

# Machine Learning:  Plant ID

**Machine Learning:** Herbarium specimens
Classifying German trees to species
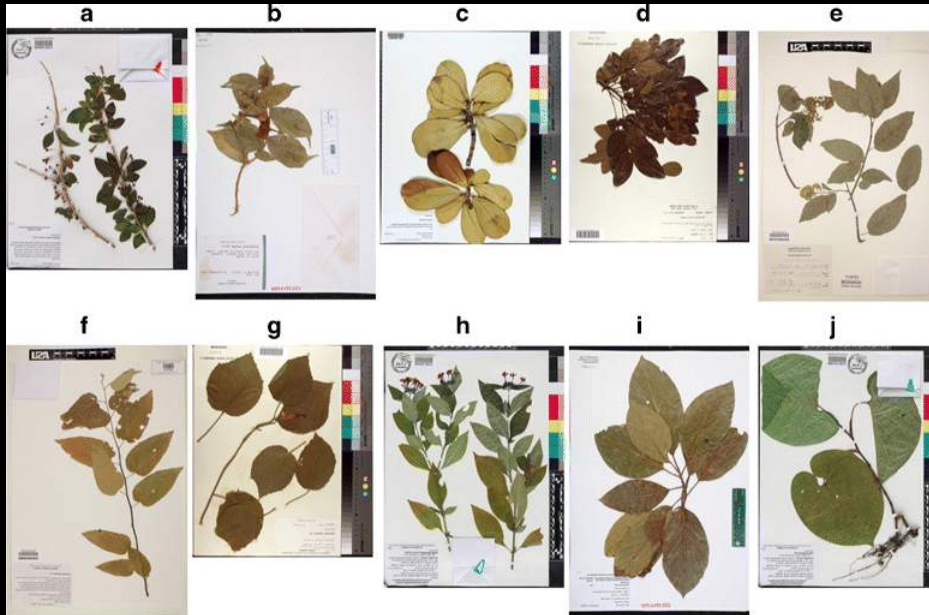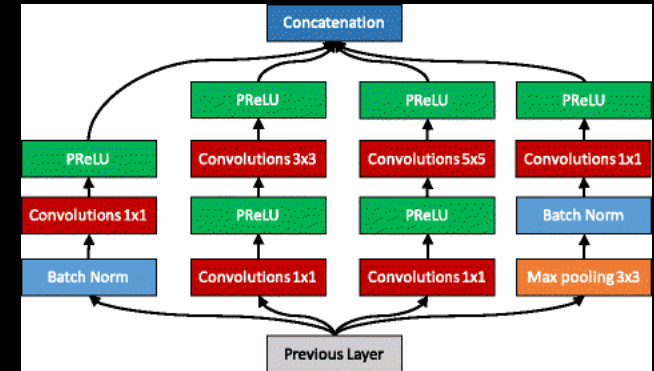Leaf shape, venation
85% accuracy



Unger et al. 2016
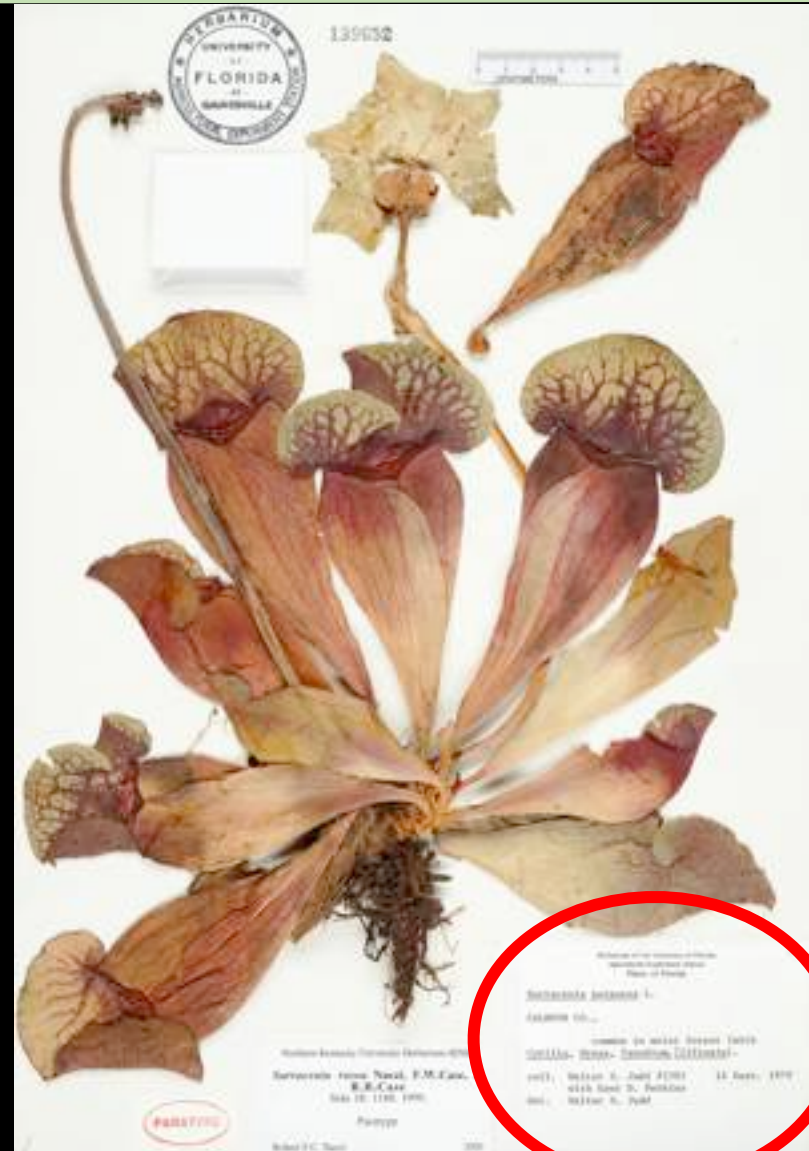
# Machine Learning:  Plant ID

Deep Learning: Herbarium specimens (2 data sets)
Classifying plants to species
>1200 species
250,000 images from iDigBio
90% accuracy





Carranza-Rojas et al. 2017
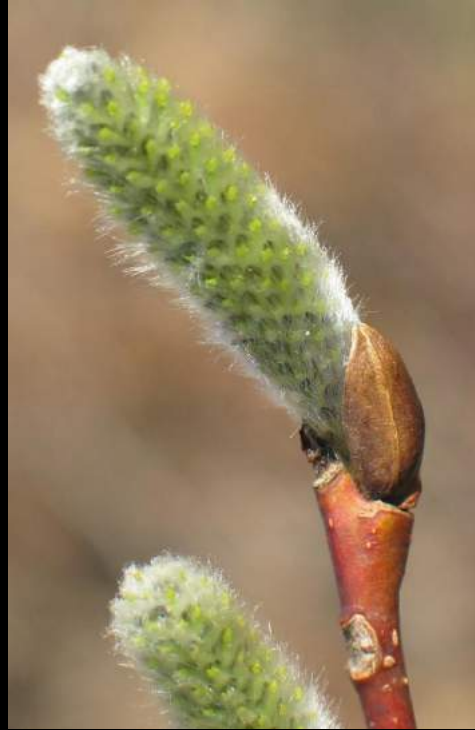
# Label Data from Herbarium Specimens

- Scientific name –
  including authority
- Date
- Collector
- Location – state, county,
  specific site,
  GPS coordinates
- Associated species
- Notes

# Phenology:  "Nature's Calendar"
## Bud Burst, Flowering, Fruiting



www.sbs.utexas.edu   www.sites.psu.edu

Phenological data – as described in label notes "tree in full fruit…"

or from image itself:

CrowdCurio, in Willis et al. 2017



Trends in Ecology & Evolution

# Machine Learning and Phenology



Maples (*Acer*): unfolded leaves present or absent?

unfolded leaves present

unfolded leaves absent

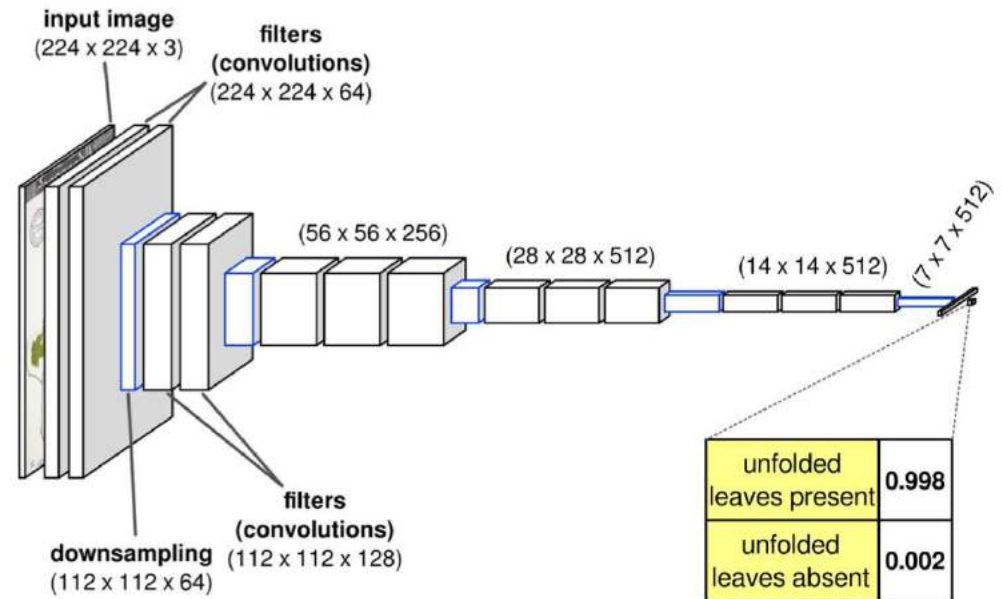*Prunus*: flowers present or absent?

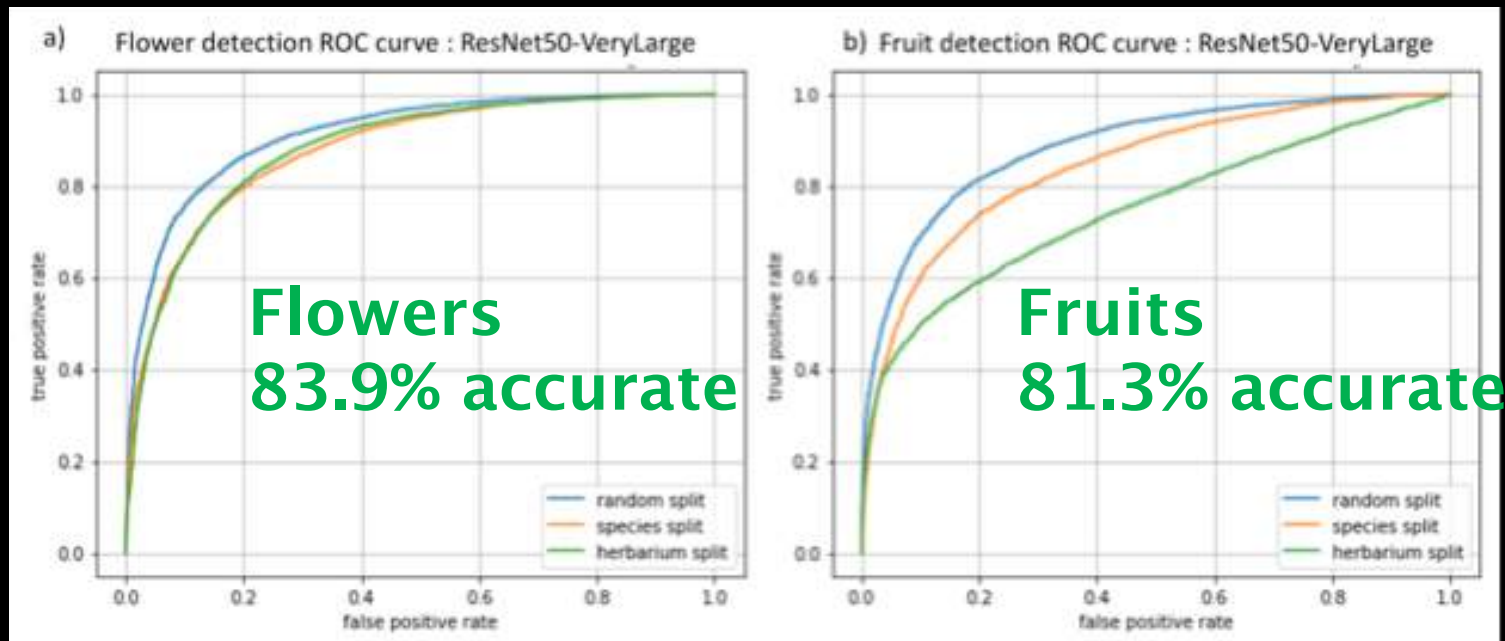flowers present

flowers absent

Deep convolutional neural network

input image (224 x 224 x 3)

filters (convolutions) (224 x 224 x 64)

downsampling (112 x 112 x 64)

filters (convolutions) (112 x 112 x 128)

(56 x 56 x 256)

(28 x 28 x 512)

(14 x 14 x 512)

(7 x 7 x 512)

| | |
|---|---|
| unfolded leaves present | 0.998 |
| unfolded leaves absent | 0.002 |

Architecture based on Simonyan and Zisserman (2015), arXiv:1409.1556v6

Brian Stucky

# Machine Learning and Phenology

Deep learning annotation
3 herbarium data sets:
    163,233 specimens
    7782 species, 236 families
4th data set: Asteraceae



**Flowers**
**83.9% accurate**

**Fruits**
**81.3% accurate**

Lorieul et al. 2019

# Phenological mismatches between species

Insect Emergence, Bird Nesting & Migration, and Phenological Synchrony/Asynchrony

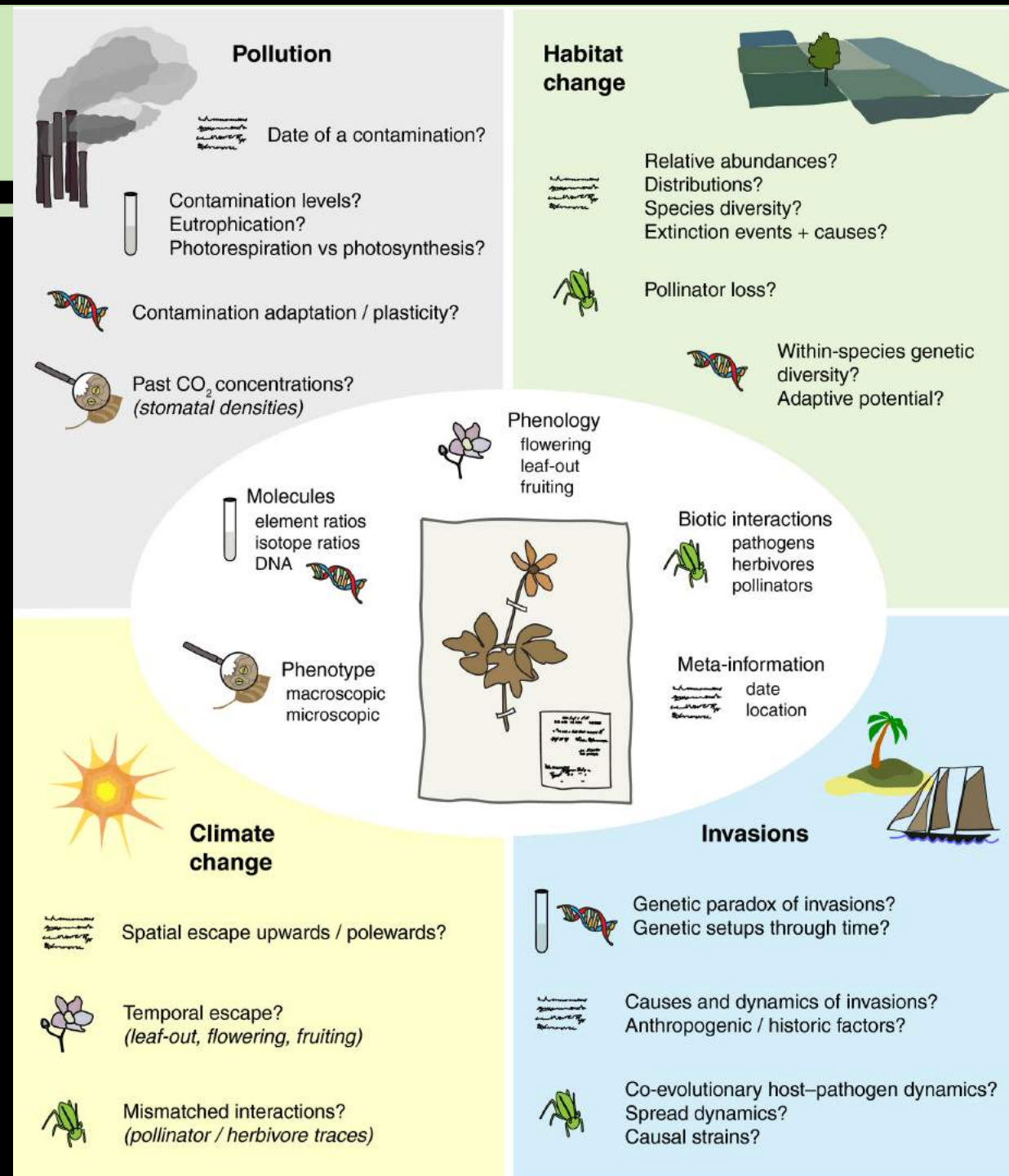# How will organisms respond to climate change?

Natural history collections contain data with….

**Temporal Range**
**Taxonomic Breadth**
**Geographic Diversity**
**Intraspecific Diversity**

…. invaluable to make predictions
about how species will behave
in the future.

# Other uses of specimens

Lang et al. 2019, New Phytologist

# Linking Specimen Data to other Data Sources

**Climate**
- World-Clim

**Genetic Data**
- Genbank
- 1KP project

**Functional Trait Databases**
- TRY

**Character Traits**
- Extract from literature
  - (eg. Biodiversity Heritage Library)
- Extract from specimens (images, etc)

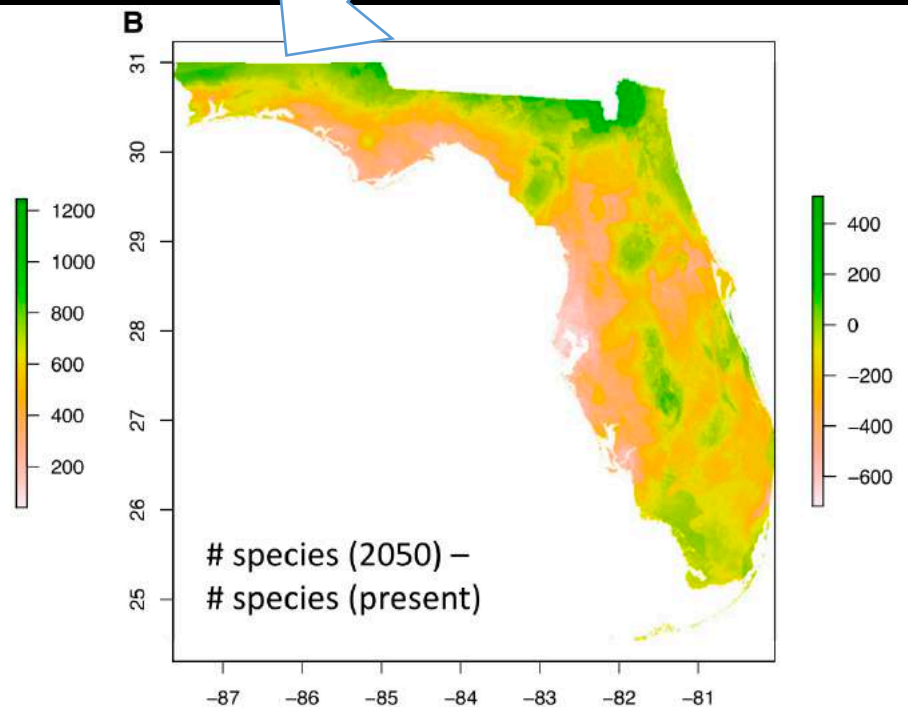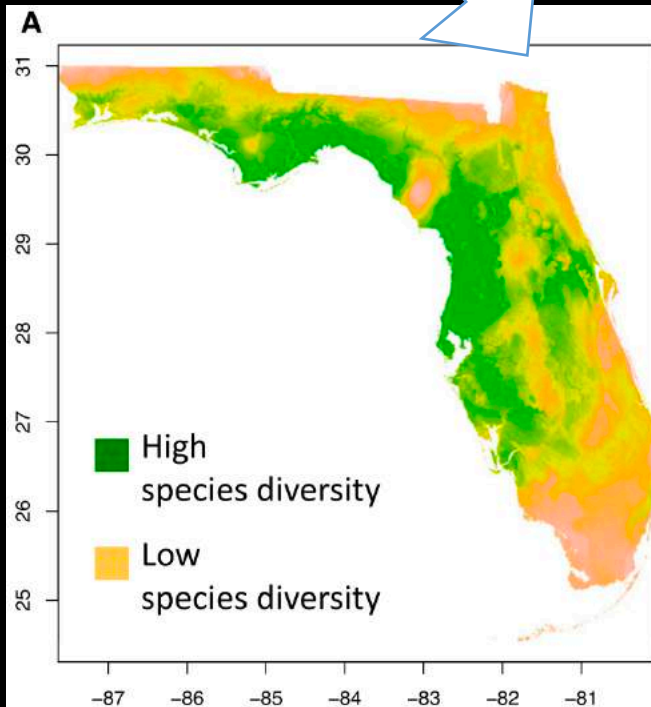… and more (eg. soil)

# Specimen Data linked with Climate Data

Vascular Plant Diversity of Florida— Ecological Niche Modeling with 1500 species of Florida Plants, based on >500,000 specimens

Soltis P. 2017,
American Journal of Botany

# Specimen Data linked with Climate and Genetic Data



WorldClim - Global Climate Data
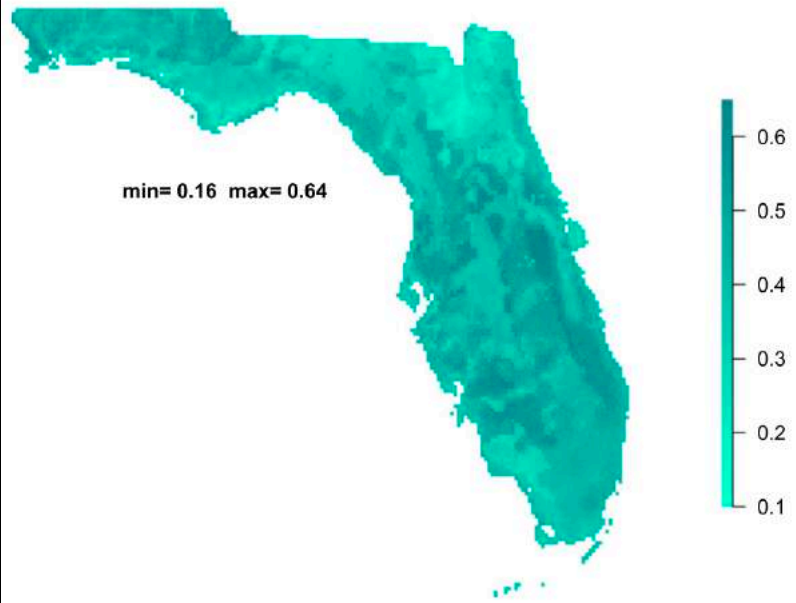
Free climate data for ecological modeling and GIS

iDigBio
Integrated Digitized Biocollections

NCBI    Resources    How To

GenBank                    Nucleotide

**A  Proportional Phylogenetic Diversity**

*Non−ultrametric tree*

min= 0.16  max= 0.64

0.6
0.5
0.4
0.3
0.2
0.1

Soltis P. 2017,
American Journal of Botany

# Specimen Data Linked with Functional Trait Data



**Plant Trait Database**

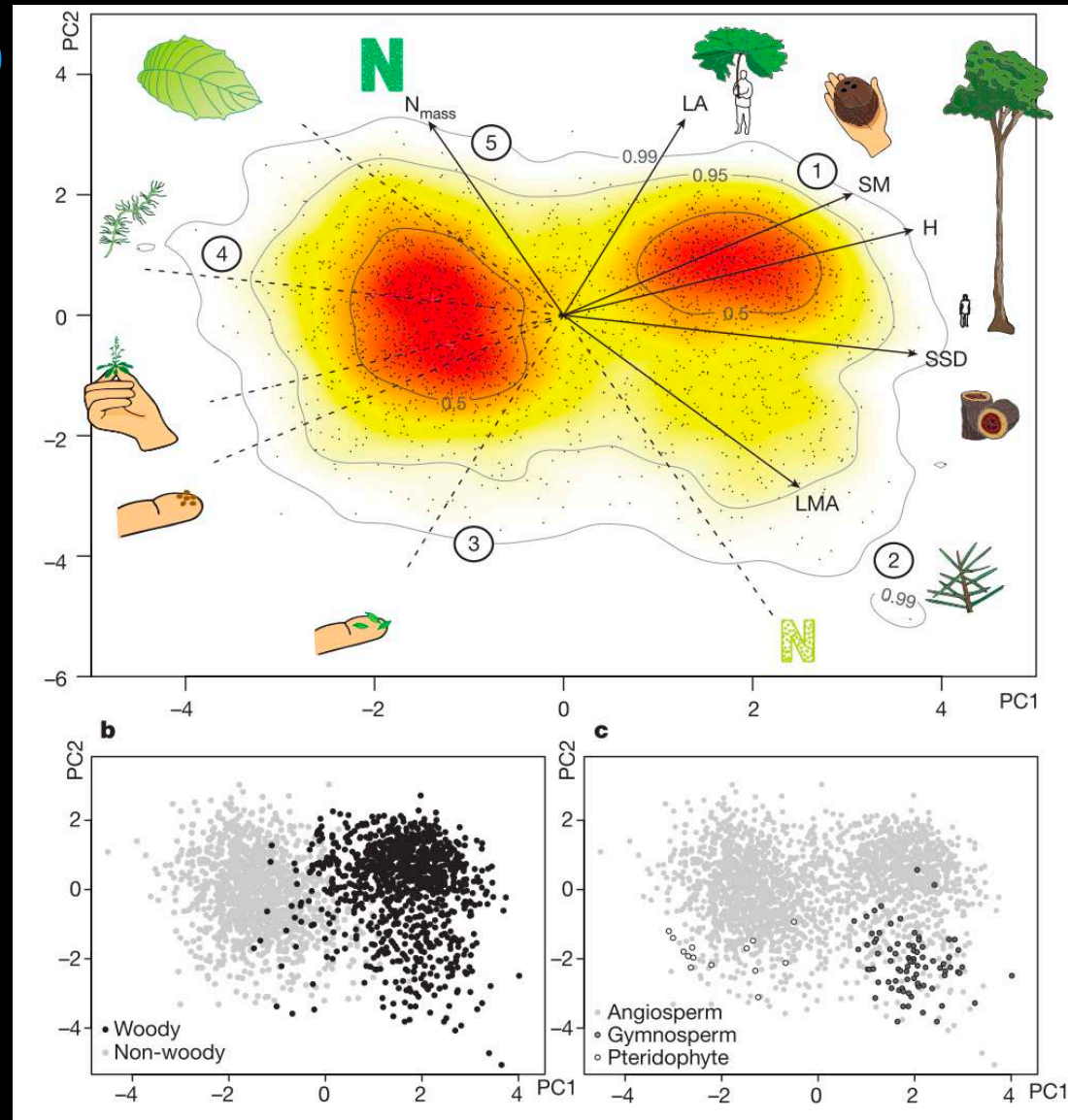Daphne mezereum
91 traits
Photo by A. Günther

PhotosyntheticPathway
Respiration LeafArea NfixationCapacity
SLA RegenerationCapacity PlantLifespan
WoodDensity GrowthForm
PhenologyType LeafN
LeafP LeafLongevity PhotosyntheticCapacity
MaxPlantHeight SeedMass

# Specimen Data Linked with Functional Trait Data

## Diaz et al. Nature (2016)

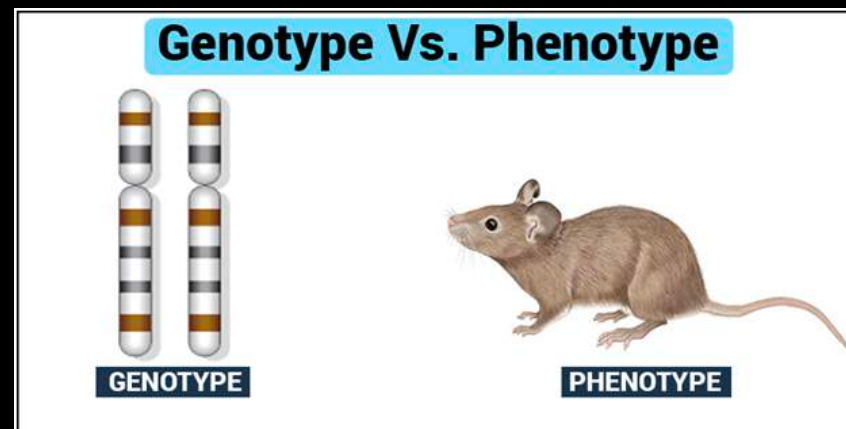Data from 46,085 plant species from 423 families

Reveals "hot spots" in the trait-space of vascular plants— some combinations of traits have evolved repeatedly across many different clades

# Potential for AI— Phenotype Data

Big data in genomics, ecology (data layers), geography... but what about phenotype?

Phenotype-- observable traits above a molecular level (anatomy, morphology, behavior)

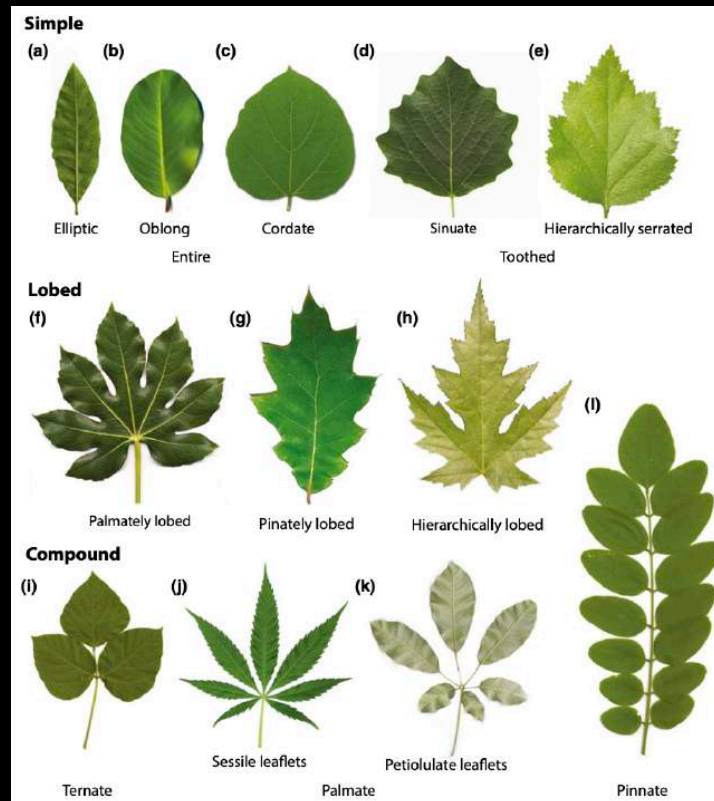# Potential for AI— Phenotype Data

The phenotype bottleneck:

"While phenotype data are as complex, diverse, and nuanced as genomic data, they have not seen data standardization and analyses applied with the same broad strokes as we have seen for genomics."

–Deans et al. 2015, PLOS Biology

In order to understand the evolution and ecology of biodiversity, it's necessary to develop ways of quantifying phenotypic diversity in a way that is comparable across species.



Runions, Tsiantis, & Prusinkiewicz
2018, New Phytologist

# Why is leaf shape interesting?

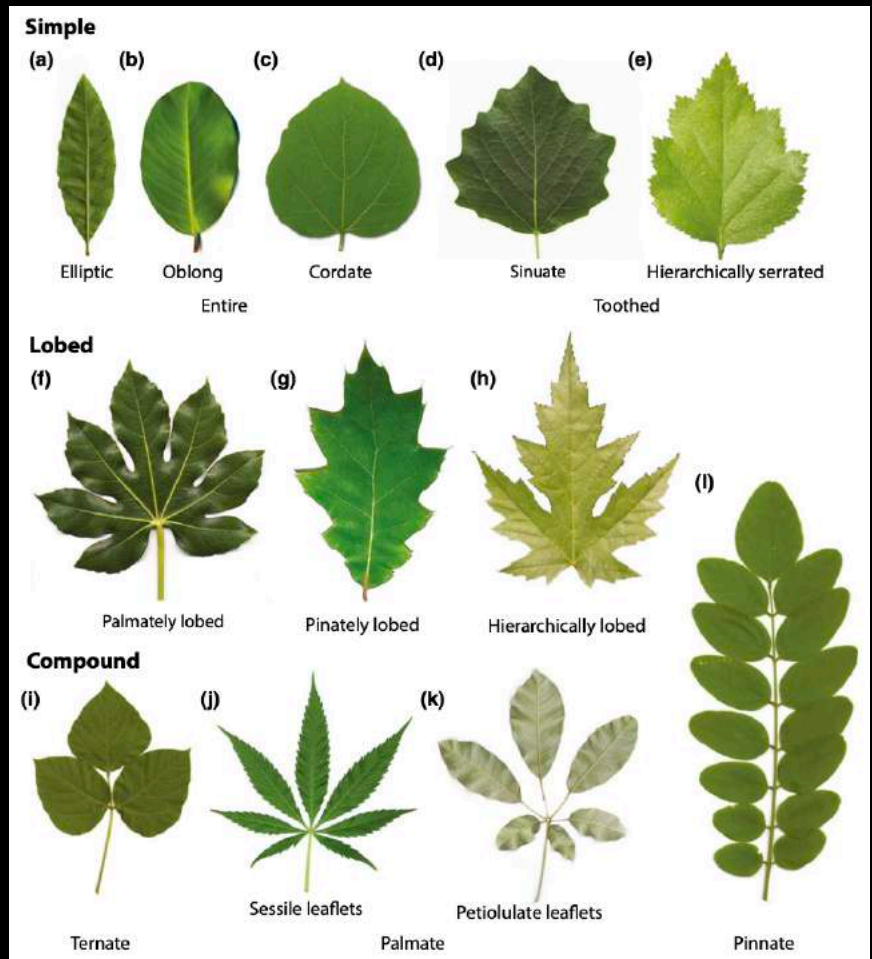**Research at the interface of Ecology-Evolution-Development (Eco-Evo-Devo)**

• Leaf shape linked to functional diversity, physiology

      eg. dissection, toothineess of leaf linked to
      thermoregulation of leaves and water balance

• Evolutionary constraints and patterns

• Ecological constraints and patterns (eg. habitat filtering)

• Plasticity and adaptability

• Developmental and genetic mechanisms of leaf shape
variability largely understood for many plant groups.

# AI and Images— Future Research Potential

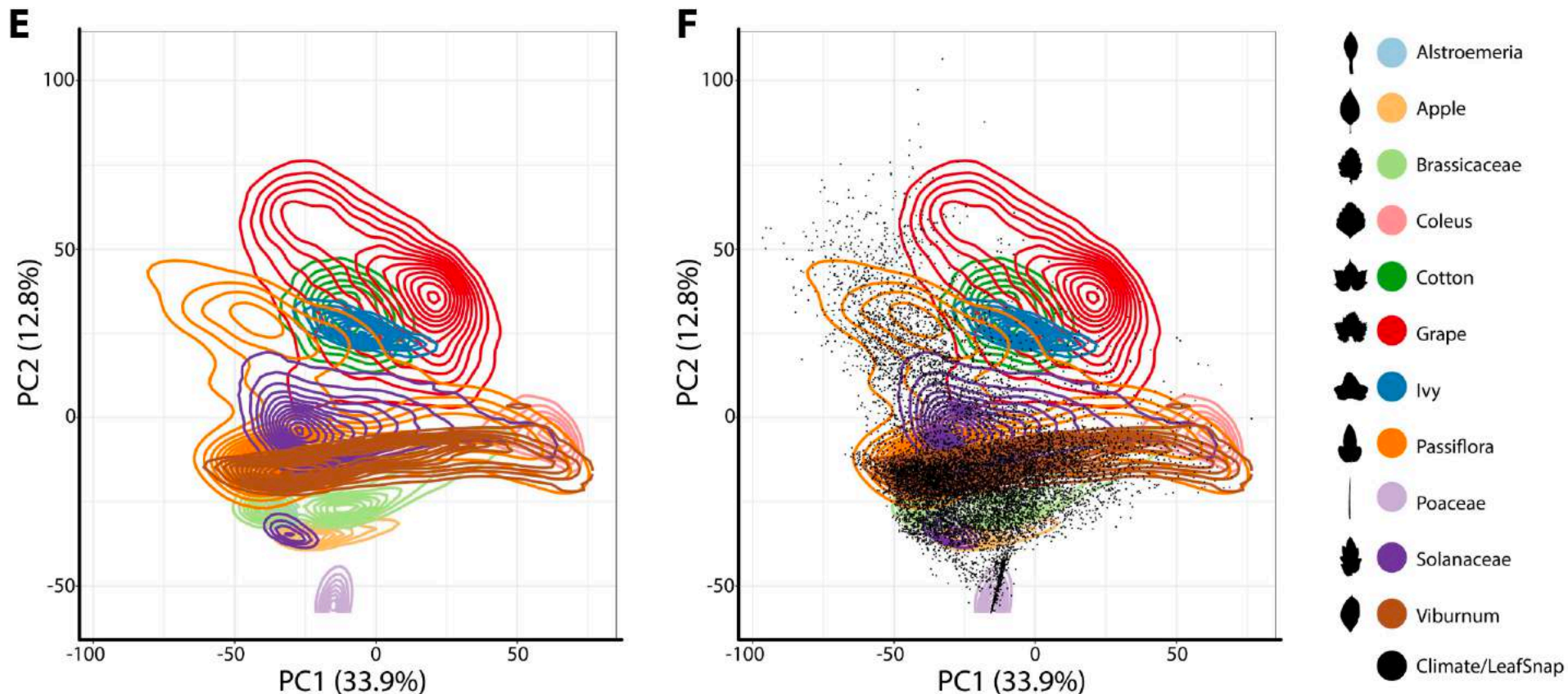Historically, botanists have described leaf shape in qualitative terms: lobed, compound, pinnate, etc.

Challenging to quantify shape:

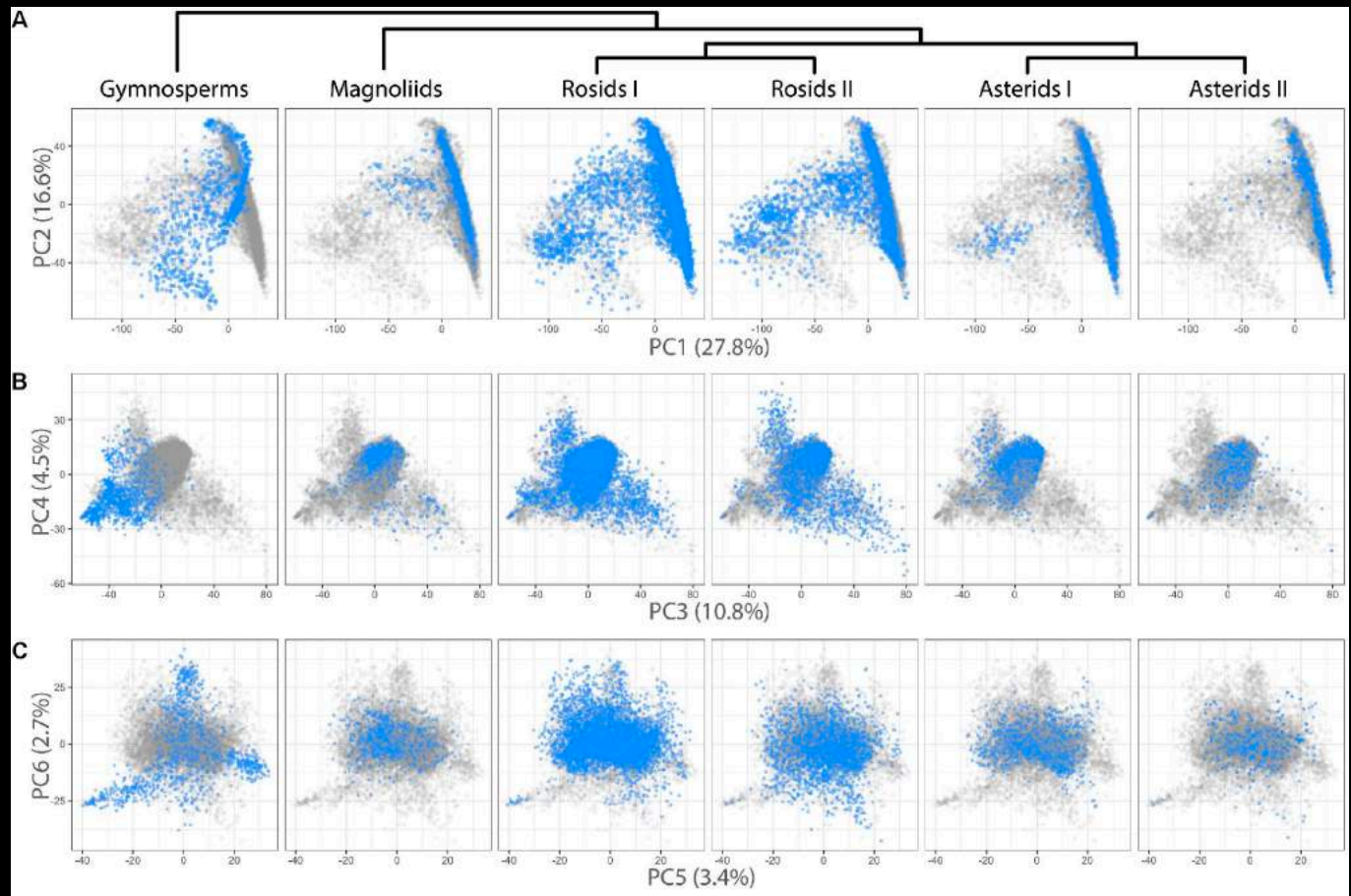Traditional land-mark based morphometric methods are only useful for closely related plants.



Runions, Tsiantis, & Prusinkiewicz
2018, New Phytologist

# Advances in quantifying leaf shape

Comparing the morphospace of leaves in different clades using persistent homology
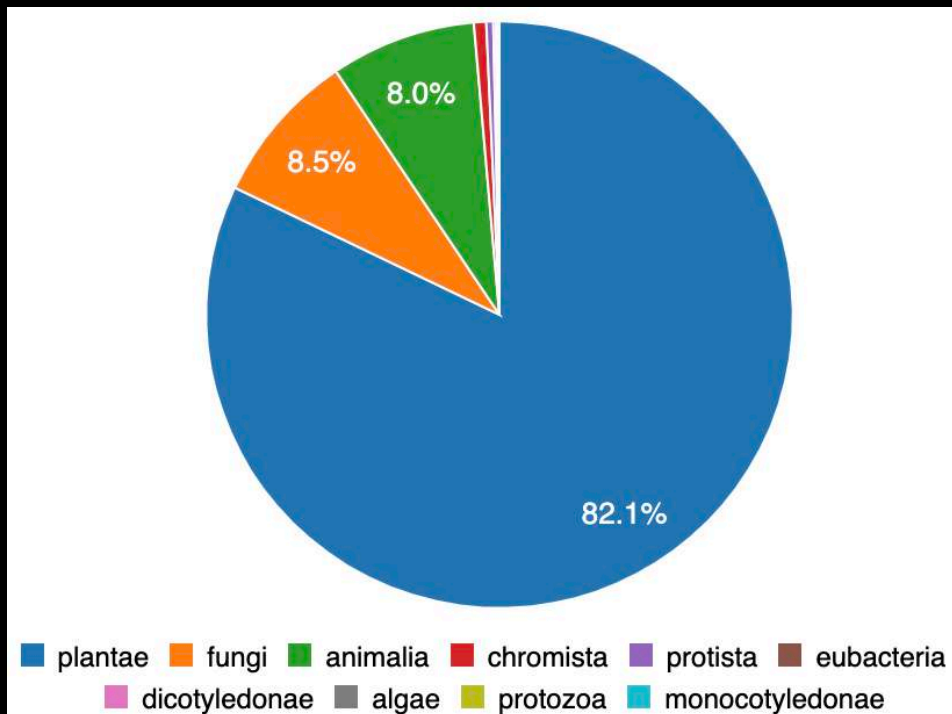


Li et al. 2018, Frontiers in Plant Science

# Advances in quantifying leaf shape

Comparing the morphospace of leaves in different clades using persistent homology— linking shape and phylogeny
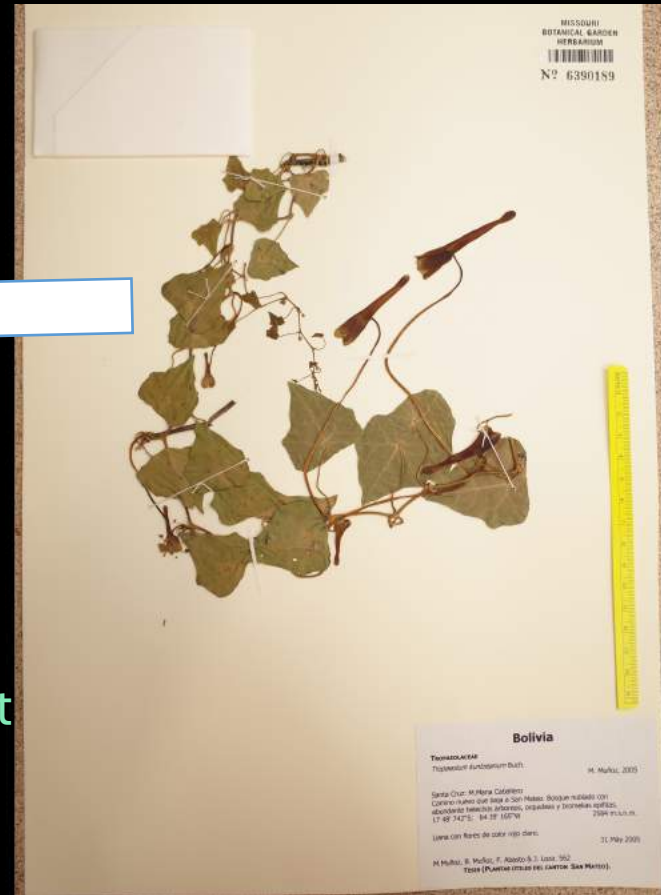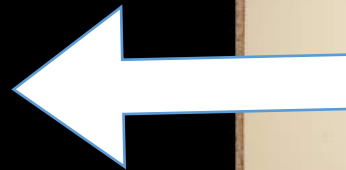


Li et al. 2018, Frontiers in Plant Science

The BIG question: Can AI unlock natural history collections as a resource for phenotype data?



29,523,938 Media Records

The BIG question: Can AI unlock natural history collections as a resource for phenotype data?



A single image contains a great deal of information about the shape and size of plant organs— could AI methods (eg. semantic segmentation) help extract and quantify this?

# Leaf and flower diversity in the nasturtiums (Tropaeolaceae)
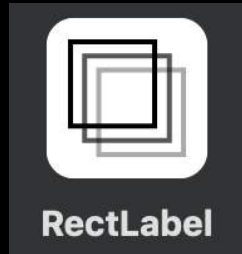
**Complications:**

**Overlapping leaves and flowers**

**Some parts broken or with insect damage**

**Plant organs of various types and developmental stages**

# Building a training data-set for semantic segmentation

The Challenge: Labeling individual pixels is exceptionally time consuming....
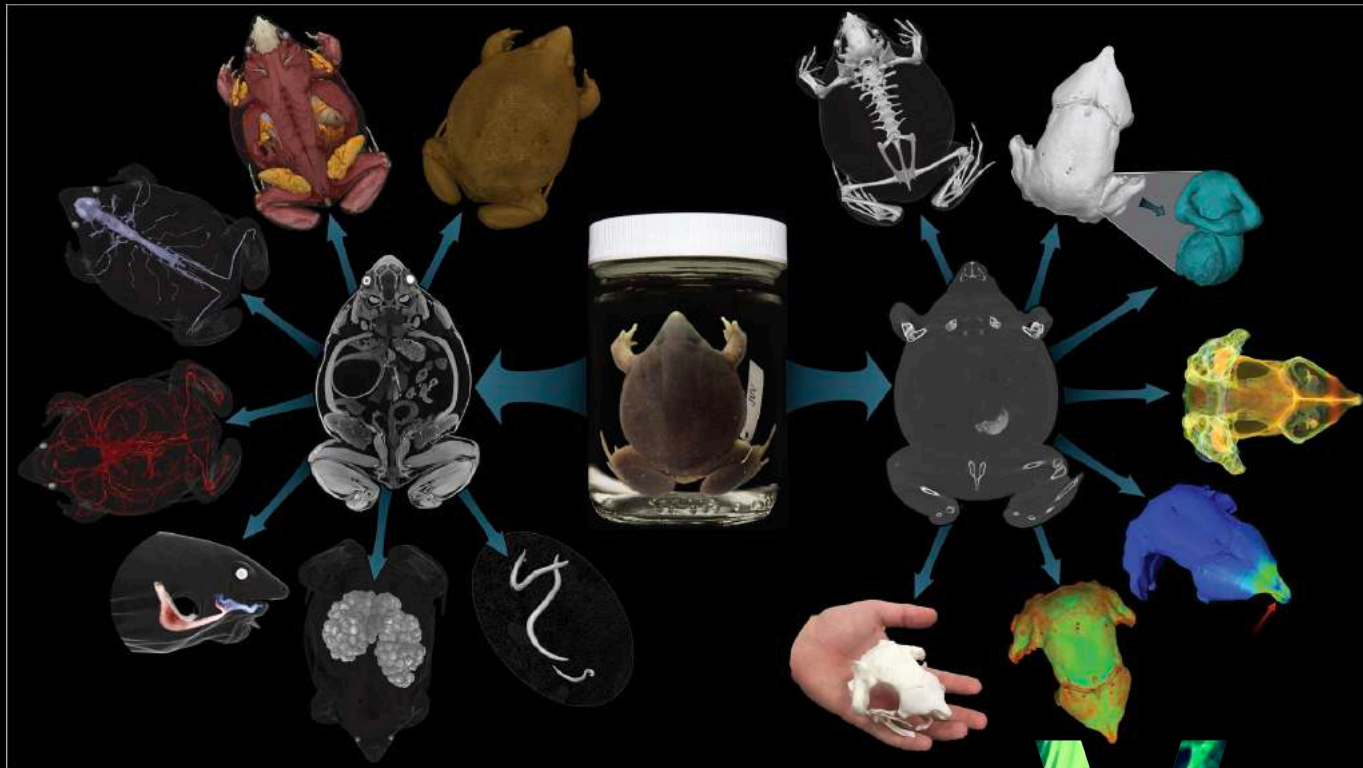
# Some Challenges: AI and Trait Data

- Vast diversity of the structure of organisms— need massive training data sets?
- Biases in specimen collection
- Images: Overlapping or damaged structures.
- Some characters too small or hidden.
- Changes between fresh and dry plants?
- Data quality— the balance between the abundance of data vs imperfect data

# Vertebrate Phenotype Data– oVert

oVert (iDigBio TCN). CT scans of vertebrate specimens, freely available through Morphosource.



3D printing for research, education, and outreach

https://www.floridamuseum.ufl.edu/overt/

# Acknowledgments



Pamela Soltis

# Thank you! Questions?

Contact:
Annika Smith
annikals@ufl.edu
@meristemming