

## Webinar 1

# Introduction: Scope and Research Potential for Multidisciplinary Biodiversity Modeling and Analysis

**Introduce biological and computational workflow concepts for integrating: biological specimen (species occurrence) data, phylogenetic trees, and multi-species distribution models!**

## **Biological Objectives:**

- ✓ **Provide context, rationale and scope for the biological research to be presented.**

## **Technical Objectives :**

- ✓ **Provide context and differentiate BiotaPhy technology from desktop approaches (software, scaling, capabilities, integration)**

1. **Webinar Series Overview: topics and dates**
2. **Biological Scope Overview:**
  - a) **Types of Data**
  - b) **BiotaPhy One**
    - i. **Background and Workflows**
  - c) **BiotaPhy Two**
    - i. **Launched Workflows**
    - ii. **Proposed Workflows**

3. **Technical Scope Overview:**
  - a) **BiotaPhy's architecture**
  - b) **Tutorial Overview**
4. **Session Summary, Q&A and Discussion**

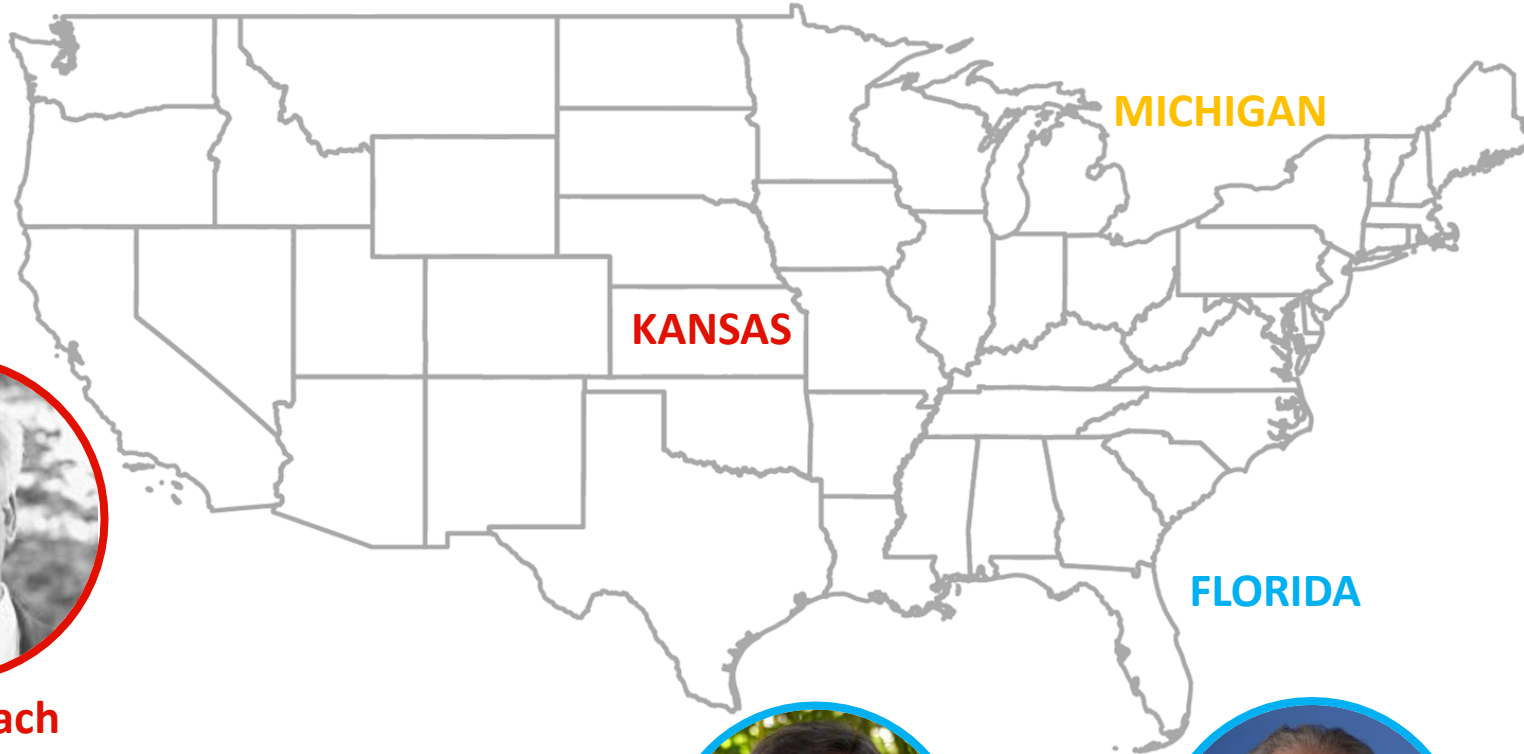
- ➔ **Introduction: Scope and Research Potential for Multidisciplinary Biodiversity Modeling and Analysis (09/21/2022)**
- ➔ **Resolving Nomenclature: Making Appropriate Taxonomic Choices (09/28/2022)**
- ➔ **Clean Your Dirty Data (10/05/2022)**

- ➔ **Georeferencing with GEOLocate** [10/12/2022]
- ➔ **Big Data Munging (a.k.a. Splitting and merging occurrence data by taxa from multiple sources)** [10/19/2022]
- ➔ **Species Distribution Modeling 1** [10/26/2022]
- ➔ **Species Distribution Modeling 2** [11/02/2022]

- ➔ **Introducing Presence-Absence Matrices for Large-Scale Analyses**  
**[11/9/2022]**
- ➔ **Phylogenetic Diversity: Integrating Phylogenies with Species and Biogeographic Data** **[11/16/2022]**
- ➔ **Hypothesis Testing and Randomization** **[11/30/2022]**



# BiotaPhy Crew – Principal Investigators



James Beach



Stephen Smith



José Fortes



Douglas Soltis



Pamela Soltis

# BiotaPhy Crew – Technical Crew



CJ Grady



Aimee Stewart



Michael Elliott



Srivattsan Sridharan

# BiotaPhy Crew – Biological Crew



Dr. Hannah Marx



Hector Figueroa



Elizabeth White



Nicholas Engle-Wrye



Dr. Ryan Folk



Dr. Anthony Melton



Dr. Jon Spoelhof



Maria Cortez

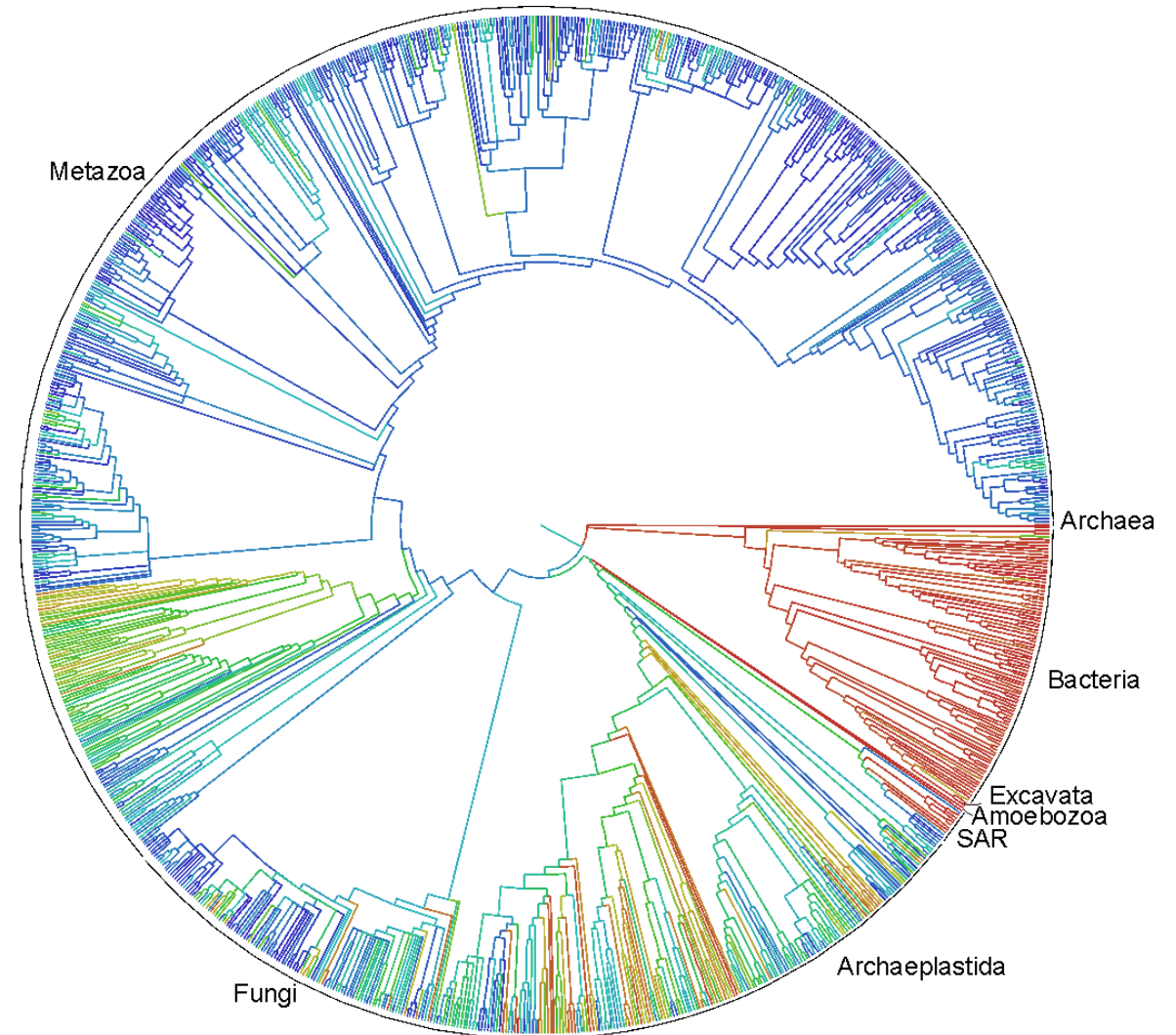
## Natural history specimens contain a

wealth of data: 

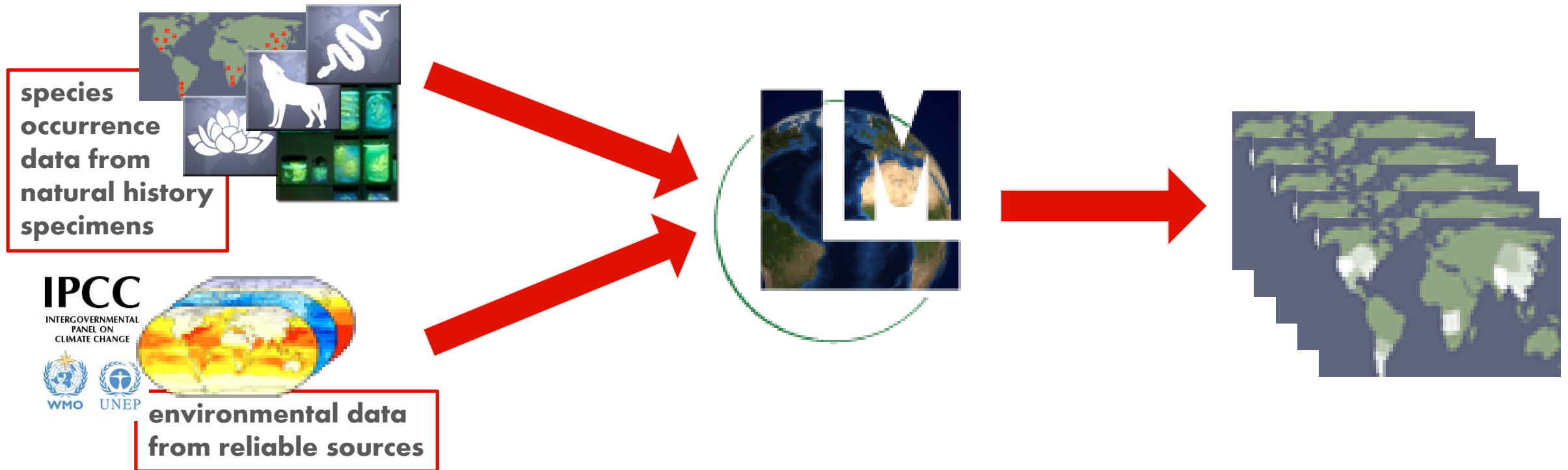
- ✓ Scientific name
- ✓ Date
- ✓ Collector
- ✓ Location – state, county, specific site, GPS coordinates
- ✓ Associated species



**Phylogenetic trees** tell a possible tale about evolutionary history



**Species Distribution Models (SDMs)** can be created by combining occurrence and environmental data



# Biological Scope: Types of Data

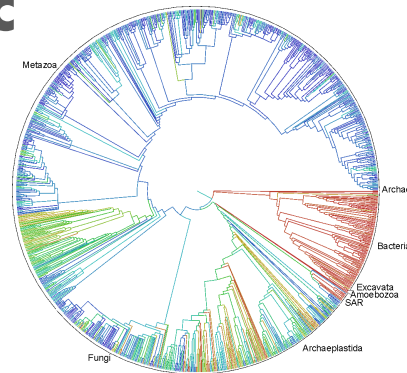
Integration of ...

Natural history  
specimen data

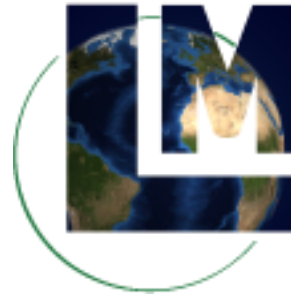


Allen et al.. 2019.  
Biodiversity synthesis  
across the green  
branches of the tree of  
life. *Nature Plants* 5: 11-13.

Phylogenetic  
trees



Species distribution  
models



**BIOTAPHY**

## BIOTAPHY 1 ...

and its 5 possible workflows:

### RESOURCES:

**Lm** Lifemapper

- ecological niche modeling
- biodiversity and range analysis
- visualization

**A** Analysis Tool Kit

- evolutionary models
- comparative methods
- visualization

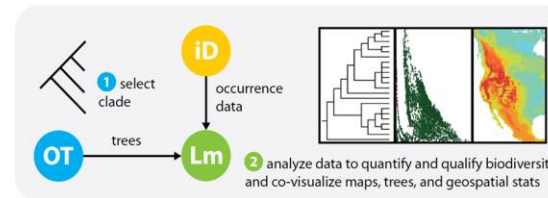
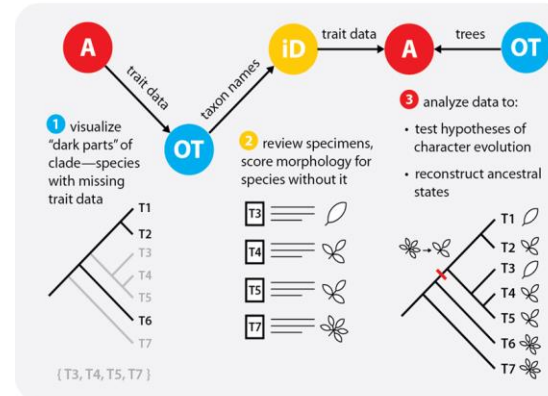
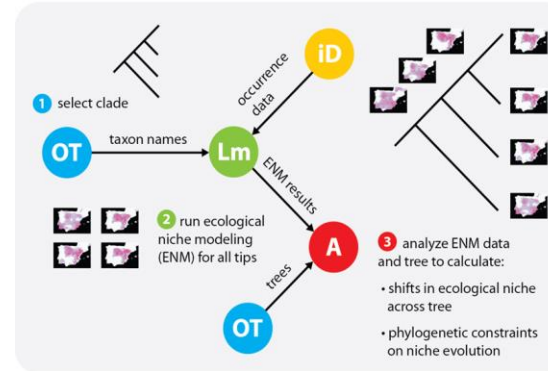
**OT** Open Tree of Life

- phylogenies
- taxonomy / names
- visualization

**iD** iDigBio

- trait data
- specimen data / images
- fossil data / images

### EXAMPLE WORKFLOWS:



### RESOURCES:

**Lm** Lifemapper

- ecological niche modeling
- biodiversity and range analysis
- visualization

**A** Analysis Tool Kit

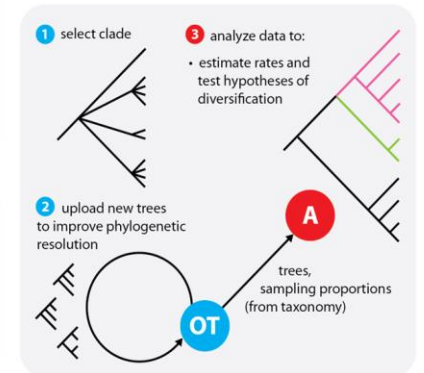
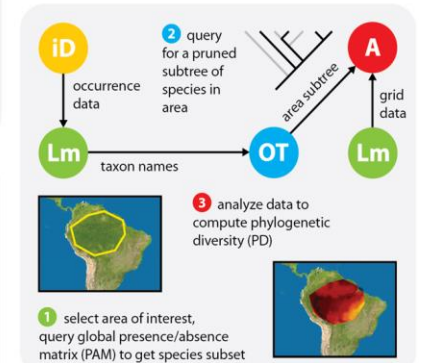
- evolutionary models
- comparative methods
- visualization

**OT** Open Tree of Life

- phylogenies
- taxonomy / names
- visualization

**iD** iDigBio

- trait data
- specimen data / images
- fossil data / images





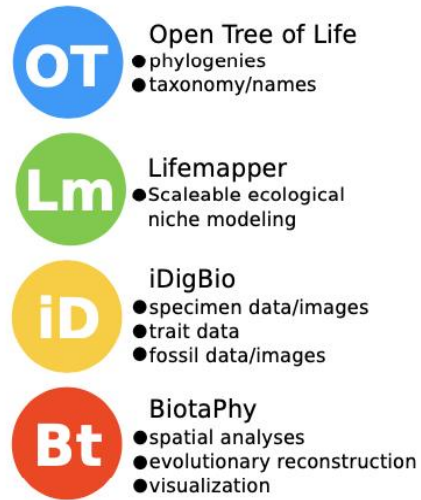
# Biological Scope: BiotaPhy Two



## BIOTAPHY 2 ...

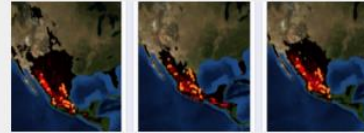
Evolved workflows:

- ✓ Launched
- ✓ Proposed (most have been implemented already!)



Launched workflows:

Generate and visualize multispecies distributions



Subset PAM with these filters

Matching species

Kingdom:

Phylum:

Class:

Order:

Family:

Genus:

Species:

Algorithm:

Model:

Projection:

B-Box:

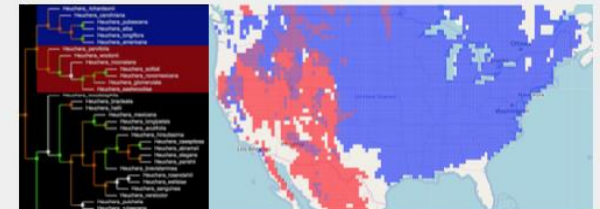
Abraxia flavinotata (De Geer, 1773)  
 Achurum carinatum (Walker, F., 1870)  
 Acrida cognata (Friedemann, 1915)  
 Acutina ferruginea (Osaka, 1899)  
 Acutina maculata (Cohn, 1919)  
 Acutina trancosa Curnutt, 1932  
 Acanemia nitidicollis (Meyen, 1818)  
 Acrida reticulata Guérin-Méneville, 1832  
 Acrocephalus tristis (Gyll., 1825)  
 Acrocephalus maculipennis (Scudder, S.H., 1890)  
 Apterops crinitus (Burmeister, H., 1838)

Subset by tree/names

MCPA—unravel history and environment

0.276 (0.000)	<b>BIOCLIM_1</b>	0.000 (0.992)	lgm_global
-0.133 (1.000)	BIOCLIM_7	0.421 (0.000)	<b>mississippi</b>
-0.179 (1.000)	BIOCLIM_12	0.000 (1.000)	continental_divide
-0.167 (1.000)	BIOCLIM_17	-0.146 (0.914)	eastern_divide
0.677 (0.000)	<b>ENV - Adjusted R-squared</b>	0.329 (0.000)	<b>BG - Adjusted R-squared</b>

Covisualize phylogeny and distribution



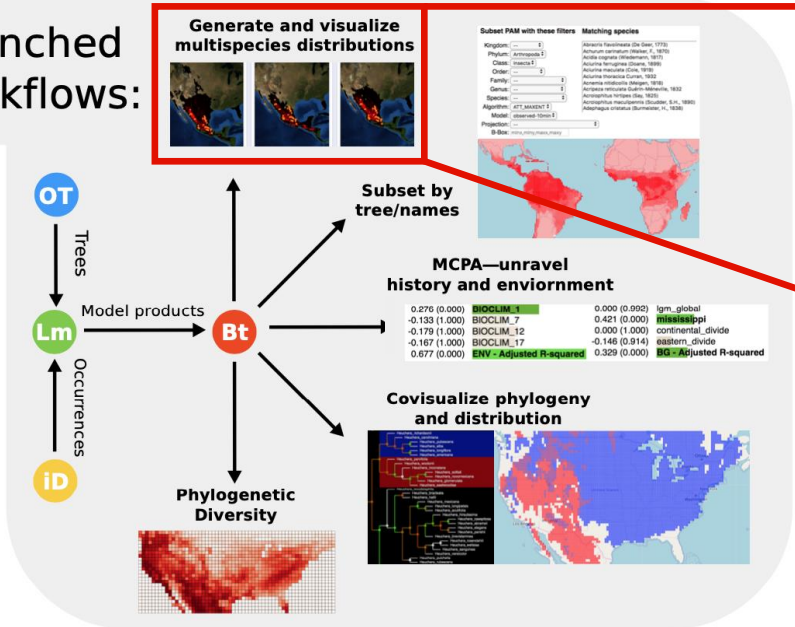
Phylogenetic Diversity



# BiotaPhy Two: Launched Workflows

Launched workflows:

- OT** Open Tree of Life
  - phylogenies
  - taxonomy/names
- Lm** Lifemapper
  - Scaleable ecological niche modeling
- iD** iDigBio
  - specimen data/images
  - trait data
  - fossil data/images
- Bt** BiotaPhy
  - spatial analyses
  - evolutionary reconstruction
  - visualization

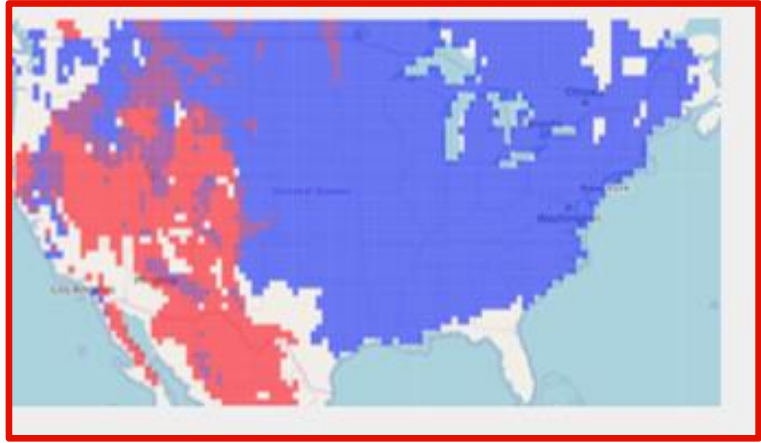
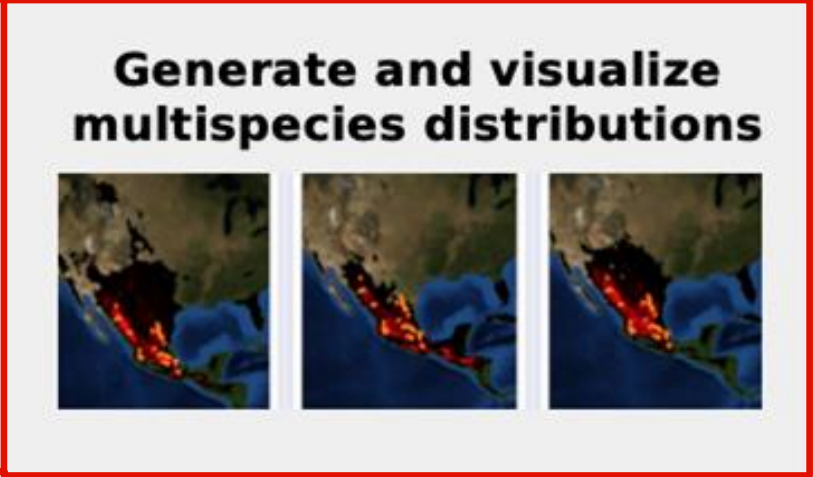


Subset PAM with these filters

Matching species	Matching species
Kingdom: Eukaryota	Artemisia tridentata (Gray) 1819
Phylum: Arthropoda	Artemisia tridentata (Gray) 1819
Class: Insecta	Artemisia tridentata (Gray) 1819
Order: Diptera	Artemisia tridentata (Gray) 1819
Family: Cecidomyiidae	Artemisia tridentata (Gray) 1819
Genus: Cecidomyia	Artemisia tridentata (Gray) 1819
Species: Cecidomyia artemisiae (Clausen, 1932)	Artemisia tridentata (Gray) 1819
Alphabet: A-Z	Artemisia tridentata (Gray) 1819
Model: observed (100%)	Artemisia tridentata (Gray) 1819
Filter: none	Artemisia tridentata (Gray) 1819
B-Box: none	Artemisia tridentata (Gray) 1819

MCPA—unravel history and environment

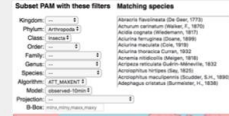
0.276 (0.000)	BIOCLIM_1	0.000 (0.002)	lgn_global
-0.133 (1.000)	BIOCLIM_7	0.421 (0.000)	mississippi
-0.179 (1.000)	BIOCLIM_12	0.000 (1.000)	continental_divide
-0.167 (1.000)	BIOCLIM_17	-0.146 (0.914)	eastern_divide
0.677 (0.000)	ENV - Adjusted R-squared	0.329 (0.000)	BG - Adjusted R-squared



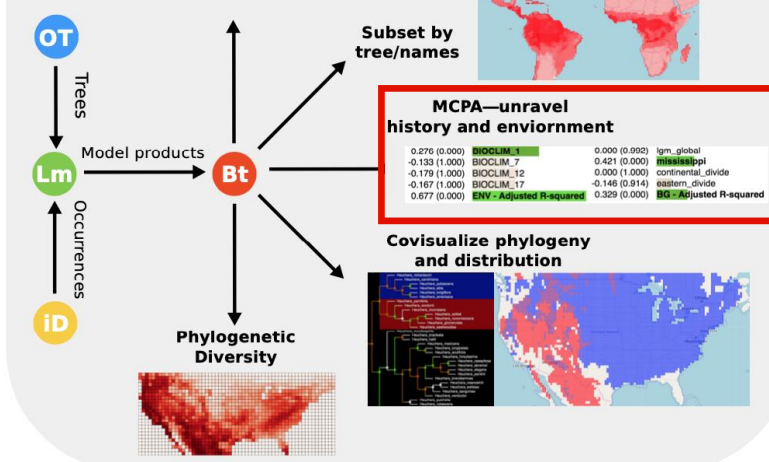
Saxifragaceae

# BiotaPhy Two: Launched Workflows

Launched workflows:



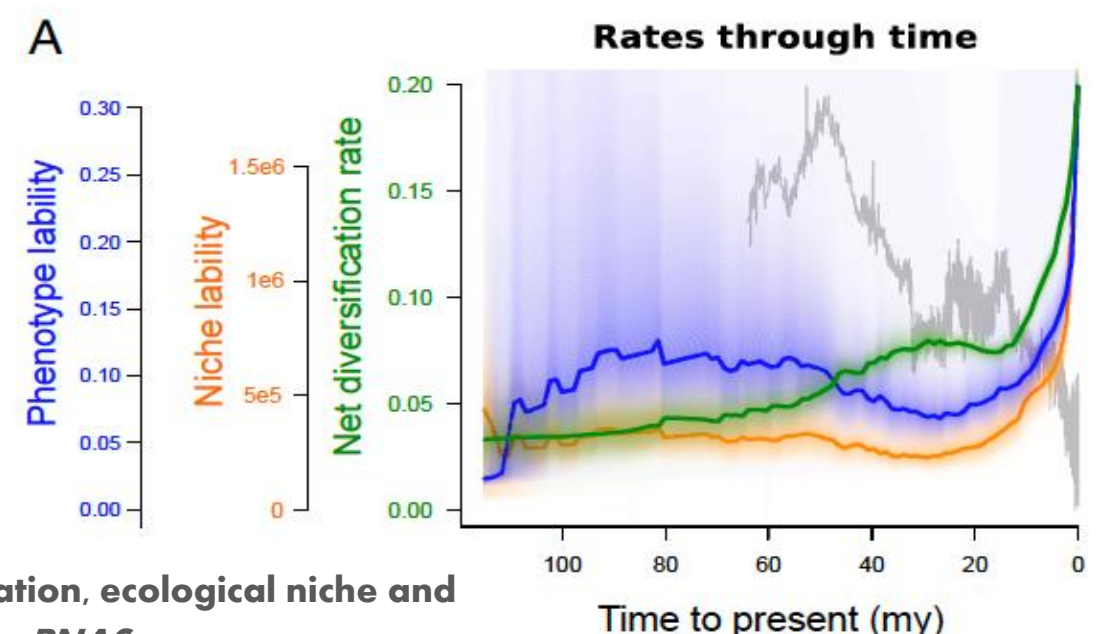
- OT** Open Tree of Life
  - phylogenies
  - taxonomy/names
- Lm** Lifemapper
  - Scaleable ecological niche modeling
- iD** iDigBio
  - specimen data/images
  - trait data
  - fossil data/images
- Bt** BiotaPhy
  - spatial analyses
  - evolutionary reconstruction
  - visualization



### MCPA—unravel history and environment

0.276 (0.000)	<b>BIOCLIM_1</b>	0.000 (0.992)	lgm_global
-0.133 (1.000)	BIOCLIM_7	0.421 (0.000)	<b>mississippi</b>
-0.179 (1.000)	BIOCLIM_12	0.000 (1.000)	continental_divide
-0.167 (1.000)	BIOCLIM_17	-0.146 (0.914)	eastern_divide
0.677 (0.000)	<b>ENV - Adjusted R-squared</b>	0.329 (0.000)	<b>BG - Adjusted R-squared</b>

## Meta-Community Phylogenetic Analysis Test the role of biogeography and environmental niche in the diversity of Saxifragales



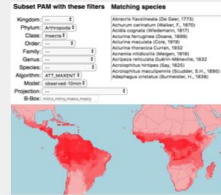
Folk et al.. 2019. Uncoupled diversification, ecological niche and phenotype in a temperate radiation. *PNAS* 116: 10874-10882.

# BiotaPhy Two: Launched Workflows

Launched workflows:



Generate and visualize multispecies distributions



Subset by tree/names

MCPA—unravel history and environment

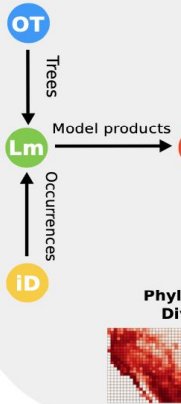
0.276 (0.000)	BIOCLIM_1	0.000 (0.000)	lgm_global
-0.133 (1.000)	BIOCLIM_7	0.421 (0.000)	maxaltpp1
-0.179 (1.000)	BIOCLIM_12	0.000 (1.000)	continental_divide
-0.167 (1.000)	BIOCLIM_17	-0.146 (0.914)	eastern_divide
0.677 (0.000)	ENV_Adjusted R-squared	0.329 (0.000)	RG_Adjusted R-squared

Covisualize phylogeny and distribution

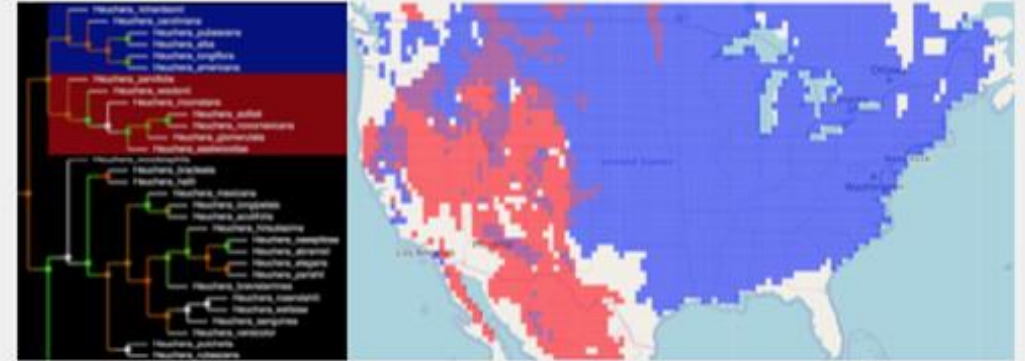


Phylogenetic Diversity

- OT** Open Tree of Life
  - phylogenies
  - taxonomy/names
- Lm** Lifemapper
  - Scaleable ecological niche modeling
- iD** iDigBio
  - specimen data/images
  - trait data
  - fossil data/images
- Bt** BiotaPhy
  - spatial analyses
  - evolutionary reconstruction
  - visualization

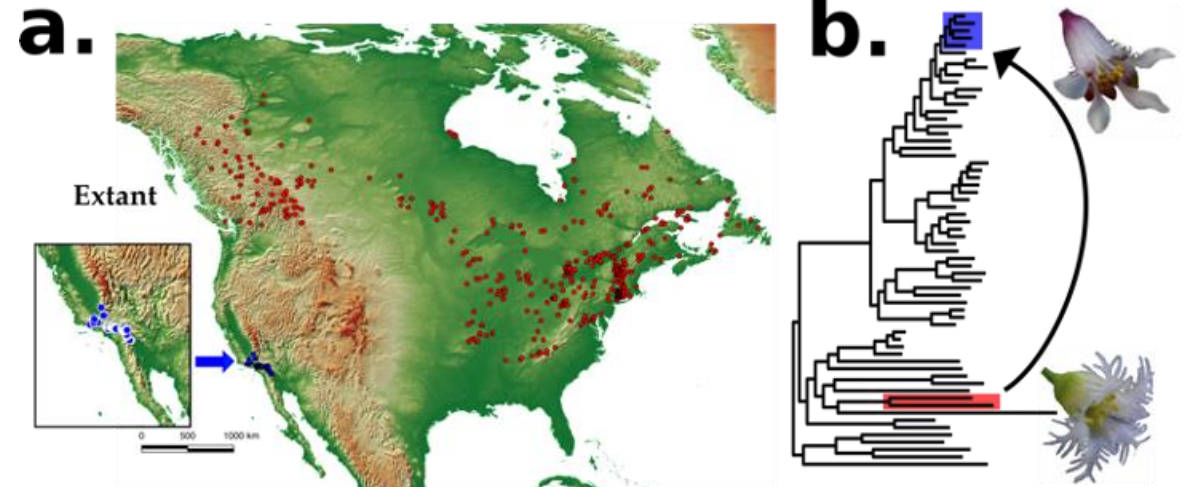


Covisualize phylogeny and distribution



Folk et al.. 2018. Assessing ancestral niche suitability and geographic range dynamics as drivers of hybridization in *Heuchera* (Saxifragaceae). *American Naturalist*

192: 171–187.

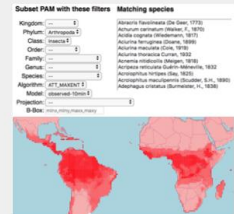


# BiotaPhy Two: Launched Workflows

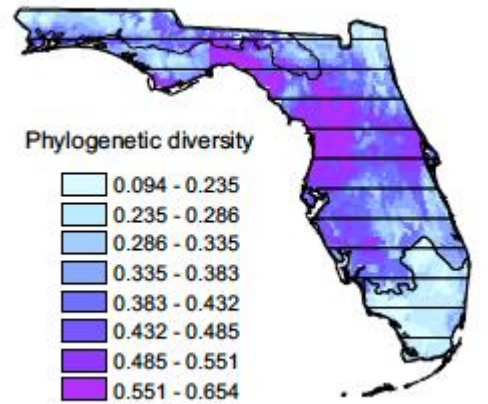
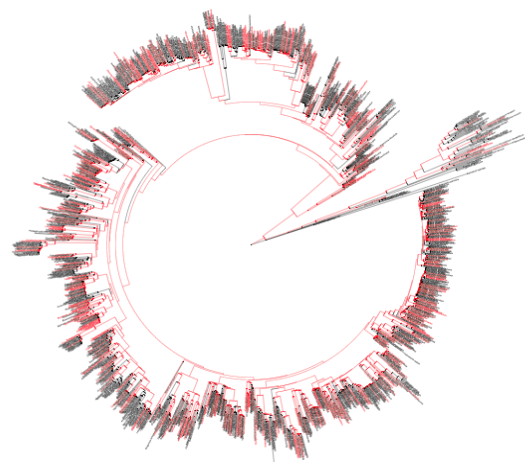
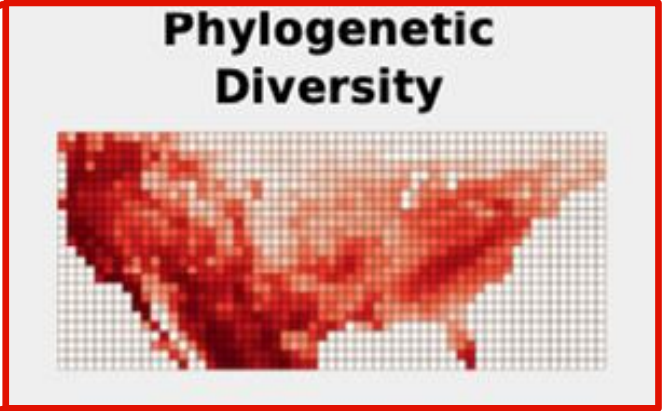
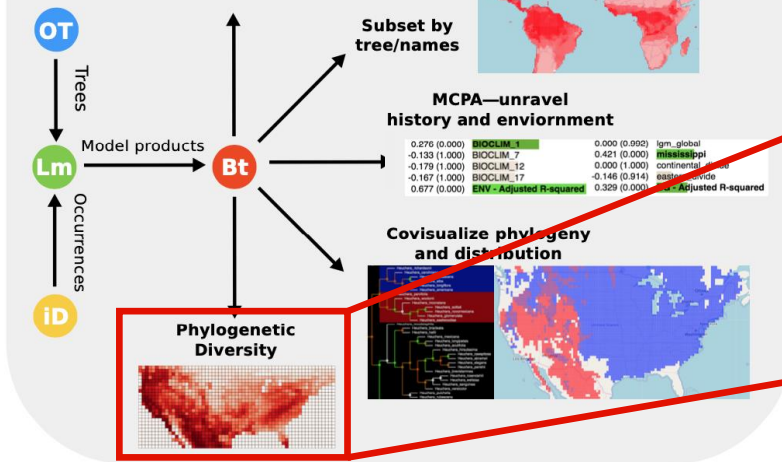


Launched workflows:

Generate and visualize multispecies distributions

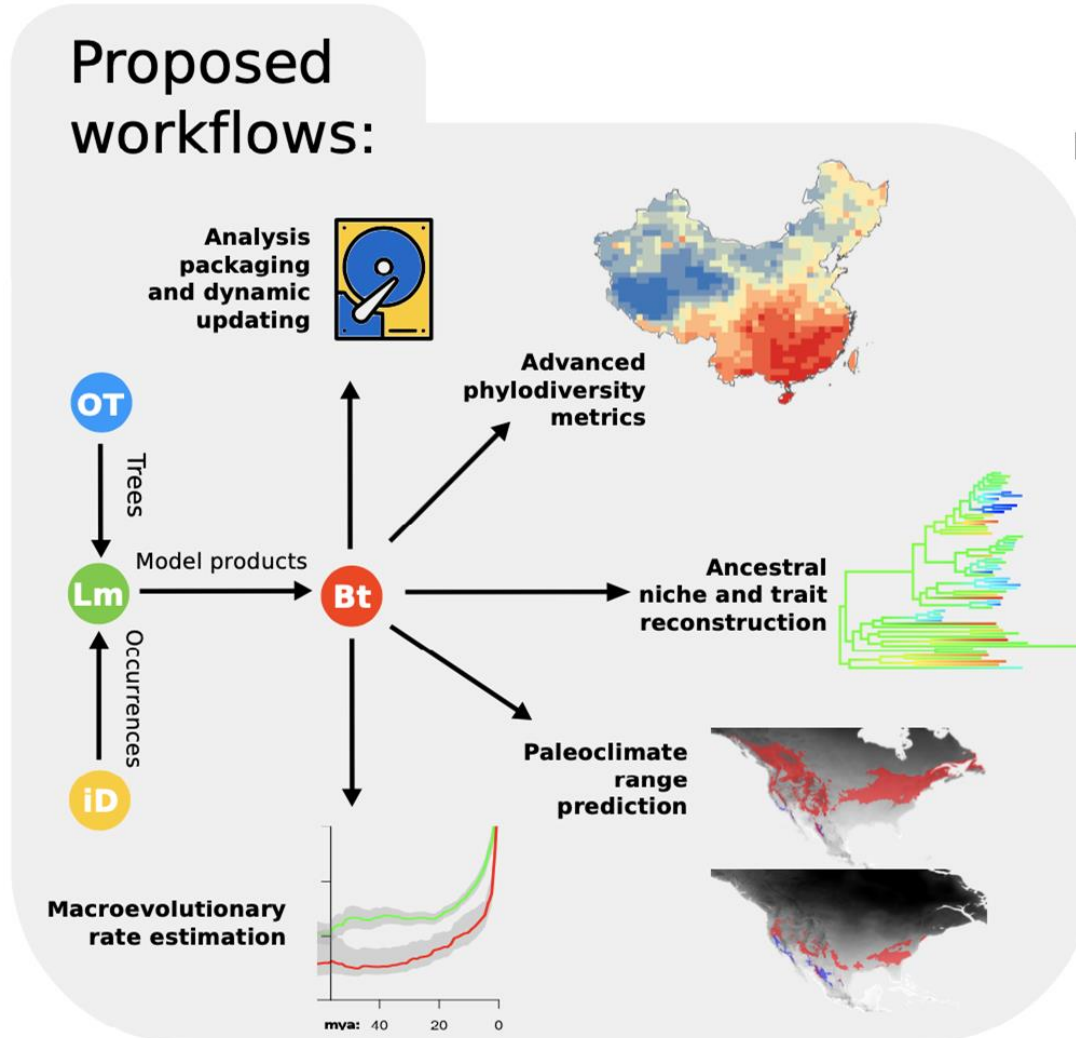


- OT** Open Tree of Life
  - phylogenies
  - taxonomy/names
- Lm** Lifemapper
  - Scaleable ecological niche modeling
- iD** iDigBio
  - specimen data/images
  - trait data
  - fossil data/images
- Bt** BiotaPhy
  - spatial analyses
  - evolutionary reconstruction
  - visualization



Allen, J., et al. 2019. Spatial phylogenetics of Florida vascular plants: The effects of tree uncertainty and ultrametricity. *iScience* 11: 57–70  
<https://doi.org/10.1016/j.isci.2018.12.002>

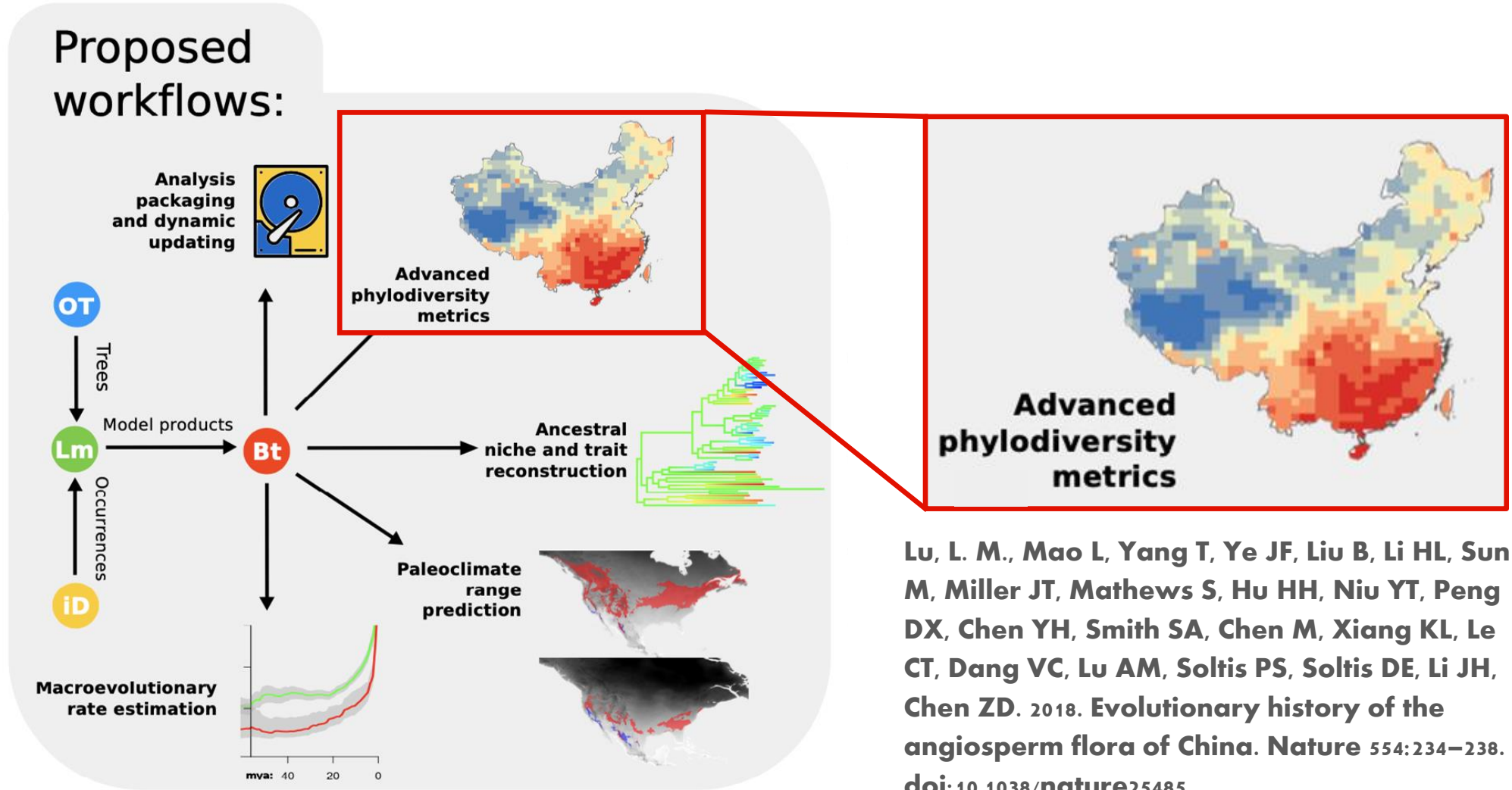
## Proposed workflows:



## New Approaches

- Trait reconstructions, correlations with ecological factors
- Ancestral niche comparisons
- Conservation of ecological niche through time
- Predictions of ancestral suitable habitat under paleoclimate scenarios
- Community divergence times and other advanced community metrics
- Species diversification rates, trait-associated diversification
- Niche evolutionary rates

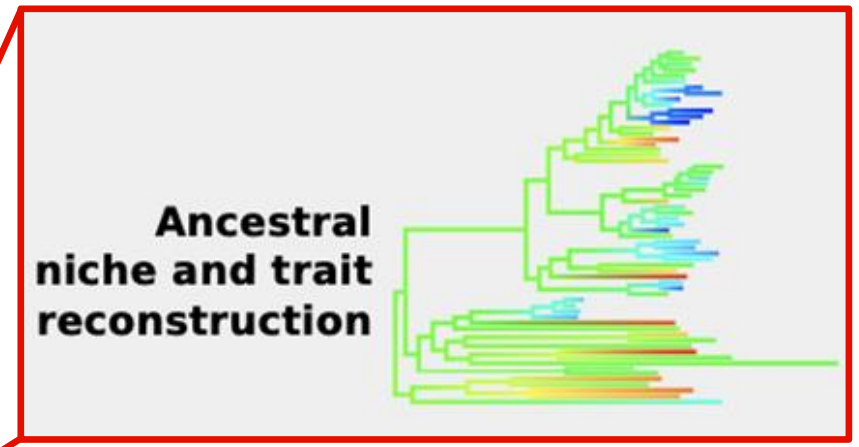
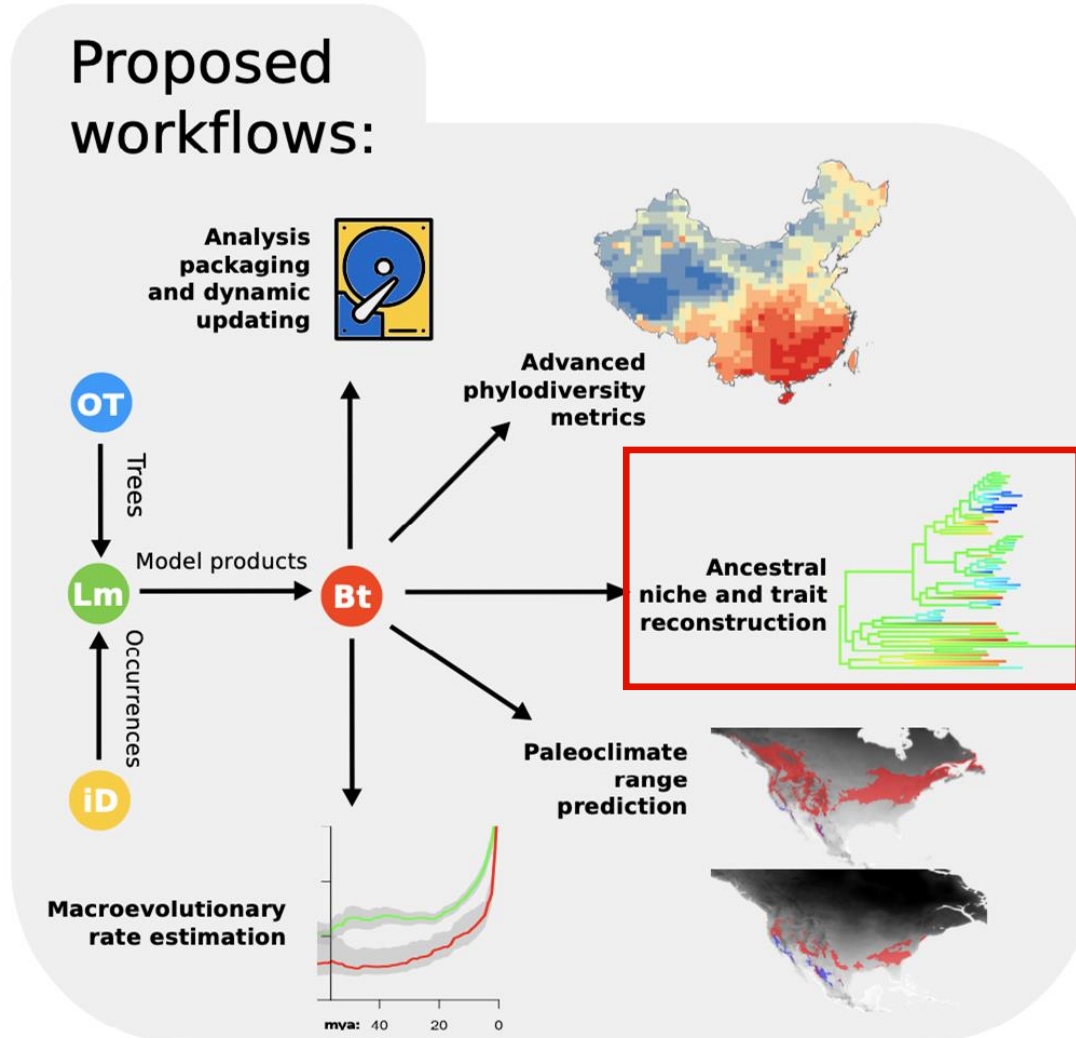
# BiotaPhy Two: Proposed Workflows



Lu, L. M., Mao L, Yang T, Ye JF, Liu B, Li HL, Sun M, Miller JT, Mathews S, Hu HH, Niu YT, Peng DX, Chen YH, Smith SA, Chen M, Xiang KL, Le CT, Dang VC, Lu AM, Soltis PS, Soltis DE, Li JH, Chen ZD. 2018. Evolutionary history of the angiosperm flora of China. *Nature* 554:234–238. doi:10.1038/nature25485.

# BiotaPhy Two: Proposed Workflows

## Proposed workflows:

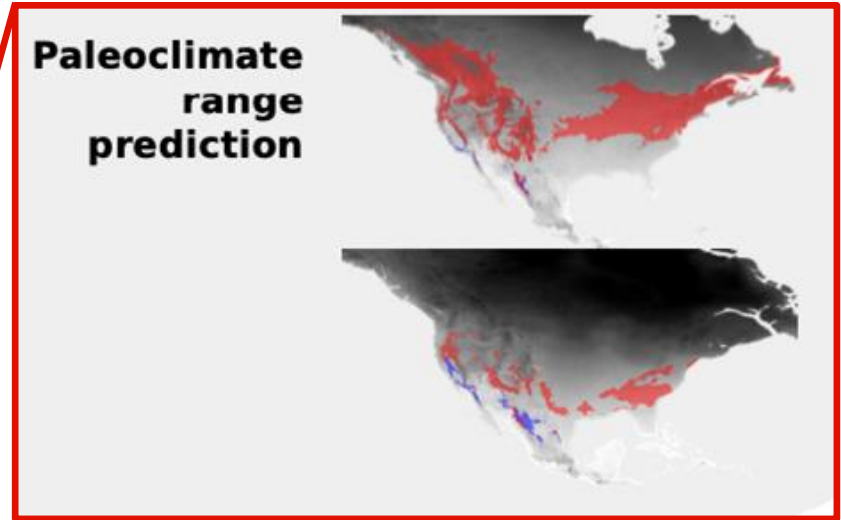
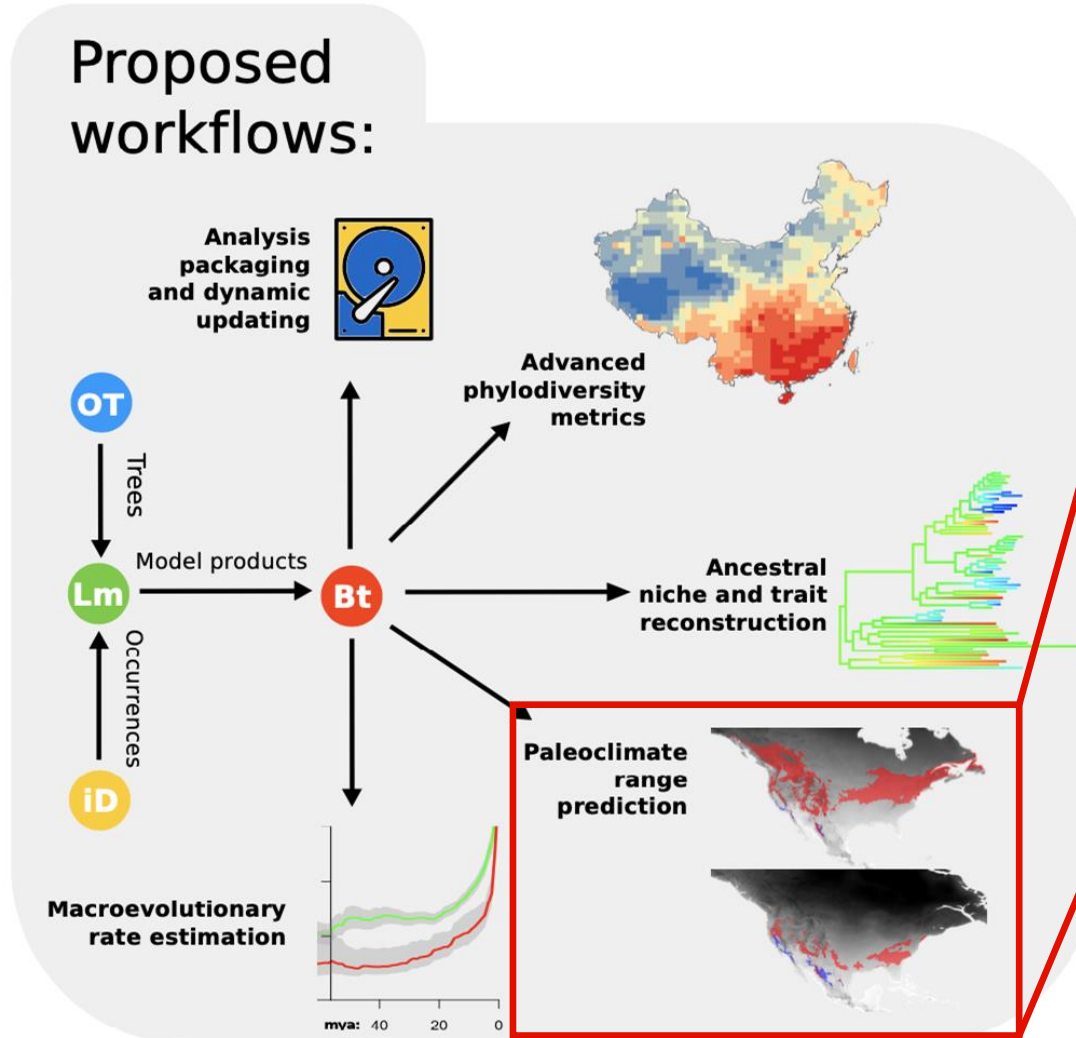


Folk et al.. 2018. Assessing ancestral niche suitability and geographic range dynamics as drivers of hybridization in *Heuchera* (Saxifragaceae). *American Naturalist* 192: 171-187.



# BiotaPhy Two: Proposed Workflows

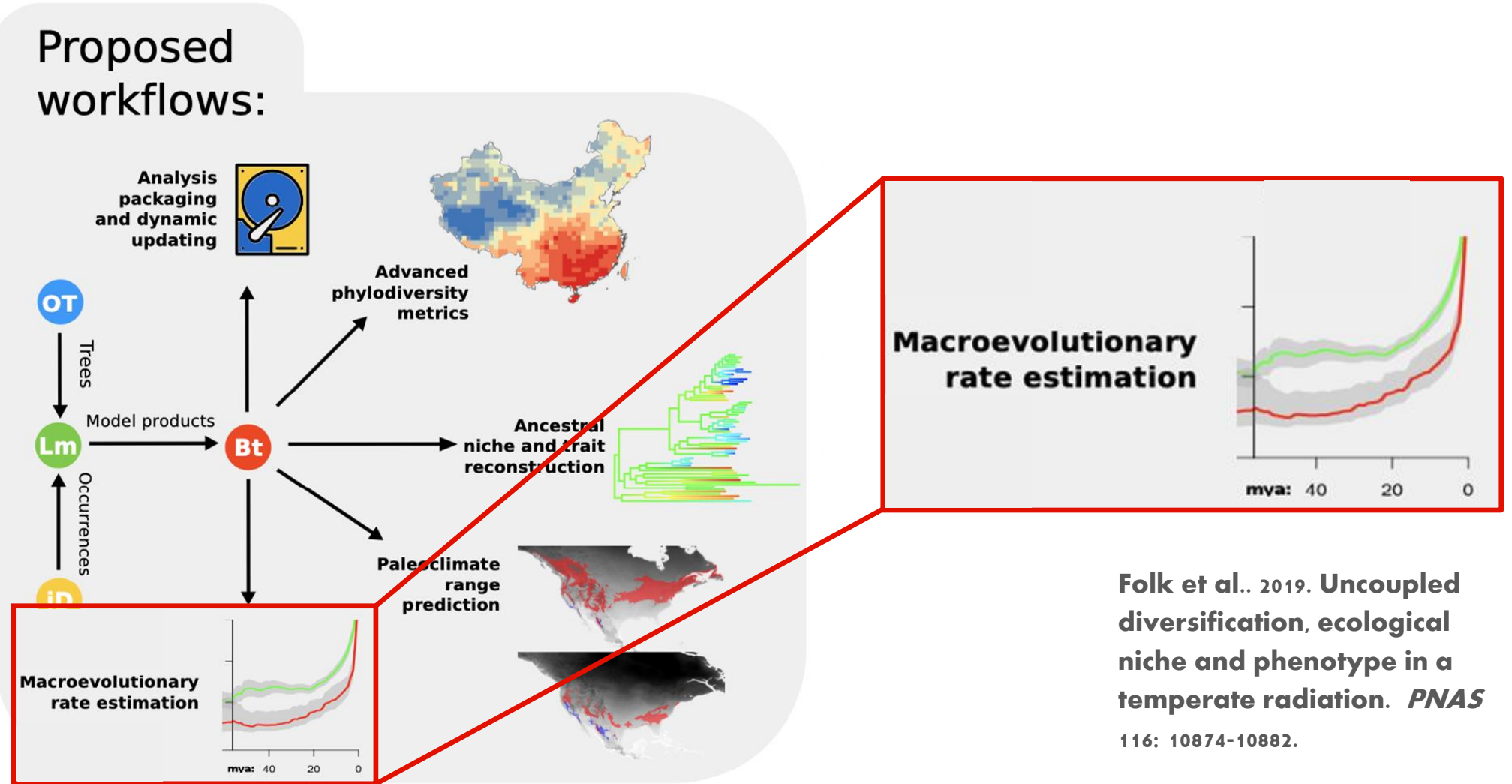
## Proposed workflows:



Folk et al.. 2018. Assessing ancestral niche suitability and geographic range dynamics as drivers of hybridization in *Heuchera* (Saxifragaceae). *American Naturalist* 192: 171-187.

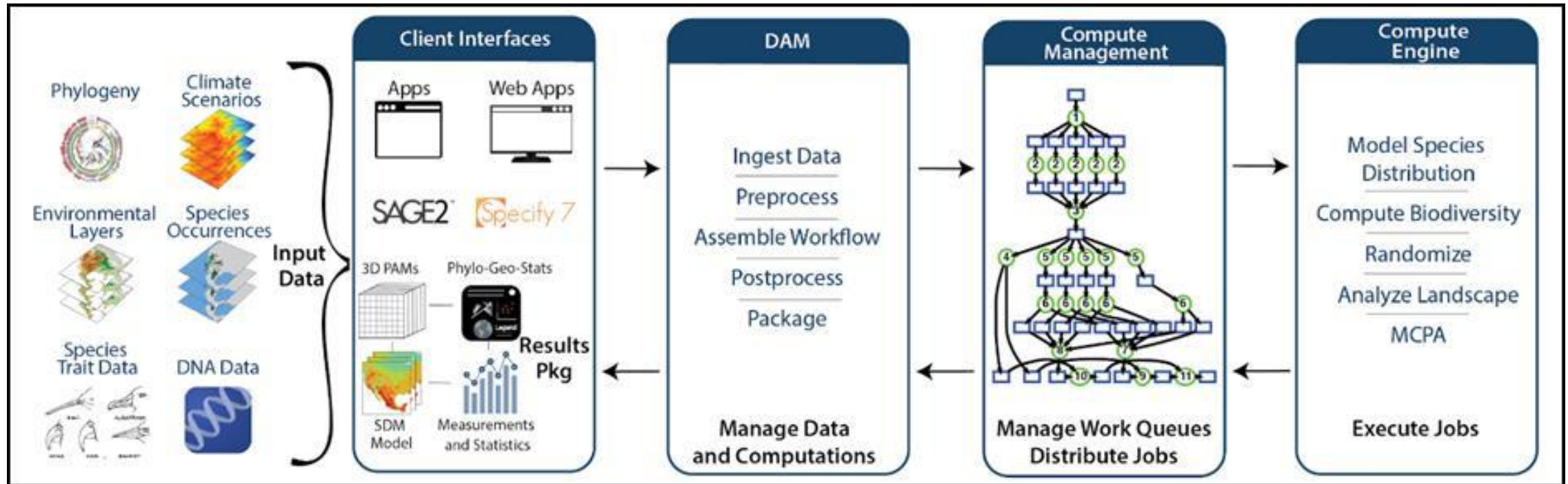
# BiotaPhy Two: Proposed Workflows

Proposed workflows:



Folk et al.. 2019. Uncoupled diversification, ecological niche and phenotype in a temperate radiation. *PNAS* 116: 10874-10882.

# Technical Scope: BiotaPhy's architecture



## Technology Context and Goals:

- ✓ **Scaling: geographically, taxonomically, phylogenetically**
- ✓ **Community gateway**
- ✓ **High throughput networks**
- ✓ **Parallelization by species**
- ✓ **Algorithms for creating and analyzing matrices, p/a and other matrices**

- ✓ **Dedicated workflow environments, VisTrails and Kepler in earlier projects**
- ✓ **Monolithic engine with some support for APIS, Web client, CC Tools**
- ✓ **Modules with stronger emphasis on APIs and libraries**
- ✓ **Docker modules connected by simple scripts**

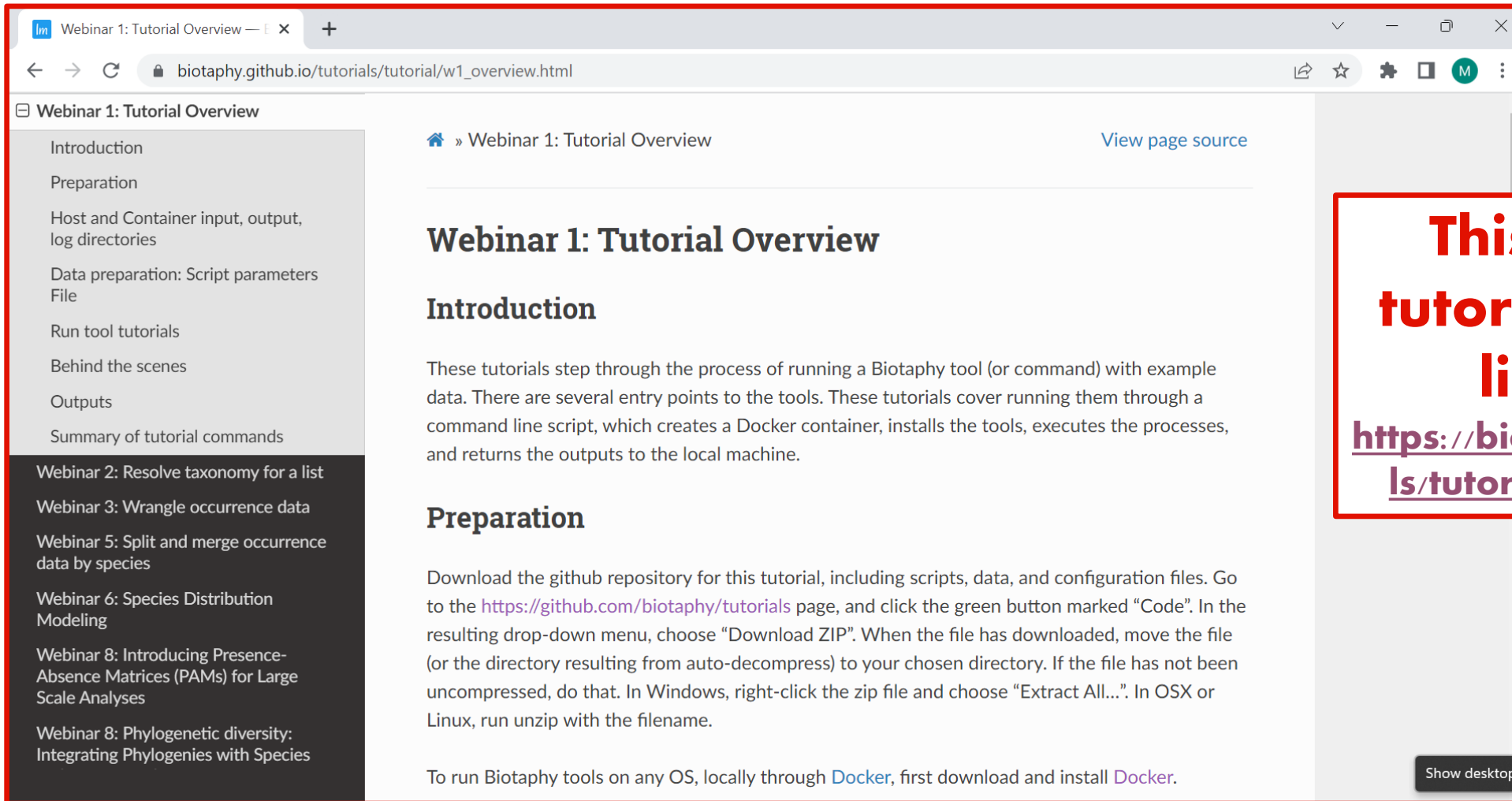
## Present and Future Work:

- ✓ **Projects, recent and ongoing with the BiotaPhy biologists**
- ✓ **Specify Software integration**
- ✓ **Fomenting new collaborations**

## Present and Future Approaches for Computing with Species

### Occurrence Data:

- ✓ **Software tools here**
- ✓ **Future of big data integrative processing**
- ✓ **Digital Object Architecture**
- ✓ **Invitation for feedback and collaboration**
- ✓ **Full documentation at: [biotaphy.github.io/tutorials](https://biotaphy.github.io/tutorials)**



Webinar 1: Tutorial Overview

- Introduction
- Preparation
- Host and Container input, output, log directories
- Data preparation: Script parameters File
- Run tool tutorials
- Behind the scenes
- Outputs
- Summary of tutorial commands

Webinar 2: Resolve taxonomy for a list

Webinar 3: Wrangle occurrence data

Webinar 5: Split and merge occurrence data by species

Webinar 6: Species Distribution Modeling

Webinar 8: Introducing Presence-Absence Matrices (PAMs) for Large Scale Analyses

Webinar 8: Phylogenetic diversity: Integrating Phylogenies with Species

## Webinar 1: Tutorial Overview

### Introduction

These tutorials step through the process of running a Biotaphy tool (or command) with example data. There are several entry points to the tools. These tutorials cover running them through a command line script, which creates a Docker container, installs the tools, executes the processes, and returns the outputs to the local machine.

### Preparation

Download the github repository for this tutorial, including scripts, data, and configuration files. Go to the <https://github.com/biotaphy/tutorials> page, and click the green button marked "Code". In the resulting drop-down menu, choose "Download ZIP". When the file has downloaded, move the file (or the directory resulting from auto-decompress) to your chosen directory. If the file has not been uncompressed, do that. In Windows, right-click the zip file and choose "Extract All...". In OSX or Linux, run unzip with the filename.

To run Biotaphy tools on any OS, locally through [Docker](#), first download and install [Docker](#).

**This is what the tutorial page looks like! Access:**  
[https://biotaphy.github.io/tutorials/tutorial/w1\\_overview.html](https://biotaphy.github.io/tutorials/tutorial/w1_overview.html)



## Webinar 1: Tutorial Overview

### Introduction

These tutorials step through the process of running a Biotaphy tool (or command) with example data. There are several entry points to the tools. These tutorials cover running them through a command line script, which creates a Docker container, installs the tools, executes the processes, and returns the outputs to the local machine.

✓ **Several entry points**



## Webinar 1: Tutorial Overview

### Introduction

These tutorials step through the process of running a Biotaphy tool (or command) with example data. There are several entry points to the tools. These tutorials cover running them through a command line script, which creates a Docker container, installs the tools, executes the processes, and returns the outputs to the local machine.

- ✓ **Several entry points**
- ✓ **Command line based**

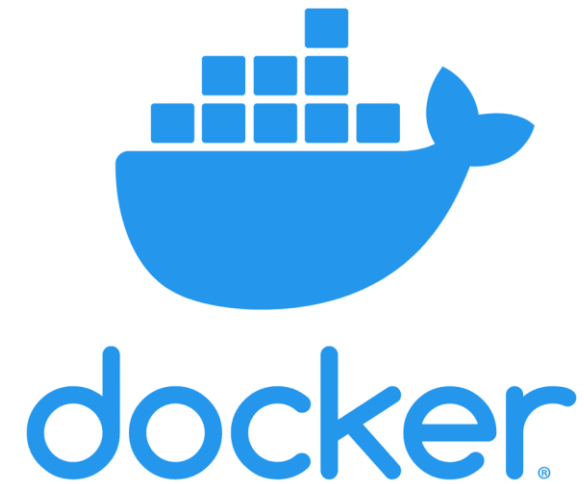
```
[mariacortez@login6 mariacortez]$ ls
ancient_angiosperms      campos_rupestres.sqlite  Ocbil
campos_rupestres         favorite_plant_project   Open_Tree_9.1
campos_rupestres.33862038.out  MAFFT                    OpenTree_SSB2020
[mariacortez@login6 mariacortez]$ module load python
[mariacortez@login6 mariacortez]$ cd campos_rupestres/
[mariacortez@login6 campos_rupestres]$ ls
mata_atlantica.R  OpenTree  paththo  SQLite  SQLite_2  SQLite_Jan_2022  SQL
[mariacortez@login6 campos_rupestres]$ cd SQLite_Jan_2022/
[mariacortez@login6 SQLite_Jan_2022]$ la
-bash: la: command not found
[mariacortez@login6 SQLite_Jan_2022]$ ls
campos_rupestres_2.sqlite  campos_rupestres_4.py      campos_rupestre
campos_rupestres_3.py     campos_rupestres_4_utf.txt  campos_rupestre
campos_rupestres_3.sqlite  campos_rupestres_jan_22.sqlite  distribution_me
[mariacortez@login6 SQLite_Jan_2022]$ cd campos_rupestres_SQLite/
[mariacortez@login6 campos_rupestres_SQLite]$ ls
```

## Webinar 1: Tutorial Overview

### Introduction

These tutorials step through the process of running a Biotaphy tool (or command) with example data. There are several entry points to the tools. These tutorials cover running them through a command line script, which creates a Docker container, installs the tools, executes the processes, and returns the outputs to the local machine.

- ✓ **Several entry points**
- ✓ **Command line based**
- ✓ **Uses Docker**




## Webinar 1: Tutorial Overview

### Introduction

These tutorials step through the process of running a Biotaphy tool (or command) with example data. There are several entry points to the tools. These tutorials cover running them through a command line script, which creates a Docker container, installs the tools, executes the processes, and returns the outputs to the local machine.

- ✓ **Several entry points**
- ✓ **Command line based**
- ✓ **Uses Docker**
- ✓ **Outputs Saved in local machine**

main ▾ [tutorials](#) / [data](#) / [easy\\_bake](#) / [heuchera\\_accepted1.txt](#)

 zzeppozz clarity ✓

🔍 1 contributor

124 lines (124 sloc) | 2.32 KB

```
1 Heuchera puberula
2 Heuchera scapigera
3 Heuchera lucida
4 Heuchera hirtiflora
5 Heuchera barbarossa
6 Heuchera richardsonii
7 Heuchera hirsuticaulis
8 Heuchera wootonii
9 Heuchera foliosa
10 Heuchera duranii
```

# Technical Scope: Tutorial Overview | BiotaPhy

## Preparation

Download the github repository for this tutorial, including scripts, data, and configuration files. Go to the <https://github.com/biotaphy/tutorials> page, and click the green button marked "Code". In the resulting drop-down menu, choose "Download ZIP". When the file has downloaded, move the file (or the directory resulting from auto-decompress) to your chosen directory. If the file has not been uncompressed, do that. In Windows, right-click the zip file and choose "Extract All...". In OSX or Linux, run unzip with the filename.

The screenshot shows the GitHub interface for the repository 'biotaphy/tutorials'. The 'Code' button is highlighted in green, and its dropdown menu is open, showing options: 'Clone' (with sub-options for HTTPS and GitHub CLI), 'Open with GitHub Desktop', and 'Download ZIP'. A red arrow points from the text in the 'Preparation' section to the 'Code' button. Below the screenshot, a table lists the files and folders in the repository.

File/Folder	Description
zzeppozz link	
.github	added pre-commit checks; unfini
_sphinx_config	link
data	filename fixes
.gitignore	Merge branch 'main' of github.co
.pre-commit-config.yaml	updates
.readthedocs.yaml	start of sphinx tutorial

# Technical Scope: Tutorial Overview

To run Biotaphy tools on any OS, locally through Docker, first download and install Docker.

**Docker works similarly to a virtual machine, but more efficiently!!!!**

**Google 'Docker', access the website and download the appropriate version!**



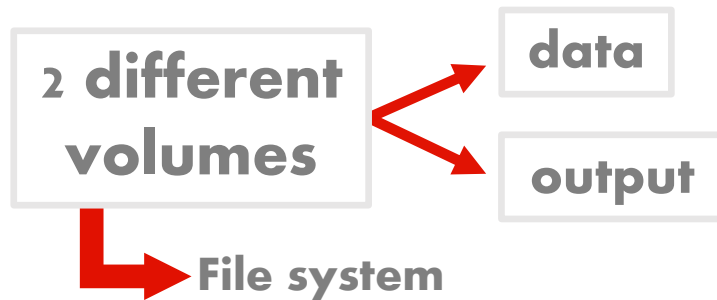
The screenshot shows the Docker documentation website. The navigation bar includes 'docker docs', a search bar, and links for 'Home', 'Guides', 'Manuals', 'Reference', and 'Samples'. The breadcrumb trail is 'Home / Guides / Get started / Part 1: Getting started'. The left sidebar lists various guides, with 'Part 1: Getting started' selected. The main content area is titled 'Download and install Docker' and contains the following text: 'This tutorial assumes you have a current version of Docker installed on your machine. If you do not have Docker installed, choose your preferred operating system below to download Docker:'. Below this text are four blue buttons: 'Mac with Intel chip', 'Mac with Apple chip', 'Windows', and 'Linux'. At the bottom, there is a section for 'For Docker Desktop installation instructions, see:' followed by a bulleted list: 'Install Docker Desktop on Mac', 'Install Docker Desktop on Windows', and 'Install Docker Desktop on Linux'.

## Host and Container input, output, log directories

A data **Docker Volume** is created on the [Host machine](#), and the `tutorials/data/input`, `tutorials/data/config`, and `tutorials/data/wrangers` directories in this repository are copied to it. These directories are then made available on the Docker [Container](#) under the `/volumes/data` volume (directory). If modified, the docker “data” volume must be re-created to propagate those changes to the containers.

Another **Docker volume**, *output*, is created on the [Host machine](#) and mounted at `/volumes/output` in the Docker [Container](#). Changes in this directory are saved in the volume, and copied back to the host machine, under the `data` directory.

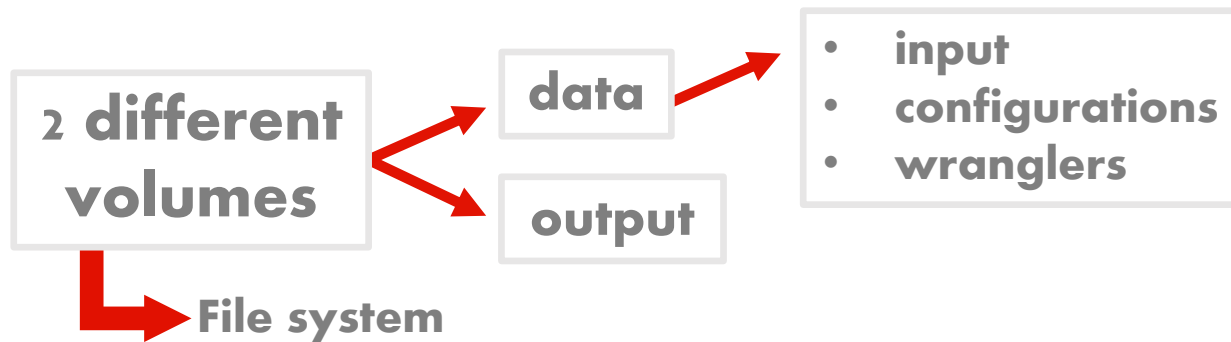
**Host machine**



## Host and Container input, output, log directories

A data **Docker Volume** is created on the [Host machine](#), and the `tutorials/data/input`, `tutorials/data/config`, and `tutorials/data/wrangers` directories in this repository are copied to it. These directories are then made available on the Docker [Container](#) under the `/volumes/data` volume (directory). If modified, the docker “data” volume must be re-created to propagate those changes to the containers.

Another **Docker volume**, `output`, is created on the [Host machine](#) and mounted at `/volumes/output` in the Docker [Container](#). Changes in this directory are saved in the volume, and copied back to the host machine, under the `data` directory.



**Host machine**



# Technical Scope: Tutorial Overview

## Host and Container input, output, log directories

A data Docker Volume, is created on the [Host machine](#), and the `tutorials/data/input`, `tutorials/data/config`, and `tutorials/data/wrangers` directories in this repository are copied to it. These directories are then made available on the [Docker Container](#) under the `/volumes/data` volume (directory). If modified, the docker “data” volume must be re-created to propagate those changes to the containers.

Another Docker volume, *output*, is created on the [Host machine](#) and mounted at `/volumes/output` in the Docker [Container](#). Changes in this directory are saved in the volume, and copied back to the host machine, under the data directory.



# Technical Scope: Tutorial Overview

## Host and Container input, output, log directories

A data Docker Volume, is created on the [Host machine](#), and the `tutorials/data/input`, `tutorials/data/config`, and `tutorials/data/wrangers` directories in this repository are copied to it. These directories are then made available on the [Docker Container](#) under the `/volumes/data` volume (directory). If modified, the docker “data” volume must be re-created to propagate those changes to the containers.

Another Docker volume, *output*, is created on the [Host machine](#) and mounted at `/volumes/output` in the Docker [Container](#). Changes in this directory are saved in the volume, and copied back to the host machine, under the data directory.



# Technical Scope: Tutorial Overview

## Host and Container input, output, log directories

A data Docker Volume, is created on the [Host machine](#), and the `tutorials/data/input`, `tutorials/data/config`, and `tutorials/data/wrangers` directories in this repository are copied to it. These directories are then made available on the [Docker Container](#) under the `/volumes/data` volume (directory). If modified, the docker “data” volume must be re-created to propagate those changes to the containers.

Another Docker volume, *output*, is created on the [Host machine](#) and mounted at `/volumes/output` in the Docker [Container](#). Changes in this directory are saved in the volume, and copied back to the host machine, under the data directory.



# Technical Scope: Tutorial Overview | BiotaPhy

Host machine – **DOCKER HAS TO BE INSTALLED!**

Run tutorial script:

```
$run_tutorial
```

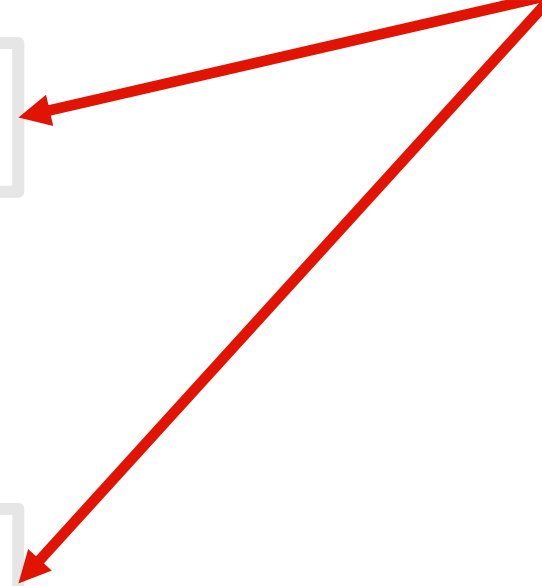
```
<command><configuration_file>
```

Empty Docker  
Volume

output

Empty Docker  
Volume

data



# Technical Scope: Tutorial Overview | BiotaPhy

Host machine – **DOCKER HAS TO BE INSTALLED!**

Run tutorial script:

```
$run_tutorial
```

```
<command><configuration_file>
```

Empty Docker  
Volume

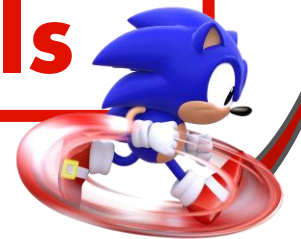
output

Empty Docker  
Volume

data

Docker  
Image

**BiotaPhy**  
**tools**



# Technical Scope: Tutorial Overview | BiotaPhy

Host machine – **DOCKER HAS TO BE INSTALLED!**

File system



Empty Docker  
Volume

output

Docker  
Volume

data

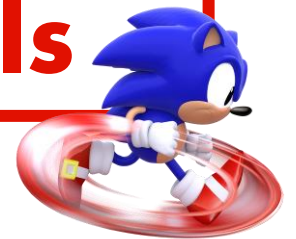
Run tutorial script:

```
$run_tutorial
```

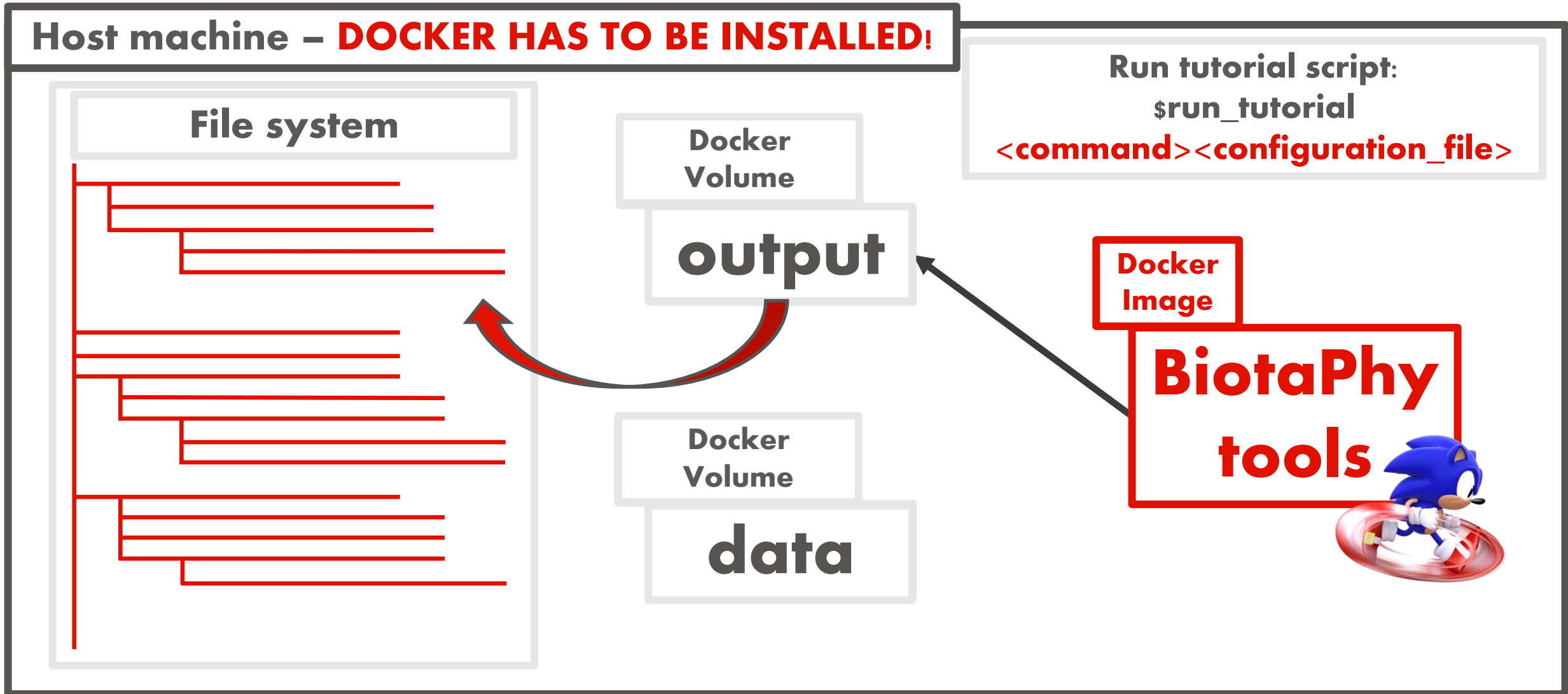
```
<command><configuration_file>
```

Docker  
Image

**BiotaPhy  
tools**



# Technical Scope: Tutorial Overview | BiotaPhy



## Data preparation: Script parameters File

All commands require a parameters file with tool-specific parameters. The file must be in [JSON](#) format. More information is [\[here\]](#)(script\_params.rst).

Each tutorial contains one or more example parameters files in the tutorials/data/config directory. These parameters files reference example data and parameters reasonable for that command. All required and optional parameters are described in individual tutorial pages.

Some tools will require an additional [JSON](#) format configuration file. In these cases, the additional JSON file will be named in the parameters file.

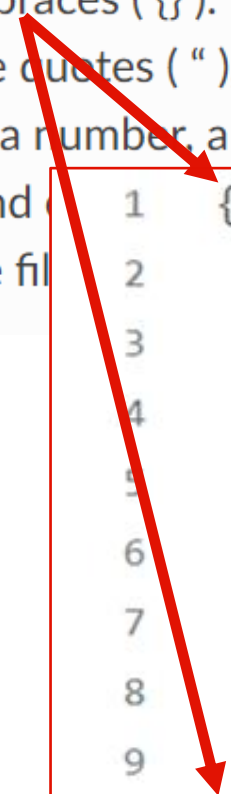
- ✓ **Parameter files are essential to properly run the commands!**
- ✓ **Stored in the data/config directory**



## Script parameters file

1. Each command has required and/or optional parameters.
2. The group of options are enclosed in curly braces ( {} ).
3. Each option keyword is quoted with double quotes ( " ), is followed by a colon ( : ) and a value.
4. Each value may be a double-quoted string, a number, a boolean (true or false, not quoted), or a list of these items, separated by commas and
5. Output filenames must be a full path to the file

```
1  {
2      "max_open_writers": 100,
3      "key_field_ignore": ["genus", "scientificName"],
4      "dwca": [
5          ["/volumes/data/input/occ_heuchera_gbif.zip",
6            "/volumes/data/wrangers/no_wrangers.json"
7          ]
8      ],
9      "out_dir": "/volumes/output/heuchera_dwca"
10 }
```

Two red arrows originate from the list of instructions. One arrow points from instruction 2 to the opening curly brace on line 1 of the code block. The other arrow points from instruction 3 to the closing curly brace on line 10 of the code block.

## Script parameters file

1. Each command has required and/or optional parameters.
2. The group of options are enclosed in curly braces ( {} ).
3. Each option keyword is quoted with double quotes ( " ), is followed by a colon ( : ) and a value.
4. Each value may be a double-quoted string, a number, a boolean (true or false, not quoted), or a list of these items, separated by commas and
5. Output filenames must be a full path to the file

```
1  {
2    "max_open_writers": 100,
3    "key_field_ignore": ["genus","scientificName"],
4    "dwca": [
5      ["/volumes/data/input/occ_heuchera_gbif.zip",
6        "/volumes/data/wrangers/no_wrangers.json"
7      ]
8    ],
9    "out_dir": "/volumes/output/heuchera_dwca"
10 }
```

## Script parameters file

1. Each command has required and/or optional parameters.
2. The group of options are enclosed in curly braces ( {} ).
3. Each option keyword is quoted with double quotes ( " ), is followed by a colon ( : ) and a value.
4. Each value may be a double-quoted string, a number, a boolean (true or false, not quoted), or a list of these items, separated by commas and
5. Output filenames must be a full path to the file

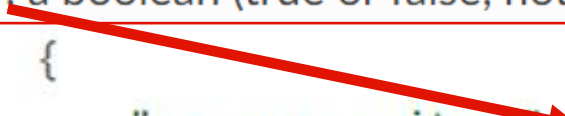
```
1  {
2    "max_open_writers": 100,
3    "key_field_ignore": ["genus","scientificName"],
4    "dwca": [
5        ["/volumes/data/input/occ_heuchera_gbif.zip",
6          "/volumes/data/wrangers/no_wrangers.json"
7        ]
8    ],
9    "out_dir": "/volumes/output/heuchera_dwca"
10 }
```



## Script parameters file

1. Each command has required and/or optional parameters.
2. The group of options are enclosed in curly braces ( {} ).
3. Each option keyword is quoted with double quotes ( " ), is followed by a colon ( : ) and a value.
4. Each value may be a double-quoted string, a number, a boolean (true or false, not quoted), or a list of these items, separated by commas and
5. Output filenames must be a full path to the file

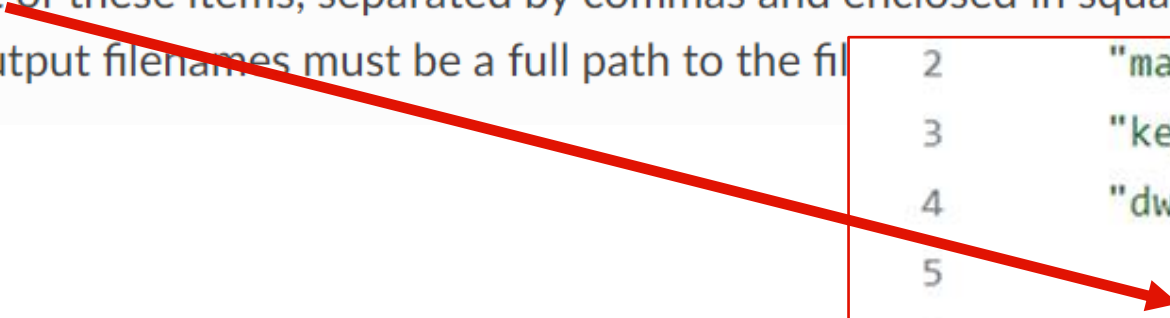
```
1  {
2    "max_open_writers": 100,
3    "key_field_ignore": ["genus", "scientificName"],
4    "dwca": [
5      "/volumes/data/input/occ_heuchera_gbif.zip",
6      "/volumes/data/wrangers/no_wrangers.json"
7    ]
8  },
9    "out_dir": "/volumes/output/heuchera_dwca"
10 }
```



## Script parameters file

1. Each command has required and/or optional parameters.
2. The group of options are enclosed in curly braces ( {} ).
3. Each option keyword is quoted with double quotes ( " ), is followed by a colon ( : ) and a value.
4. Each value may be a double-quoted string, a number, a boolean (true or false, not quoted), or a list of these items, separated by commas and enclosed in square brackets ( [] ).
5. Output filenames must be a full path to the file.

```
2     "max_open_writers": 100,  
3     "key_field_ignore": ["genus","scientificName"],  
4     "dwca": [  
5         ["/volumes/data/input/occ_heuchera_gbif.zip",  
6           "/volumes/data/wrangers/no_wrangers.json"  
7         ]  
8     ],  
9     "out_dir": "/volumes/output/heuchera_dwca"  
10 }
```

A red arrow originates from the text in rule 4 of the list above and points to the list of file paths within the "dwca" parameter in the code block.

## Script parameters file

1. Each command has required and/or optional parameters.
2. The group of options are enclosed in curly braces ( {} ).
3. Each option keyword is quoted with double quotes ( " ), is followed by a colon ( : ) and a value.
4. Each value may be a double-quoted string, a number, a boolean (true or false, not quoted), or a list of these items, separated by commas and enclosed in square brackets ( [] ).
5. Output filenames must be a full path to the file (or directory containing files)

A red arrow originates from the fifth item in the list above and points to the "out\_dir" parameter in the code block below.

```
8     ],
9     "out_dir": "/volumes/output/heuchera_dwca"
10  }
```

# Technical Scope: Tutorial Overview

- log\_console: this is a boolean value. The value *true*, causes the command to print logging lines in the command line window, to show the processes and progress.
- log\_filename: a full path to the output log file to be created. The file is an *Output filename* as described above. It contains all logging output from the process, and may be useful for identifying what processes were executed and their outcomes.
- report\_filename: a full path to the report log file to be created. The file is an *Output filename* as described above. It contains a summary of the modifications made to the output data, and may be useful for quantifying or comparing them.

**Continuing with  
Script parameters  
file! Optional  
parameters**

```
8 lines (8 slots) | 386 Bytes
1  {
2    "log_filename": "/volumes/output/wrangle_species_list1.log",
3    "log_console": true,
4    "report_filename": "/volumes/output/wrangle_species_list1_rpt.json",
5    "in_species_list_filename": "/volumes/data/input/heuchera.txt",
6    "wrangler_configuration_file": "/volumes/data/wrangers/splist_wranglers_gbif.json",
7    "out_species_list_filename": "/volumes/output/heuchera_accepted1.txt"
8  }
```

- `log_console`: this is a boolean value. The value `true`, causes the command to print logging lines in the command line window, to show the processes and progress.
- `log_filename`: a full path to the output log file to be created. The file is an *Output filename* as described above. It contains all logging output from the process, and may be useful for identifying what **processes** were executed and their outcomes.
- `report_filename`: a full path to the report log file to be created. The file is an *Output filename* as described above. It contains a summary of the modifications made to the output data, and may be useful for quantifying or comparing them.

**Continuing with  
Script parameters  
file! Optional  
parameters**

```
8 lines (8 sloc) | 386 bytes
1 {
2   "log_filename": "/volumes/output/wrangle_species_list1.log",
3   "log_console": true,
4   "report_filename": "/volumes/output/wrangle_species_list1_rpt.json",
5   "in_species_list_filename": "/volumes/data/input/heuchera.txt",
6   "wrangler_configuration_file": "/volumes/data/wrangers/splist_wranglers_gbif.json",
7   "out_species_list_filename": "/volumes/output/heuchera_accepted1.txt"
8 }
```



- `log_console`: this is a boolean value. The value `true`, causes the command to print logging lines in the command line window, to show the processes and progress.
- `log_filename`: a full path to the output log file to be created. The file is an *Output filename* as described above. It contains all logging output from the process, and may be useful for identifying what processes were executed and their outcomes.
- `report_filename`: a full path to the report log file to be created. The file is an *Output filename* as described above. It contains a summary of the modifications made to the output data, and may be useful for quantifying or comparing them.

**Continuing with  
Script parameters  
file! Optional  
parameters**

```
8 lines (8 sloc) | 386 Bytes
1  {
2  "log_filename": "/volumes/output/wrangle_species_list1.log",
3  "log_console": true,
4  "report_filename": "/volumes/output/wrangle_species_list1_rpt.json",
5  "in_species_list_filename": "/volumes/data/input/heuchera.txt",
6  "wrangler_configuration_file": "/volumes/data/wrangers/splist_wranglers_gbif.json",
7  "out_species_list_filename": "/volumes/output/heuchera_accepted1.txt"
8  }
```

- `log_console`: this is a boolean value. The value `true`, causes the command to print logging lines in the command line window, to show the processes and progress.
- `log_filename`: a full path to the output log file to be created. The file is an *Output filename* as described above. It contains all logging output from the process, and may be useful for identifying what processes were executed and their outcomes.
- `report_filename`: a full path to the report log file to be created. The file is an *Output filename* as described above. It contains a summary of the modifications made to the output data, and may be useful for quantifying or comparing them.

**We will learn more about wranglers later in the series!**

8 lines (8 sloc) | 386 Bytes

```
1 {
2   "log_filename": "/volumes/output/wrangle_species_list1.log",
3   "log_console": true,
4   "report_filename": "/volumes/output/wrangle_species_list1_rpt.json",
5   "in_species_list_filename": "/volumes/data/input/heuchera.txt",
6   "wrangler_configuration_file": "/volumes/data/wranglers/splist_wranglers_gbif.json",
7   "out_species_list_filename": "/volumes/output/heuchera_accepted1.txt"
8 }
```

## Run tool tutorials

The “run\_tutorial” script will run each tutorial with two arguments, the 1) command name and 2) parameters file. The parameters file will be a path on the local machine, in the tutorials/data/config directory. The script will translate that to the container path, and execute the command in the container with the container’s copy of the file. For example, the `wrangle_species_list` tutorial can be initiated with:

```
./run_tutorial.sh wrangle_species_list ./data/config/wrangle_species_list_gbif.json
```



**Command name**

## Run tool tutorials

The “run\_tutorial” script will run each tutorial with two arguments, the 1) command name and 2) parameters file. The parameters file will be a path on the local machine, in the tutorials/data/config directory. The script will translate that to the container path, and execute the command in the container with the container's copy of the file. For example, the `wrangle_species_list` tutorial can be initiated with:

```
./run_tutorial.sh wrangle_species_list ./data/config/wrangle_species_list_gbif.json
```



**Parameters file**

## Run tool tutorials

The “run\_tutorial” script will run each tutorial with two arguments, the 1) command name and 2) parameters file. The parameters file will be a path on the local machine, in the tutorials/data/config directory. The script will translate that to the container path, and execute the command in the container with the container's copy of the file. For example, the `wrangle_species_list` tutorial can be initiated with:

```
./run_tutorial.sh wrangle_species_list ./data/config/wrangle_species_list_gbif.json
```



**Command name**



**Parameters file**



## Run tool tutorials

The “run\_tutorial” script will run each tutorial with two arguments, the 1) command name and 2) parameters file. The parameters file will be a path on the local machine, in the tutorials/data/config directory. The script will translate that to the container path, and execute the command in the container with the container’s copy of the file. For example, the [wrangle\\_species\\_list](#) tutorial can be initiated with:

```
./run_tutorial.sh wrangle_species_list ./data/config/wrangle_species_list_gbif.json
```

## Outputs

All outputs are specified in the Tool Configuration File provided to the command, and will be copied to the data/outputs directory on completion.



## Behind the scenes

The “run\_tutorial” script will execute the following functions, unless their outputs have already been created:

1. Create several Docker Volumes to share data between the host and Docker container.
2. Build a [Docker image](#).
3. Start a Docker [Container](#) from the image, with volumes attached. A Container is similar to a fully functioning computer with data and applications.
4. Execute the specified command with the parameters in the specified configuration file. The process will execute using a parameter file and data in the *data* data volume and write outputs to the output volume, executing code in the [Docker container](#).
5. Copy the container `/volumes/output` directory back to the local data directory.
6. Stop and delete the container. All outputs in the docker volume are preserved and accessible the next time it is attached to a container.



- **Advances in biodiversity science, combined with emerging technologies and the ability to handle “big data” have greatly improved and expanded our capability to explore biodiversity in an unprecedented fashion. We can now link data from growing repositories (trees, occurrences, traits) and computational tools/approaches to integrate evolution and ecology at broad scale. This new synthesis is reshaping views of ecology/evolution with important conservation implications.**
- **Scaling biogeographic analyses from a small number of species or genera to explore patterns of evolution and diversity for thousands or tens of thousands of taxa on continental and global scales requires software tools that automate and parallelize computational tasks to make analyses practically feasible, efficient, with reasonable timelines.**

Allen et al.. 2019. Biodiversity synthesis across the green branches of the tree of life. *Nature Plants* 5:11-13.



**Any questions??**

**Please use the link to the Jamboard to write your question!**