



Molecular phylogenetics, phylogenomics, and phylogeography

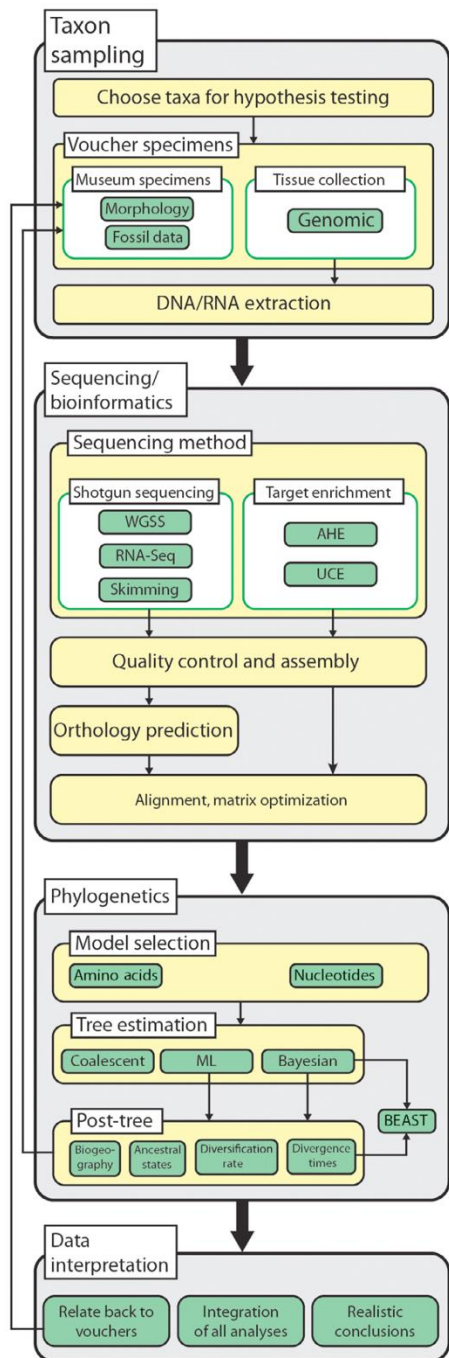
## Assessment of targeted enrichment locus capture across time and museums using odonate specimens

Aaron Goodman<sup>1,2,\*</sup>, Ethan Tolman<sup>1,2,t</sup>, Rhema Uche-Dike<sup>1,2,t</sup>, John Abbott<sup>3</sup>,  
Jesse W. Breinholt<sup>4,5</sup>, Seth Bybee<sup>6</sup>, Paul B. Frandsen<sup>7</sup>, J. Stephen Gosnell<sup>2,8</sup>,  
Rob Guralnick<sup>5</sup>, Vincent J. Kalkman<sup>9</sup>, Manpreet Kohli<sup>1,8</sup>, Judicael Fomekong Lontchi<sup>6</sup>,  
Pungki Lupiyaningdyah<sup>6,10</sup>, Lacie Newton<sup>1</sup>, Jessica L. Ware<sup>1</sup>



# Objectives for this study

- 1. Assess the effectiveness of targeted enrichment techniques in capturing genomic data from Odonata specimens preserved in museums
- 2. Evaluate the impact of specimen age and preservation methods on the success of locus capture and DNA quality
- 3. Determine whether high-quality genomic data can be generated from both recent and very old museum specimens, including those over 100 years old
- 4. Explore the potential for using historical museum specimens in phylogenetic and evolutionary studies, particularly for studying Odonata



## Museum specimen



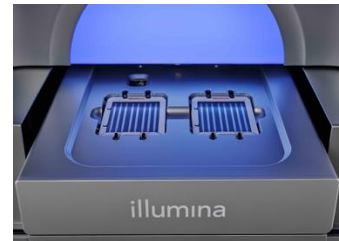
94 specimen  
64 Dragonflies  
30 Zygoptera

48 AMNH  
46 RMNH

## DNA extraction



DNA Quantification



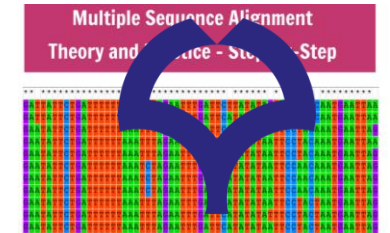
AHE illumina  
sequencing (Bybee et  
al.2021)



(Wingett & Andrew,  
2018)



Assembly with  
spades  
(Bankevitch,20  
12)



Multiple Sequence  
Alignment with  
MAFFT (Kato,2002)





B

### Hybrid-capture based Enrichment

Fragmented DNA

Adapter Ligation

Fragmented DNA Library

Biotin-labelled  
Probes

Probe Hybridization

Capture by  
Streptavidin-magnetic Beads  
Wash steps (remove non-specific  
background)

Target-Enriched DNA for  
NGS

# Targeted enrichment Overview

- **Why use AHE?** Sometimes, scientists don't need all the DNA from an organism—just specific regions that are informative for their study. TAE allows us to focus on capturing and sequencing these targeted regions
- **What is the “anchored” part?:** The method uses short pieces of DNA called "probes" designed to match and "anchor" to certain genome regions across many species. These probes "enrich" or increase the target regions' presence during sequencing, making them easier to study

B

## Hybrid-capture based Enrichment

Fragmented DNA

Adapter Ligation

Fragmented DNA Library

Biotin-labelled Probes

Probe Hybridization

Capture by  
Streptavidin-magnetic Beads  
Wash steps (remove non-specific  
background)

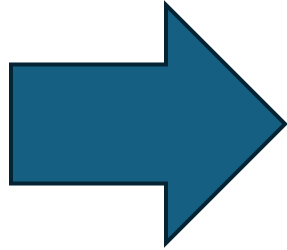
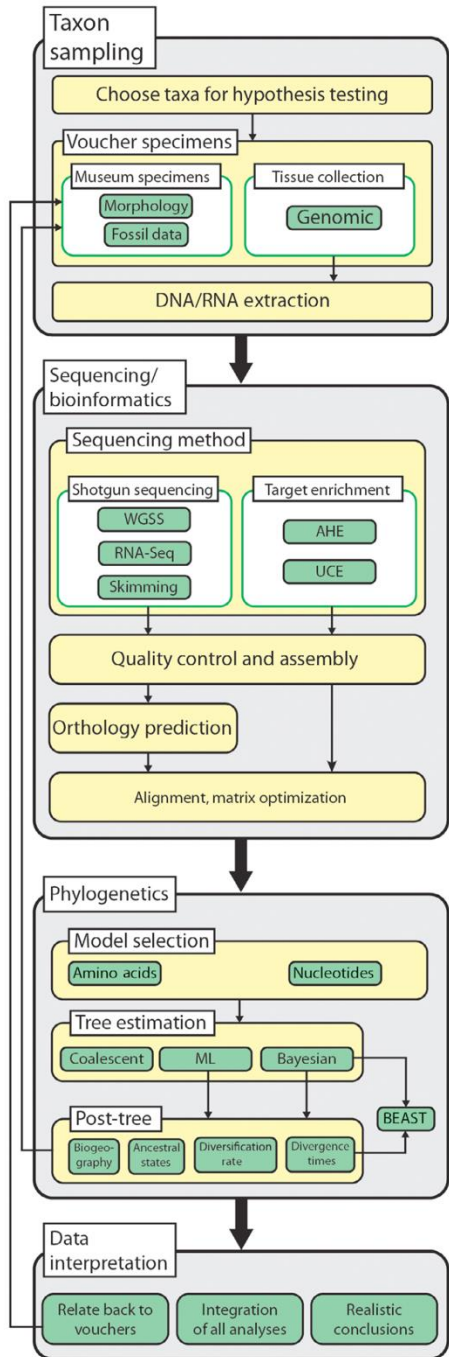
Target-Enriched DNA for  
NGS

# Targeted enrichment Overview

## How the Process Works

- **Designing Probes:** Scientists design DNA probes that will stick to the regions of the genome they are interested in (the target regions)
- **Capturing the Targets:** DNA is extracted from the samples (e.g., dragonflies or damselflies), and the probes are added. The probes latch onto the target regions in the DNA
- **Sequencing:** Once the targeted regions are captured, they can be sequenced. This means scientists can read the DNA letters (A, T, C, G) from these regions
- **Analyzing the Data:** The sequenced DNA is used to compare different species, build phylogenetic trees, and explore evolutionary relationships.

# Overview of the Bioinformatics Workflow for AHE



- 1.Data Preparation
- 2.Read Alignment
- 3.Assembly or Reference Mapping
- 4.Quality Control
- 5.Phylogenetic Analysis

#### Install **fastp** with bioconda:

**fastp** is a program that performs several pre-processing quality control steps for fastQ files (e.g., trim adapters, filter out bad reads, trim ends)

```
## create a conda environment (-n fastp)
## install fastp with bioconda (-c bioconda fastp)
conda create -n fastp -c bioconda fastp

## may have to reset shell settings before activating environment
source ~/.bashrc
```

#### Log onto the supercomputer:

```
## logging onto the Mendel AMNH cluster
ssh <username>@mendel.sdmz.amnh.org
```

#### Download and install Miniconda on the cluster:

```
## download the script to install miniconda
wget
https://repo.anaconda.com/miniconda/Miniconda3-latest-Linux-x86_64.sh

## install miniconda
bash Miniconda3-latest-Linux-x86_64.sh

## reset shell settings; this will allow us to install and activate
## conda environments without logging off and back on
source ~/.bashrc

## lists info about current conda install, check that install worked
conda info
```

# Step 1: Data Preparation

- **Input Data:** The sequencing process provides raw DNA sequences, often in a file format called FASTQ. These files contain millions of short DNA sequences (reads) generated from the specimens

- **FastQC:** A tool for checking the quality of the raw sequence data. It looks for errors or biases in the data.
- **Trimmomatic:** A tool used to trim low-quality sequences or adapters that may have been attached during sequencing.

## Step 2: Read Alignment (Mapping Reads to the Target Regions)

- **What's Happening?:** After cleaning the data, the next step is to align the reads to a reference genome or the regions that the AHE probes targeted.
- **Why?:** We need to see where each small DNA sequence (read) belongs in the genome to reconstruct the target regions.
- **Tools:**
  - **BWA (Burrows-Wheeler Aligner):** Aligns the reads to a reference genome or target regions.
  - **Bowtie2:** Another alignment tool, faster in some cases, for mapping reads to the reference.

To view multiple fastp results in one summary report, we will install and run [multiqc](#):

```
## deactivate fastp environment
conda deactivate

## create a new environment to install multiqc; installing multiple
## programs into one environment can result in conflicting
## dependencies, thus new environment prevents that
conda create -n multiqc -c conda-forge -c bioconda multiqc python=3

## make directory for multiqc in directory above FASTQ file directory
mkdir ../multiqc

## copy json files to multiqc directory
cp *.json ../multiqc

## change to multiqc directory
cd ../multiqc

## Make sure that all files end in "fastp.json"
rename .json _fastp.json *.json

## activate multiqc environment
conda activate multiqc
```



## Step 3: Assembly or Reference Mapping

- **De novo assembly:** If there is no reference genome, you can assemble the short reads into longer sequences (contigs)
- **Reference mapping:** If you have a reference genome, you can map your reads to it to reconstruct the target regions
- **Tools:**
  - **SPAdes:** A popular tool for assembling genomes from scratch (de novo assembly)
  - **SAMtools:** A tool for working with aligned reads, helping organize them into a usable format (BAM files)

### Assembly & Orthology Filter Steps:

In the assemblies folder: tanypteryx genome database files, unique assembly list (FASTQ names coupled with sample IDs, e.g. GEODE# or RMNH#), assembly slurm array job file, python script that does assembly steps: ASS\_F4\_fast.py, and the program [usearch](#).

\*\*\* genome and database files unique across datasets; *Tanypteryx* is a dragonfly genome \*\*\*

**NOTE TO SELF: include information about making a blast database for a genome file**

```
## create new environment and install assembly and alignment programs
conda create -n ahe_pipeline -c conda-forge -c bioconda python=3
mafft blast spades=3.15.4 biopython

## go to assemblies directory and list contents
cd ../assemblies
ls
```

# Step 4: Quality Control

---

- **What's Happening?:** After assembling or mapping, it's important to check if the process worked well and if the data is accurate.
- **Why?:** We want to ensure that the sequences we've reconstructed are reliable and represent the actual regions of DNA we targeted.
- **Tools:**
  - **QUAST:** A tool that checks how well the assembly or mapping performed.
  - **Coverage analysis:** Determines how well each region of the genome is represented by reads
  - Check Contamination

## Post-assembly steps to clean out contamination and choose the best sequence for each taxon:

In the after\_assembly folder: singleline.pl, split.py, SELFBLAST\_array.slurm, contamination\_filter.py, getlist.py, removelist.py, usearch program

```
## change into after assembly directory
cd ../after_assembly

## concatenate all orthologs from each sample into one fasta file
cat ../assemblies/*targetsFULL_ORTHO.fasta > ALL_FULL_ORTHO.fa
```

# Step 5: Phylogenetic Analysis

```
## submit alignment job
sbatch mafft_array.job

## NOTE TO SELF BEFORE ALIGNMENT
## write script to remove loci with less than 4 taxa or other cutoff
## to avoid including loci with no sequence data (will get alignment
## errors in array job and with FASconCAT)
## for loop to count number of taxa
## might have to type '>' in terminal; weird symbol issue?
for i in `cat locus_names.txt`; do grep '>' $i | wc -l >>
taxa_number.txt; done

##REMEMBER TO DELETE taxa_number.txt each time, the code does not
overwrite
```

## Concatenating our alignments

The next step in the process will be to concatenate our single locus files into a “concatenated supermatrix”.

```
## change directory
cd ../concatenated_loci_full

## copy alignments
```

- Once you have your sequences, the next step is to compare them across different species or samples to infer evolutionary relationships.

### Tools:

**MAFFT:** A tool for aligning multiple DNA sequences to see where they match or differ.

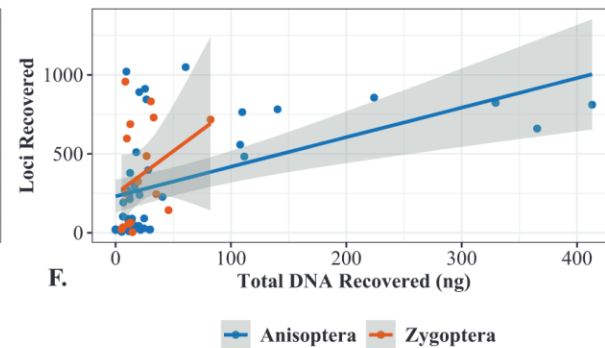
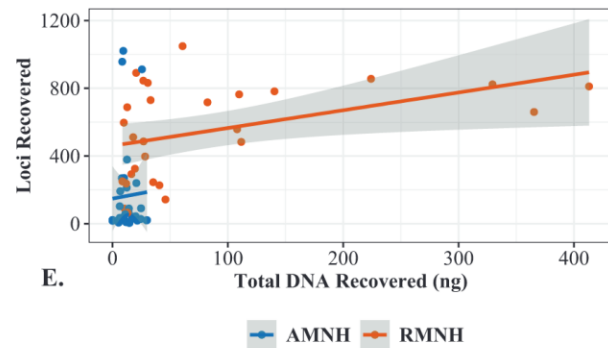
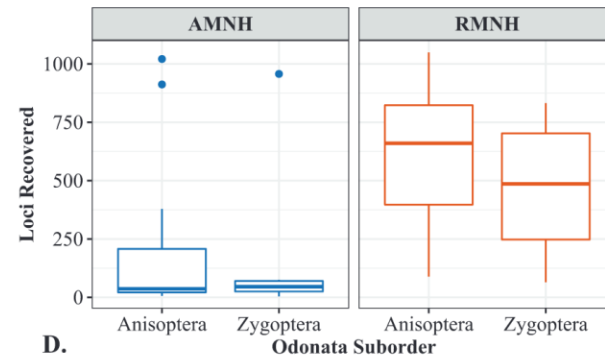
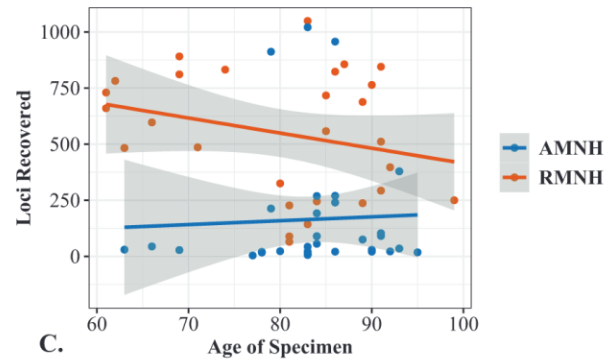
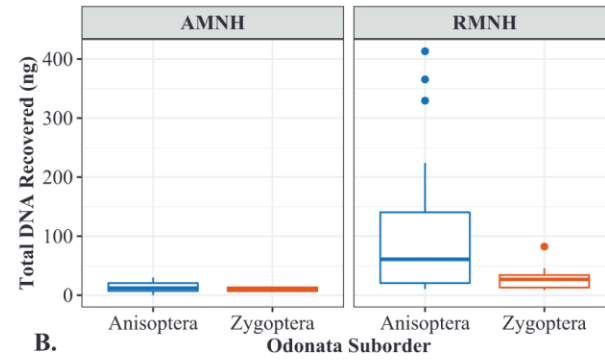
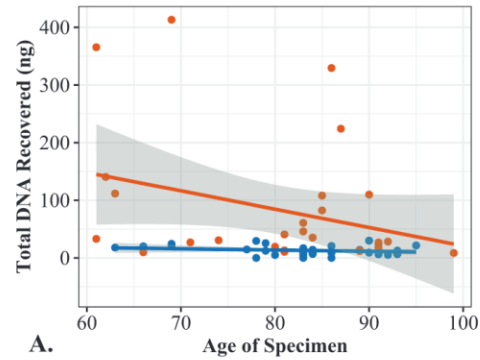
**IQTREE:** A tool for building phylogenetic trees, which show evolutionary relationships.

## Select models and tree search with IQtree:

```
## edit job file with email
nano iqtree_model_selection_and_tree_search.job

## run job
sbatch iqtree_model_selection_and_tree_search.job
```

# RESULTS





# CONCLUSION

---



Targeted enrichment successfully captures genomic data from recent and historical Odonata specimens, including those over 100 years old



Museum collections are valuable resources for phylogenetic and evolutionary studies, even with degraded DNA



The technique can be applied to other taxa, expanding research opportunities beyond Odonata



Museum specimens hold great potential for biodiversity, conservation, and evolutionary research