

---

# **phyluce documentation**

*Release 1.5.0*

**Brant C. Faircloth**

**Jul 12, 2017**



---

## Contents

---

<b>1</b>	<b>Contributions</b>	<b>3</b>
<b>2</b>	<b>Issues</b>	<b>5</b>
<b>3</b>	<b>Guide</b>	<b>7</b>
3.1	Purpose . . . . .	7
3.2	Installation . . . . .	9
3.3	Quality Control . . . . .	16
3.4	Assembly . . . . .	19
3.5	UCE Processing for Phylogenomics . . . . .	23
3.6	Tutorial I: UCE Phylogenomics . . . . .	43
3.7	Tutorial III: Harvesting UCE Loci From Genomes . . . . .	73
3.8	Tutorial IV: Identifying UCE Loci and Designing Baits To Target Them . . . . .	79
<b>4</b>	<b>Project info</b>	<b>117</b>
4.1	Citing . . . . .	117
4.2	License . . . . .	117
4.3	Changelog . . . . .	118
4.4	Attributions . . . . .	118
4.5	Funding . . . . .	120
4.6	Acknowledgements . . . . .	120
<b>5</b>	<b>Supporting documents</b>	<b>121</b>
5.1	List of Programs . . . . .	121
	<b>Bibliography</b>	<b>125</b>



Release v1.5.0. (*Changelog*)

**Author** Brant C. Faircloth

**Date** 12 July 2017 20:43 UTC (+0000)

**Copyright** This documentation is available under a Creative Commons (CC-BY) license.

phyluce (phy-**loo**-chee) is a software package that was initially developed for analyzing data collected from ultraconserved elements in organismal genomes (see *References* and <http://ultraconserved.org> for additional information).

The package now includes a number of tools spanning:

- the assembly of raw read data to contigs
- the separation of UCE loci from assembled contigs
- parallel alignment generation, alignment trimming, and alignment data summary methods in preparation for analysis
- alignment and SNP calling using UCE or other types of raw-read data.

As it stands, the phyluce package is useful for analyzing both data collected from UCE loci and also data collection from other types of loci for phylogenomic studies at the species, population, and individual levels.



# CHAPTER 1

---

## Contributions

---

`phyluce` is open-source (see [License](#)) and we welcome contributions from anyone who is interested. Please make a pull request on [github](#). The issue tracker for `phyluce` is also available on [github](#).





## CHAPTER 2

---

### Issues

---

If you have an issue, please ensure that you are experiencing this issue on a supported OS (see *Installation*) using the `conda` installation of `phyluce`. If possible, please submit a test case demonstrating the issue and indicate which platform, git checkout, and `phyluce` version you are using.



## Purpose

Phylogenomics offers the possibility of helping to resolve the *Tree of Life*. To do this, we often need either very cheap sources of organismal genome data or decent methods of subsetting larger genomes (e.g., vertebrates; 1 Gbp) such that we can collect data from across the genome in an efficient and economical fashion, find the regions of each genome that are shared among organisms, and attempt to infer the evolutionary history of the organisms in which we're interested using the data we collect.

Genome reduction techniques offer one way to collect these types of data from both small- and large-genome organisms. These “reduction” techniques include various flavors of *amplicon sequencing*, *RAD-seq* (**R**estriction site **A**ssociated **D**N**A** markers), *RNA-seq* (transcriptome sequencing), and sequence capture methods.

*phyluce* is a software package for working with data generated from sequence capture of UCE (**u**ltra-**c**onserved **e**lement) loci, as first published in [BCF2012]. Specifically, *phyluce* is a suite of programs to:

- assemble raw sequence reads from Illumina platforms into contigs
- determine which contigs represent UCE loci
- filter potentially paralagous UCE loci
- generate different sets of UCE loci across taxa of interest

Additionally, *phyluce* is capable of the following tasks, which are generally suited to any number of phylogenomic analyses:

- produce large-scale alignments of these loci in parallel
- manipulate alignment data prior to further analysis
- convert alignment data between formats
- compute statistics on alignments and other data

*phyluce* is written to process data/individuals/samples/species in parallel, where possible, to speed execution. Parallelism is achieved through the use of the *Python multiprocessing* module, and most computations are suited to workstation-class machines or servers (i.e., rather than clusters). Where cluster-based analyses are needed, *phyluce*

will produce the necessary outputs for input to the cluster/program that you are running so that you can setup the run according to your cluster design, job scheduling system, etc. Clusters are simply too heterogenous to do a good job at this part of the analytical workflow.

### Short-term goals (v1.4.x+)

We are currently working on a new release (this documentation) to:

- ease the burden of installing dependencies using `conda`
- standardize parameters input to various analyses
- improve logging of what is going on
- improve and standardize the documentation

### Longer-term goals (v2.0.0+ and beyond)

We are also working towards adding:

- simplify the `CLI` (command-line interface) of `phyluce`
- improve test coverage of the code by `unittests`
- SNP-calling pipelines (sensu *[BTS2013]*)
- sequence capture bait design
- identification of UCE loci

Much of this code is already written and in use by several of the *Contributed to the code*. As we test and improve these functions, we will add them to the code in the future.

### Who wrote this?

This documentation was written primarily by Brant Faircloth (<http://faircloth-lab.org>). Brant is also responsible for the development of most of the `phyluce` code. Bugs within the code are usually his.

You can find additional authors and contributors in the *Attributions* section.

### How do I report bugs?

To report a bug, please post an issue to <https://github.com/faircloth-lab/phyluce/issues>. Please also ensure that you are using one of the “supported” platforms:

- Apple OSX 10.9.x
- CentOS 6.x
- Ubuntu 14.04 LTS

and that you have installed `phyluce` and dependencies using `conda` as described in the *Installation* section.

## Installation

phyluce uses a number of tools that allow it to assemble data, search for UCE loci, align results reads, manipulate alignments, prepare alignments for analysis, etc. To accomplish these goals, phyluce uses wrappers around a number of programs that do each of these tasks (sometimes phyluce can use several different programs that accomplish the same task in different ways). As a result, the [dependency chain](#), or the programs that phyluce requires to run, is reasonably complex.

In previous versions (< 1.4.x), we required users to install a number of different software packages, get those into the user's *\$PATH*, cross their fingers, and hope that everything ran smoothly (it usually did not).

In the current versions (> 1.4.x), we have removed a number of dependencies, **and we very strongly suggest** that users install phyluce using either the [anaconda](#) or [miniconda](#) Python distributions.

**Attention:** We do not support installing phyluce through means other than the [conda](#) installer. This means that we do not test phyluce against any binaries, other than those we build and distribute through [conda](#). You will eventually be able to configure phyluce to use binaries of different provenance, although this will not be officially supported, other than providing a mechanism to do so.

---

**Note:** We build and test the binaries available through [conda](#) using 64-bit operating systems that include:

- Apple OSX 10.9.x
  - CentOS 6.x
  - Ubuntu 14.04 LTS
- 

## Why conda?

It may seem odd to impose a particular distribution on users, and we largely agree. However, [conda](#) makes it very easy for us to distribute both [Python](#) and non-Python packages (e.g. [velvet](#), [ABYSS](#), etc.), setup identical environments across very heterogeneous platforms (linux, osx), make sure all the *\$PATHs* are correct, and have things run largely as expected. Using [conda](#) has several other benefits, including environment separation similar to [virtualenv](#). In short, using [conda](#) gets us as close to a “one-click” install that we will probably ever get.

## Install Process

**Attention:** We do not support [phyluce](#) on Windows.

---

**Note:** We build and test the binaries available through [conda](#) using 64-bit operating systems that include:

- Apple OSX 10.9.x
  - CentOS 6.x
  - Ubuntu 14.04 LTS
- 

The installation process is now a 3-step process. You need to:

1. Install JAVA
2. Install `conda` (either `anaconda` or `miniconda`)
3. Install `phyluce`

Installing `phyluce` will install all of the required binaries, libraries, and `Python` dependencies.

### Install JAVA

Although we're using `conda`, you need to install a JAVA distribution for your platform. You should install **JAVA7**, as that will ensure a number of future tools in `phyluce` will work on your platform. Installing JAVA is a little tricky across different platforms (and somewhat beyond the scope of this document), but we describe, below, how we usually do it.

### Apple OS X

To install JAVA 1.7, download and install the JAVA 1.7 package from Oracle here: <http://www.java.com/en/download/manual.jsp>

### CentOS 6.5.x linux

You can install the JRE with the following `yum` command:

```
su -c "yum update"  
su -c "yum install java-1.7.0-openjdk"
```

### Ubuntu 14.04 linux

You can install the JRE with the following `apt-get` command:

```
sudo apt-get update  
sudo apt-get install openjdk-7-jre
```

### Install Anaconda or miniconda

After you installed `JAVA`, you need to install `anaconda` or `miniconda`. Which one you choose is up to you, your needs, how much disk space you have, and if you are on a fast/slow connection.

**Attention:** You can easily install `anaconda` or `miniconda` in your `$HOME`, although you should be aware that this setup can cause problems in some cluster-computing situations.

---

**Tip:** Do I want `anaconda` or `miniconda`?

The major difference between the two python distributions is that `anaconda` comes with many, many packages pre-installed, while `miniconda` comes with almost zero packages pre-installed. As such, the beginning `anaconda` distribution is roughly 200-500 MB in size while the beginning `miniconda` distribution is 15-30 MB in size.

---

## anaconda

Follow the instructions here for your platform: <http://docs.continuum.io/anaconda/install.html>

## miniconda

Find the correct *miniconda-x.x.x* file for your platform from <http://repo.continuum.io/miniconda/> and download that file. Be sure you **do not** get one of the packages that has a name starting with *miniconda3-*. When that has completed, run one of the following:

```
bash Miniconda-x.x.x-Linux-x86_64.sh [linux]
bash Miniconda-x.x.x-MacOSX-x86_64.sh [osx]
```

**Note:** Once you have installed Miniconda, we will refer to it as **anaconda** throughout the remainder of this documentation.

## Checking your \$PATH

Regardless of whether you install **anaconda** or **miniconda**, you need to check that you've installed the package correctly. To ensure that the correct location for **anaconda** or **miniconda** are added to your \$PATH (this occurs automatically on the \$BASH shell), run the following:

```
$ python -V
```

The output should look similar to (*x* will be replaced by a version):

```
Python 2.7.x :: Anaconda x.x.x (x86_64)
```

Notice that the output shows we're using the *Anaconda x.x.x* version of **Python**. If you do not see the expected output (or something similar), then you likely need to edit your \$PATH variable to add **anaconda** or **miniconda**.

The easiest way to edit your path, if needed is to open `~/ .bashrc` with a text editor (if you are using ZSH, this will be `~/ .zshrc`) and add, as the last line:

```
export PATH=$HOME/path/to/conda/bin:$PATH
```

where `$HOME/path/to/conda/bin` is the location of **anaconda/miniconda** on your system (usually `$HOME/anaconda/bin` or `$HOME/miniconda/bin`).

**Warning:** If you have previously set your \$PYTHONPATH elsewhere in your configuration, it may cause problems with your **anaconda** or **miniconda** installation of **phyluce**. The solution is to remove the offending library (-ies) from your \$PYTHONPATH.

## Add the faircloth-lab repository to conda

You need to add the location of the packages we need to your conda distributions. to do taht, you have to add the **faircloth-lab conda repository** to conda. You can do that with the following command, which automatically edits your `~/ .condarc` file):

```
conda config --add channels https://conda.anaconda.org/faircloth-lab
```

## Install phyluce

Now, you install [phyluce](#) and all of its dependencies by running:

```
conda install phyluce
```

This step will add the [phyluce](#) library to your \$PYTHONPATH and also add a number of scripts from the [phyluce](#) code to your *\$HOME/anaconda/bin* directory.

## What conda installs

When you install [phyluce](#), it specifies a number of dependencies that it needs to run. [conda](#) is great because it will pull specific **versions** of the required programs from the [faircloth-lab conda repository](#) and install those on your machine, setup the paths, etc.

Below is a list of what [phyluce](#) currently (v1.4.x) requires for installation.

### 3rd-party dependencies

- abyss 1.3.7 (max kmer=96; built with boost and google-sparsehash)
- bowtie 1.1.1
- bedtools 2.18.1
- bwa 0.7.7
- bx-python 0.7.1
- dendropy 3.12.0
- gatk-lite 2.3.0
- gblocks 0.91b
- illumiprocessor 2.0.7
- lastz 1.02.00
- mafft 7.130
- muscle 3.8.31
- picard 1.106
- pysam 0.7.7
- pyvcf 0.6.4
- raxml 8.0.19
- samtools 0.1.19
- trinity 2.0.6
- velvet 1.2.10 (max kmer=96)



## Python packages

- Python 2.7 (sets conda default Python to 2.7)
- numpy 1.7
- BioPython 1.63
- dendropy 3.12.0
- illumiprocessor 2.0.7
- phyluce 1.4.x

## Added benefits

An added benefit of using `conda` and installing packages in this way is that you can also run all of the 3rd-party binaries without worrying about setting the correct `$PATH`, etc.

For example, `phyluce` required `MUSCLE` for installation, and `MUSCLE` was installed by `conda` as a dependency of `phyluce`. Because `$HOME/anaconda/bin` (which we will now call `$CONDA`) is part of our path now, and because `phyluce` installed `MUSCLE`, we can also just run `MUSCLE` on the command-line, with:

```
$ muscle

MUSCLE v3.8.31 by Robert C. Edgar

http://www.drive5.com/muscle
This software is donated to the public domain.
Please cite: Edgar, R.C. Nucleic Acids Res 32(5), 1792-97.

Basic usage

    muscle -in <inputfile> -out <outputfile>

Common options (for a complete list please see the User Guide):

    -in <inputfile>      Input file in FASTA format (default stdin)
    -out <outputfile>   Output alignment in FASTA format (default stdout)
    -diags              Find diagonals (faster for similar sequences)
    -maxiters <n>      Maximum number of iterations (integer, default 16)
    -maxhours <h>     Maximum time to iterate in hours (default no limit)
    -html              Write output in HTML format (default FASTA)
    -msf               Write output in GCG MSF format (default FASTA)
    -clw               Write output in CLUSTALW format (default FASTA)
    -clwstrict         As -clw, with 'CLUSTAL W (1.81)' header
    -log[a] <logfile> Log to file (append if -loga, overwrite if -log)
    -quiet             Do not write progress messages to stderr
    -version           Display version information and exit

Without refinement (very fast, avg accuracy similar to T-Coffee): -maxiters 2
Fastest possible (amino acids): -maxiters 1 -diags -sv -distance1 kbit20_3
Fastest possible (nucleotides): -maxiters 1 -diags
```

This is true for other binaries you install from our repository (e.g. `velveth`, `velvetg`, `abyss-pe`, `mafft`) or any other `conda` repository - those binaries are all stored in `$CONDA/bin`.

We have setup conda to install other files in a standard location as well. So JAR files are stored in `$CONDA/jar`; libraries that you install from our repo are stored in `$CONDA/lib`, etc. The locations and versions are standardized within our `conda` distribution so that we always know where things are installed, hopefully avoiding lots of the problems with `dependency hell` and making our lives easier.

The structure of the conda repository that we use looks like the following:

```
$CONDA/ (<= can be main CONDA or CONDA env)
  bin/
    Binaries we build
  config/
    Holds a default configuration file for phyluce
  docs/
    Pertinent documentation, if any.
  include/
    Header files (Qt, Boost)
  jar/
    Java Archive (JAR) Files
  lib/
    Libraries
  libexec/
    Helpers called by other programs (e.g. mafft)
  share/
    Additional support files
  test/
    Test data required by some programs.
```

## \$PATH configuration

As of v1.5, `phyluce` uses a configuration file to keep track of paths to relevant binaries, as well as some configuration information. This file is located at `$CONDA/config/phyluce.conf`. Although you can edit this file directly, you can also create a user-specific configuration file at `~/.phyluce.conf` (**note the preceding dot**), which will override the default values for different paths. So, if you need to use a slightly different binary or you want to experiment with new binaries (e.g. for assembly), then you can change the paths in this file rather than deal with hard-coded paths.

**Attention:** You do NOT need to do anything with this file - \$PATHs should automatically resolve.

**Warning:** Changing the \$PATHs in the config file can break things pretty substantially, so please use with caution (and edit the copy at `~/.phyluce.conf`) rather than the default copy.

The format of the config file as of v1.5 looks like the following:

```
[abyss]
abyss:$CONDA/bin/ABYSS
abyss-pe:$CONDA/bin/abyss-pe

[bowtie]
bowtie:$CONDA/bin/bowtie

[bwa]
bwa:$CONDA/bin/bwa
```

```
[gblocks]
gblocks:$CONDA/bin/gblocks

[java]
executable:java
mem:-Xmx8g
jar:$CONDA/jar
gatk:GenomeAnalysisTKLite.jar

[lastz]
lastz:$CONDA/bin/lastz

[mafft]
mafft:$CONDA/bin/mafft

[muscle]
muscle:$CONDA/bin/muscle

[raxml]
raxmlHPC-SSE3:$CONDA/bin/raxmlHPC-SSE3
raxmlHPC-PTHREADS-SSE3:$CONDA/bin/raxmlHPC-PTHREADS-SSE3

[samtools]
samtools:$CONDA/bin/samtools

[trinity]
trinity:$CONDA/bin/Trinity
max_memory:8G
kmer_coverage:2

[velvet]
velvetg:$CONDA/bin/velvetg
velveth:$CONDA/bin/velveth

#-----
#   Advanced
#-----

[headers]
trinity:comp\d+_c\d+_seq\d+|c\d+_g\d+_i\d+|TR\d+\\|c\d+_g\d+_i\d+
velvet:node_\d+
```

## Other useful tools

You will need to be familiar with the command-line/terminal, and it helps to have a decent text editor for your platform:

- gedit [linux]
- Sublime Text [linux, osx]
- TextMate [osx]

## Quality Control

When you receive your data from the sequencer, typically they are already demultiplexed (split by sequence tags) for you. If your data are from the [MiSeq](#), they may also have been trimmed of adapter contamination.

---

**Note:** If your data are not demultiplexed, they can come in a vast array of different types, although one common type are so-called BCL-files. Regardless, if your data are not demultiplexed fastq files, you will need to talk to your sequencing provider about how to accomplish demultiplexing. I provide an **unsupported** guide to demultiplexing *BCL* files using [Casava](#) or [bcl2fastq](#) here: <https://gist.github.com/brantfaircloth/3125885>.

---

Regardless, you need to do a fair bit of quality control on your read data. **At a minimum**, this includes:

- getting some idea of how much data you have
- trimming adapter contamination of reads
- trimming low quality bases from reads

Although the MiSeq may trim some adapter contamination, running your reads through an additional round of trimming won't hurt. There is also [lots of evidence](#) showing that quality control of your read data has a **large** effect on your overall success. Bottom line is: *garbage in, garbage out*.

## Read Counts

The first thing to do once you have your read data in hand is to get an idea of the split of reads among your samples. This does two things:

1. Allows you to determine how well you split the run among your sequence tags
2. Shows you which samples may be suboptimal (have very few reads)

Really unequal read counts mean that you may want to switch up your library quantification process (or your pooling steps, prior to enrichment). Suboptimal read counts for particular libraries may or may not mean that the enrichments of those samples “failed”, but it's generally a reasonable indication.

You can get read counts in one of two way - the first is very simple, the second uses [phyluce](#).

### Count reads using shell tools

You can get a quick and dirty idea of the number of reads you have for each sample using simple shell or “terminal” tools - counting lines in lots of files is a task that's really suited to UNIX-like operating systems.

Because FASTQ files contain 4 lines for each read, we can count all the lines in every file and divide each by 4 to get the count of reads in a given file.

Assuming your reads are in a directory structure like:

```
some-directory-name/  
  sample1_R1.fastq.gz  
  sample1_R2.fastq.gz  
  sample2_R1.fastq.gz  
  sample2_R2.fastq.gz
```

you can count lines across all files in all directories, using:

```
for i in some-directory-name/*_R1.fastq.gz; do echo $i; gunzip -c $i | wc -l; done
```

This will output a list of results like:

```
bfidt-000_R1.fastq.gz
10000000
bfidt-000_R2.fastq.gz
16000000
```

If you **divide these numbers by 4**, then that is the count of R1 reads. The number of reads in the R2 files, if you have paired-end data, should **always** be equal.

### Get read counts using phyluce

There is a bit of code in phyluce that lets you count reads. To use it, you pass the code the directory containing reads you want to summarize and run:

```
phyluce_assembly_get_fastq_lengths --input /directory/containing/reads/ --csv; done
```

You can run this across many directories of reads as described in *Tutorial I: UCE Phylogenomics*.

### Adapter- and quality-trimming

Generally speaking, *Casava* and *bcl2fastq*, the two programs used to demultiplex sequence data from *Illumina* platforms, output fastq files in a format similar to the following:

```
Project_name/
  Sample_BFIDT-000
    BFIDT-000_AGTATAGCGC_L001_R1_001.fastq.gz
    BFIDT-000_AGTATAGCGC_L001_R2_001.fastq.gz
  Sample_BFIDT-001
    BFIDT-001_TTGTTGGCGG_L001_R1_001.fastq.gz
    BFIDT-001_TTGTTGGCGG_L001_R2_001.fastq.gz
```

What you want to do is to clean these reads of adapter contamination and trim low-quality bases from all reads (and probably also drop reads containing “N” (ambiguous) bases. Then you want to interleave the resulting data, where read pairs are maintained, and also have an output file of singleton data, where read pairs are not.

You can do this however you like. *phyluce* assumes that the results structure of your data after trimming will look like the following (replace *genus\_species* with your taxon names):

```
genus_species1/
  split-adapter-quality-trimmed/
    genus_species1-READ1.fastq.gz
    genus_species1-READ2.fastq.gz
    genus_species1-READ-singleton.fastq.gz
genus_species2/
  split-adapter-quality-trimmed/
    genus_species2-READ1.fastq.gz
    genus_species2-READ2.fastq.gz
    genus_species2-READ-singleton.fastq.gz
```

**Note:** We no longer create *interleaved* fastq files because most programs assume that PE (paired-end) reads are maintained as read pairs. Rather than create an interleaved set of data and a non-interleaved set of data, we have decided to keep files for each of R1, R2, and singleton reads (those that have lost their mates) and never produce an interleaved file.

## Trimming with illumiprocessor

You can run your adapter and quality trimming and output the files in the correct format using a program I wrote called `illumiprocessor`. It automates these processes over hundred of files and produces output in the format we want downstream.

**Attention:** Prior to running `illumiprocessor`, if you used Casava, it will be easiest to place all of the demultiplexed reads into a single directory.

You need to generate a configuration file `your-illumiprocessor.conf`, that gives details of your reads, how you want them processed, and what renaming options to use. There are **several** variations in formatting required depending on the library preparation method that you used.

**Attention:** See the documentation for `illumiprocessor` for configuration information.

You can run `illumiprocessor` against your data (in *demultiplexed*) with the following. If you do not have a multicore machine, you may wish to run with `--cores=1`. Additionally, multicore operations require a fair amount of RAM, so if you're low on RAM, run with fewer cores:

```
illumiprocessor \  
  --input demultiplexed \  
  --output uce-clean \  
  --config your-illumiprocessor.conf \  
  --cores 12
```

The clean data will appear in `uce-clean` with the following structure:

```
uce-clean/  
  genus_species1/  
    adapters.fasta  
    raw-reads/  
      genus_species1-READ1.fastq.gz (symlink)  
      genus_species1-READ2.fastq.gz (symlink)  
    split-adapter-quality-trimmed/  
      genus_species1-READ1.fastq.gz  
      genus_species1-READ2.fastq.gz  
      genus_species1-READ-singleton.fastq.gz  
    stats/  
      genus_species1-adapter-contam.txt  
  
  genus_species2/  
    adapters.fasta  
    raw-reads/  
      genus_species2-READ1.fastq.gz (symlink)  
      genus_species2-READ2.fastq.gz (symlink)  
    split-adapter-quality-trimmed/  
      genus_species2-READ1.fastq.gz  
      genus_species2-READ2.fastq.gz  
      genus_species2-READ-singleton.fastq.gz  
    stats/  
      genus_species2-adapter-contam.txt
```

You are now ready to move onto processing the cleaned read data.

## Assembly

### Setup

Once your reads are clean, you're ready to assemble. At the moment, you can use [velvet](#), [ABySS](#), and [Trinity](#) for assembly. However, be aware that there is not a [conda](#) install package for [Trinity](#) due to some difficulties in how that package is structured.

Most of the assembly process is automated using code within [phyluce](#), specifically the following 3 scripts:

- `phyluce_assembly_assemblo_abyss`
- `phyluce_assembly_assemblo_trinity`
- `phyluce_assembly_assemblo_velvet`

The code of each of the above programs **always** expects your input directories to have the following structure (from the [Quality Control](#) section):

```
uce-clean/
  genus_species1/
    adapters.fasta
    raw-reads/
      genus_species1-READ1.fastq.gz (symlink)
      genus_species1-READ2.fastq.gz (symlink)
    split-adapter-quality-trimmed/
      genus_species1-READ1.fastq.gz
      genus_species1-READ2.fastq.gz
      genus_species1-READ-singleton.fastq.gz
    stats/
      genus_species1-adapter-contam.txt
```

And, each of these assembly helper programs take the same configuration files as input. You should format the configuration file for input according to the following scheme:

```
[samples]
name_you_want_assembly_to_have:/path/to/uce-clean/genus_species1
```

In practice, this means you need to create a configuration file that looks like:

```
[samples]
anas_platyrrhynchos1:/path/to/uce-clean/anas_platyrrhynchos1
anas_carolinensis1:/path/to/uce-clean/anas_carolinensis1
dendrocygna_bicolor1:/path/to/uce-clean/dendrocygna_bicolor1
```

The *assembly name* on the left side of the colon can be whatever you want. The *path name* on the right hand side of the colon must be a valid path to a directory containing read data in a format similar to that described above.

**Attention:** Assembly names **MUST** be unique.

### Question: How do I name my samples/assemblies?

Naming samples is a contentious issue and is also a hard thing to deal with using computer code. You should **never** have a problem if you name your samples as follows, where the genus and specific epithet are separated by an underscore, and multiple individuals of a given species are indicated using a trailing integer value:

```
anas_platyrhynchos1
anas_carolinensis1
dendrocygna_bicolor1
```

You should also not have problems if you use a naming scheme that suffixes the species binomial(s) with an accession number that is **simply** formatted (e.g. no slashes, dashes, etc.):

```
anas_platyrhynchos_KGH2267
anas_carolinensis_KGH2269
dendrocygna_bicolor_DWF4597
```

## Running the assembly

Once your configuration file is created (best to use a decent text editor) that will not cause you grief, you are ready to start assembling your read data into contigs that we will search for UCES. The code to do this for the three helper scripts is below (remember, we are using `$HOME/anaconda/bin` generically to refer to your `anaconda` or `miniconda` install).

### General process

The general process that the helper scripts use is:

1. Create the output directory (AKA `$ASSEMBLY`, below)
2. Create a `contigs` folder within the output directory
3. For each taxon create `$ASSEMBLY/genus-species` directory, based on config file entries
4. Find the correct fastq files for a given sample
5. Input those fastq files to whichever assembly program
6. Assemble reads using user-defined/static `--kmer` value
7. Strip contigs of potentially problematic bases (ABYSS-only)
8. Normalize contig names
9. Link assembly file with normalized names in `$ASSEMBLY/genus-species/` into `$ASSEMBLY/contigs/genus-species`

### velvet

```
# make a directory for log files
mkdir log
# run the assembly
phyluce_assembly_assemble_velvet \
  --config config_file_you_created.conf \
  --output /path/where/you/want/assemblies \
  --kmer 35 \
  --subfolder split-adapter-quality-trimmed \
  --cores 12 \
  --clean \
  --log-path log
```



## Results

The directory structure created for `velvet`-based assemblies looks like:

```
path-to-output-directory/
  contigs/
    genus-species1 -> ../genus-species1/out_k31/contigs.fa
  genus-species1/
    contigs.fasta -> out_k31/contigs.fa
    out_k31
    velvetg-k31.err.log
    velvetg-k31.out.log
    velveth-k31.err.log
    velveth-k31.out.log
```

## ABySS

```
# make a directory for log files
mkdir log
# run the assembly
phyluce_assembly_assemblo_abys \
  --config config_file_you_created.conf \
  --output /path/where/you/want/assemblies \
  --kmer 35 \
  --subfolder split-adapter-quality-trimmed \
  --cores 12 \
  --clean \
  --log-path log
```

**Attention:** Following assembly, `phyluce_assembly_assemblo_abys` modifies the assemblies by replacing degenerate base codes with standard nucleotide encodings. We do this because `lastz`, which we use to match contigs to targeted UCE loci, is not compatible with degenerate IUPAC codes.

The `phyluce_assembly_assemblo_abys` code makes these substitutions for every site having a degenerate code by selecting the appropriate nucleotide encoding randomly. The code also renames the ABySS assemblies using the `velvet` naming convention. The modified contigs are then symlinked into `$ASSEMBLY/contigs`. Unmodified contigs are available in `$ASSEMBLY/genus-species/out_k*-contigs.fa`

## Results

The directory structure created for `ABySS`-based assemblies looks like:

```
path-to-output-directory/
  contigs/
    genus-species1 -> ../genus-species1/out_k31-contigs-velvet.fa
  genus-species1/
    abyss-k31.err.log
    contigs.fasta -> out_k31-contigs-velvet.fa
    out_k31-contigs.fa
    out_k31-scaffolds.fa
    out_k31-unitigs.fa
    abyss-k31.out.log
    coverage.hist
```

```
out_k31-contigs-velvet.fa
out_k31-stats
```

## Trinity

```
# make a directory for log files
mkdir log
# run the assembly
phyluce_assembly_assemblo_trinity \
  --config config_file_you_created.conf \
  --output /path/where/you/want/assemblies \
  --subfolder split-adapter-quality-trimmed \
  --clean \
  --cores 12 \
  --log-path log
```

---

### Question: Why do I not use the `-kmer` flag with `phyluce_assembly_assemblo_trinity`?

Trinity assembles using a “static” or “program-defined” (i.e., not user- defined) kmer value of 25.

---

---

### Question: What is the `-clean` option?

The `-clean` option removes extraneous temporary files that Trinity creates during the assembly process. This results in a vastly smaller *assemblies* directory.

---

## Results

The directory structure created for Trinity-based assemblies looks like:

```
path-to-output-directory/
  contigs/
    genus-species1 -> ../genus-species1/Trinity.fasta
  genus-species1/
    contigs.fasta -> Trinity.fasta
    Trinity.fasta
    trinity.log
```

## Common questions

---

### Question: Which assembly program do I pick?

Generally, I would suggest that you use Trinity. In my hands, it produces reasonable contig assemblies that are longer than the assemblies built by either `velvet` or `ABYSS`. There are some caveats, however. If you want the **most accurate** assemblies possible, then it may be best to use `ABYSS`. This is because `ABYSS` runs read-based error correction prior to assembly which results in more accurate contigs.

---

---

### Question: For `ABYSS` and `velvet`, what `-kmer` value do I use?

Also a hard question. Part of the reason that it is hard is due to the fact that we are trying to assemble data of heterogenous read depth (i.e., our reads are spread across (mostly) UCE loci, but the depth of coverage of each locus is variable due to capture efficiency). Longer kmer values can give you longer (but fewer) contigs, while shorter kmer values produce fewer, more abundant contigs. In most cases, your assemblies will be decent with a kmer value around 55-65.

## UCE Processing for Phylogenomics

The workflow described below is meant for users who are analyzing UCE in phylogenetic contexts - meaning that you are interested in addressing questions at or deeper than the species-level.

### Identifying UCE loci

Once we have assembled our fastq data (see *Assembly*), we need to process those contigs to (a) determine which represent enrichend UCE loci and (b) remove any potential paralogs from the data set. Before we can do that, we need to do a little preparatory work by downloading a FASTA file representing the probe set that we used.

#### Get the probe set FASTA

To identify which of the contigs we've assembled are UCE loci (and which UCE loci they might be), we are going to match our assembled contigs to the probes we used to enrich UCE loci. Before we do that, however, we need to download a copy of probe set we used for matching purposes.

**Attention:** We archive official probe sets at <https://github.com/faircloth-lab/uce-probe-sets/>, but you need to be careful about which one you grab - probe sets can be of different sizes (e.g. 2,500 or 5,500 loci) and for different groups of taxa (e.g., amniotes, fish)

#### Download the probe set

To download a given probe set for *phyluce*, you need to figure out which probe set you need. Then, you can use a command like `wget` on the command-line (or navigate with your browser to the URL and save the file):

```
# to get the 2.5k, amniote probe set
wget https://raw.githubusercontent.com/faircloth-lab/uce-probe-sets/master/uce-2.5k-
↳probe-set/uce-2.5k-probes.fasta

# to get the 5k, amniote probe set
wget https://raw.githubusercontent.com/faircloth-lab/uce-probe-sets/master/uce-5k-
↳probe-set/uce-5k-probes.fasta
```

#### Match contigs to probes

Once we've downloaded the probe set we used to enrich UCE loci, we need to find which of our assembled contigs are the UCE loci that we enriched. During this process, the code will also remove any contigs that appear to be duplicates as a result of assembly/other problems **or** a biological event(s).

The way that this process works is that **phyluce** aligns (using **lastz**) the contigs you assembled to the probes you input on a taxon-by-taxon (or otu-by-otu) basis. Then, the code parses the alignment file to determine which contigs matched which probes, whether any probes from a single locus matched multiple contigs or whether a single contig matched probes designed from multiple UCE loci. Either of these latter two events suggests that the locus in question is problematic.

---

**Hint:    ADVANCED:** The default regular expression assumes probes in your file are named according to `uce-NNN_pN`, where `uce-` is just a text string, `NNN` is an integer value denoting each unique locus, `_p` is a text string denoting a “probe” targeting locus `NNN`, and the trailing `N` is an integer value denoting each unique probe targeting the same locus.

If you are using a custom probe file, then you will either need to ensure that your naming scheme conforms to this approach **OR** you will need to input a different regular expression to convert the probe names to locus names using the `--regex` flag.

---

To identify which of your assembled contigs are UCE contigs, run:

```
# make a directory for log files
mkdir log
# match contigs to probes
phyluce_assembly_match_contigs_to_probes \
  --contigs /path/to/assembly/contigs/ \
  --probes uce-5k-probes.fasta \
  --output /path/to/uce/output \
  --log-path log
```

When you run this code, you should see output similar to:

```
2014-04-24 14:38:15,979 - match_contigs_to_probes - INFO - ===== Starting_
↳match_contigs_to_probes =====
2014-04-24 14:38:15,979 - match_contigs_to_probes - INFO - Version: git 7aec8f1
2014-04-24 14:38:15,979 - match_contigs_to_probes - INFO - Argument --contigs: /path/
↳to/assembly/contigs/
2014-04-24 14:38:15,980 - match_contigs_to_probes - INFO - Argument --keep_
↳duplicates: None
2014-04-24 14:38:15,980 - match_contigs_to_probes - INFO - Argument --log_path: None
2014-04-24 14:38:15,980 - match_contigs_to_probes - INFO - Argument --min_coverage: 80
2014-04-24 14:38:15,980 - match_contigs_to_probes - INFO - Argument --min_identity: 80
2014-04-24 14:38:15,980 - match_contigs_to_probes - INFO - Argument --output: /path/
↳to/uce/output
2014-04-24 14:38:15,980 - match_contigs_to_probes - INFO - Argument --probes: uce-5k-
↳probes.fasta
2014-04-24 14:38:15,981 - match_contigs_to_probes - INFO - Argument --regex: ^(uce-
↳\d+)(?:_p\d+.* )
2014-04-24 14:38:15,981 - match_contigs_to_probes - INFO - Argument --verbosity: INFO
2014-04-24 14:38:16,138 - match_contigs_to_probes - INFO - Checking probe/bait_
↳sequences for duplicates
2014-04-24 14:38:19,022 - match_contigs_to_probes - INFO - Creating the UCE-match_
↳database
2014-04-24 14:38:19,134 - match_contigs_to_probes - INFO - Processing contig data
2014-04-24 14:38:19,134 - match_contigs_to_probes - INFO - -----
↳-----
2014-04-24 14:38:25,713 - match_contigs_to_probes - INFO - genus_species1: 1031 (70.14
↳%) uniques of 1470 contigs, 0 dupe probe matches, 48 UCE probes matching multiple_
↳contigs, 117 contigs matching multiple UCE probes
2014-04-24 14:38:32,846 - match_contigs_to_probes - INFO - genus_species2: 420 (68.52
↳%) uniques of 613 contigs, 0 dupe probe matches, 30 UCE probes matching multiple_
↳contigs, 19 contigs matching multiple UCE probes
```

```

2014-04-24 14:38:39,184 - match_contigs_to_probes - INFO - genus_species3: 1071 (63.15
↳%) uniques of 1696 contigs, 0 dupe probe matches, 69 UCE probes matching multiple_
↳contigs, 101 contigs matching multiple UCE probes
2014-04-24 14:49:59,654 - match_contigs_to_probes - INFO - -----
↳-----
2014-04-24 14:49:59,654 - match_contigs_to_probes - INFO - The LASTZ alignments are_
↳in /path/to/uce/output/
2014-04-24 14:49:59,654 - match_contigs_to_probes - INFO - The UCE match database is_
↳in /path/to/uce/output/probe.matches.sqlite
2014-04-24 14:49:59,655 - match_contigs_to_probes - INFO - ===== Completed_
↳match_contigs_to_probes =====

```

**Note:** The \*.log files for each operation are always printed to the screen AND also written out to the \$CWD (current working directory). You can keep these files more orderly by specifying a \$LOG on the command line using the --log-path option.

## Results

The resulting files will be in the:

```
/path/to/output
```

directory. If you look in this directory, you'll see that it contains species- specific *lastz\_* files as well as an *sqlite* database:

```

$ ls /path/to/output
genus_species1.contigs.lastz
genus_species2.contigs.lastz
genus_species3.contigs.lastz
probe.matches.sqlite

```

The \*.lastz files within the /path/to/output directory are basically for reference and individual review (they are text files that you can open using a text editor to view). The really important data from the *lastz* files are summarized in the:

```
probe.matches.sqlite
```

database file. It's probably a good idea to have some knowledge of how this database is structured, since it's basically what makes the next few steps work. So, let's go over the structure and contents of this database.

### The probe.matches.sqlite database

probe.matches.sqlite is a [relational database](#) that summarizes all **valid** matches of contigs to UCE loci across the set of taxa that you fed it. The database is created by and for a program named *sqlite*, which is a very handy, portable SQL database. For more info on SQL and SQLITE, see this [sqlite-tutorial](#). I'll briefly cover the database contents and use below.

First, take a look at the contents of the database by running:

```
sqlite3 probe.matches.sqlite
```

You'll now see something like:

```
SQLite version 3.7.3
Enter ".help" for instructions
Enter SQL statements terminated with a ";"
sqlite>
```

It's often easier to change some defaults for better viewing, so at the prompt, paste in the following:

```
sqlite> .mode columns
sqlite> .headers on
sqlite> .nullvalue .
```

---

**Tip:** For more info on `sqlite` "dot" commands, you can type `.help`.

---

Now that that's done, let's see which tables the database contains by running the `.tables` command:

```
sqlite> .tables
match_map matches
```

This tells us there's two tables in the database, named `match_map` and `matches`.

### The matches table

Let's take a look at the contents of the `matches` table. Once you've started the `sqlite` interface, run:

```
sqlite> SELECT * FROM matches LIMIT 10;
```

This query select all rows (`SELECT *`) from the `matches` table (`FROM matches`) and limits the number of returned rows to 10 (`LIMIT 10`). This will output data that look something like:

uce	genus_species1	genus_species2	genus_species3
uce-500	1	.	.
uce-501	1	.	.
uce-502	1	.	.
uce-503	1	1	1
uce-504	1	.	.
uce-505	1	.	.
uce-506	.	.	.
uce-507	1	.	.
uce-508	1	1	.
uce-509	1	1	1

Basically, what this indicates is that you enriched 9 of 10 targeted UCE loci from `genus_species1`, 3 of 10 UCE loci in the list from `genus_species2`, and 2 of 10 UCE loci from `genus_species3`. The locus name is given in the `uce` column. Remember that we've limited the results to 10 rows for the sake of making the results easy to view.

If we wanted to see only those loci that enriched in all species, we could run:

```
sqlite> SELECT * FROM matches WHERE genus_species1 = 1
...> AND genus_species2 = 1 AND genus_species3 = 1;
```

Assuming we only had those 10 UCE loci listed above in the database, if we ran this query, we would see something like:

uce	genus_species1	genus_species2	genus_species3
uce-503	1	1	1
uce-509	1	1	1

Basically, the `matches` table and this query are what we run to generate **complete** (only loci enriched in all taxa) and **incomplete** (all loci enriched from all taxa) matrices very easily and quickly (see [Creating a data matrix configuration file](#)).

### The `match_map` table

The `match_map` table shows us which species-specific, contigs match which UCE loci. Because each assembly program assigns an arbitrary designator to each assembled contig, we need to map these arbitrary designators (which also differ for each taxon/OTU) to the UCE locus to which it corresponds. Because assembled contigs are also not in any particular orientation relative to each other across taxa/OTUs (i.e., they may be 5' - 3' or 3' - 5'), the database also records the orientation of all contigs relative to orientation of each probe in the probes file.

Let's take a quick look at the `match_map` table:

```
SELECT * FROM match_map LIMIT 10;
```

This query is similar to the one that we ran against `matches` and returns the first 10 rows of the `match_map` table:

uce	genus_species1	genus_species2	genus_species3
uce-500	node_233 (+)	.	.
uce-501	node_830 (+)	.	.
uce-502	node_144 (-)	.	.
uce-503	node_1676 (+)	node_243 (+)	node_322 (+)
uce-504	node_83 (+)	.	.
uce-505	node_1165 (-)	.	.
uce-506	.	.	.
uce-507	node_967 (+)	.	.
uce-508	node_671 (+)	node_211 (-)	.
uce-509	node_544 (-)	node_297 (+)	node_37 (+)

As stated above, these results show which assembled contigs “hit” particular UCE loci. So, if we were to open the `$ASSEMBLY/contigs/genus_species1.contigs.fasta` symlink the contig named `node_1676` corresponds to UCE locus `uce-503`. Because contigs are named arbitrarily during assembly, this same UCE locus is also found in `genus_species2`, but it is named `node-243`.

Each entry in the rows also provides the orientation for particular contigs (-) or (+). This orientation is relative to the orientation of the UCE probes/locus in the source genome (e.g., chicken for tetrapod probes).

We use this table to generate a FASTA file of UCE loci for alignment (see `:ref :fasta-file`), after we've identified the loci we want in a particular data set (see [Creating a data matrix configuration file](#)). The code for this step also uses the associated orientation data to ensure that all the sequence data have the same orientation prior to alignment (some aligners will force alignment of all reads using the given orientation rather than also trying the reverse complement and picking the better alignment of the two).

Now that we know the taxa for which we've enriched UCE loci and which contigs we've assembled match which UCE loci, we're ready to generate some data matrices.

The data matrix generation process consists of two distinct parts:

1. Getting locus counts and generating a taxon set
2. Extracting FASTA data from our `ASSEMBLY/contigs` based on the taxon set

## Creating a data matrix configuration file

After we identify the UCE loci we enriched, but before we extract fasta data from our `ASSEMBLY/contigs` corresponding to those loci, we need to create a data matrix configuration file that denotes (1) which taxa we want to include in a given analysis and (2) which loci will be included with this taxon set.

The taxa included in the data matrix configuration file are determined by the user - you input a list of taxa you want to the analysis. The UCE loci included in the data matrix configuration are then determined by the software which compares the requested taxa to UCE match results in `probe.matches.sqlite` and two flags that you pass either one requesting **complete data matrix** or one requesting an **incomplete data matrix**.

**complete matrix** A phylogenetic matrix (typically sequence data) in which there are no missing data at any locus for any taxon/OTU.

**incomplete matrix** A phylogenetic matrix (typocally sequence data) in which data may be missing from a given taxon or a given loci (or both).

During the creation of the data matrix configuration file you can also include additional data from pre-existing UCE match databases and contigs (see `:ref :outgroup-data`).

We'll start very simply.

### Complete taxon set

First, let's generate a data matrix configuration file from only the current UCE enrichments that will be **complete** - meaning that we will not include loci where certain taxa have no data (either the locus was not enriched for that taxon or removed during the filtering process for duplicate loci).

To do this, you need to create a starting taxon-configuration file (a text-based file) denoting the taxa we want in the data set. The taxon-configuration file should look exactly like this (substitute in your taxon names):

```
[dataset1]
genus_species1
genus_species2
genus_species3
```

Let's assume you save this file as `datasets.conf`. Now, to create the data matrix configuration file from this taxon-configuration file, run:

```
# create the output directory for this taxon set
mkdir /path/to/uce/taxon-set1/

# create the data matrix configuration file
phyluce_assembly_get_match_counts \
  --locus-db /path/to/uce/output/probe.matches.sqlite \
  --taxon-list-config datasets.conf \
  --taxon-group 'dataset1' \
  --output /path/to/uce/taxon-set1/dataset1.conf
```

This will basically run a query against the database, and pull out those loci for those taxa in the `datasets.conf` file having UCE contigs.



## Results

The output printed to the screen and \$LOG file should look something like:

```

2014-04-24 17:25:08,145 - get_match_counts - INFO - ===== Starting get_
↳match_counts =====
2014-04-24 17:25:08,145 - get_match_counts - INFO - Version: git 7aec8f1
2014-04-24 17:25:08,145 - get_match_counts - INFO - Argument --extend_locus_db: None
2014-04-24 17:25:08,145 - get_match_counts - INFO - Argument --incomplete_matrix:
↳False
2014-04-24 17:25:08,146 - get_match_counts - INFO - Argument --keep_counts: False
2014-04-24 17:25:08,146 - get_match_counts - INFO - Argument --locus_db: /path/to/uce/
↳output/probes.matches.sqlite
2014-04-24 17:25:08,146 - get_match_counts - INFO - Argument --log_path: /path/to/uce
2014-04-24 17:25:08,146 - get_match_counts - INFO - Argument --optimize: False
2014-04-24 17:25:08,146 - get_match_counts - INFO - Argument --output: /path/to/uce/
↳taxon-set1/dataset1.conf
2014-04-24 17:25:08,146 - get_match_counts - INFO - Argument --random: False
2014-04-24 17:25:08,146 - get_match_counts - INFO - Argument --sample_size: 10
2014-04-24 17:25:08,146 - get_match_counts - INFO - Argument --samples: 10
2014-04-24 17:25:08,147 - get_match_counts - INFO - Argument --silent: False
2014-04-24 17:25:08,147 - get_match_counts - INFO - Argument --taxon_group: dataset1
2014-04-24 17:25:08,147 - get_match_counts - INFO - Argument --taxon_list_config:
↳datasets.conf
2014-04-24 17:25:08,147 - get_match_counts - INFO - Argument --verbosity: INFO
2014-04-24 17:25:08,150 - get_match_counts - INFO - There are 3 taxa in the taxon-
↳group '[dataset1]' in the config file dataset1.conf
2014-04-24 17:25:08,151 - get_match_counts - INFO - Getting UCE names from database
2014-04-24 17:25:08,407 - get_match_counts - INFO - There are 1314 total UCE loci in
↳the database
2014-04-24 17:25:11,046 - get_match_counts - INFO - Getting UCE matches by organism
↳to generate a COMPLETE matrix
2014-04-24 17:25:11,051 - get_match_counts - INFO - There are 306 shared UCE loci in
↳a COMPLETE matrix
2014-04-24 17:25:11,051 - get_match_counts - INFO - Failed to detect 428 UCE loci
↳in genus_species1
2014-04-24 17:25:11,051 - get_match_counts - INFO - Failed to detect 380 UCE loci
↳in genus_species2
2014-04-24 17:52:54,850 - get_match_counts - INFO - Writing the taxa and loci in the
↳data matrix to /path/to/uce/taxon-set1/dataset1.conf
2014-04-24 17:52:54,862 - get_match_counts - INFO - ===== Completed get_
↳match_counts =====

```

This basically says that although we've detected a total of 1,314 UCE loci in the 3 taxa in which we are interested, when we boil those down to a complete matrix, the complete matrix is only going to contain 306 UCE loci (of the 1,314). We had to drop 428 loci because we did not detect them in genus\_species1 and we had to drop another 380 loci because we did not detect them in genus\_species2.

The output written to the /path/to/uce/taxon-set1/dataset1.conf will look something like:

```

[Organisms]
genus_species1
genus_species2
genus_species3
[Loci]
uce-1005
uce-1018
uce-1025

```

```
uce-1028
uce-1042
uce-1055
uce-1060
uce-107
uce-1073
uce-1074
uce-1076
uce-108
...
```

### Taxon set membership and locus number

Now, you might think that increasing the locus count is simply a matter of removing *genus\_species1* from the list of taxa. This is not strictly true, however, given the vagaries of hits and misses among taxa. *phyluce\_assembly\_get\_match\_counts* has several other options to help you determine which taxa may be causing problems, but picking the best combination of taxa to give you the highest number of loci is a reasonably hard optimization problem.

### Incomplete data matrix

You may not always want a complete data matrix. Or generating a complete matrix drops too many loci for your tastes (it often does). In that case, you can easily generate an incomplete dataset using the following:

```
# create the data matrix configuration file
phyluce_assembly_get_match_counts \
  --locus-db /path/to/uce/output/probe.matches.sqlite \
  --taxon-list-config datasets.conf \
  --taxon-group 'dataset1' \
  --output /path/to/uce/taxon-set1/dataset1.conf \
  --incomplete-matrix
```

**Attention:** Note the addition of the `--incomplete-matrix` flag.

This will generate a dataset that includes any loci enriched across the taxa in the *datasets.conf* file.

**Note:** You do not determine the “completeness” of the final data matrix that you want to create during this stage - that happens later, after alignment (see *Finalize matrix completeness*). As a result, we are aligning data from any and all UCE loci having  $\geq 3$  taxa, which allows us to flexibly select the level of incompleteness later, without having to re-run our alignments.

### Creating additional data matrix configuration files for other analyses

If you want to generate/evaluate many data matrix configuration files containing different taxa, you can simply create new lists within the *datasets.conf* file like so:

```
[dataset1]
genus_species1
genus_species2
```

```
genus_species3

[dataset2]
genus_species2
genus_species3
genus_species4
genus_species5
genus_species6
```

And then you can run `phyluce_assembly_get_match_counts` against this new section to output the data matrix configuration files:

```
# create the data matrix configuration file
phyluce_assembly_get_match_counts \
  --locus-db /path/to/uce/output/probe.matches.sqlite \
  --taxon-list-config datasets.conf \
  --taxon-group 'dataset2' \
  --output /path/to/uce/taxon-set2/dataset2.conf
```

In this way, you can get some idea of how different taxon-set memberships affect the resulting data matrix configuration files *prior to* extracting the relevant FASTA data from `$ASSEMBLY/contigs` - which is a reasonably slow process.

### Incorporating outgroup/other data

You may want to include outgroup data from another source into your datasets. This can be from the pre-processed outgroup data files, but it doesn't need to be these outgroup data. These additional data can also be contigs previously assembled from a different set of taxa.

**Hint: ADVANCED:** If you want to include outgroup data from genome-enabled taxa, we have already created several repositories containing these data. We maintain these data under version control at: <https://github.com/faircloth-lab/uce-probe-sets>. To download these data and use them in your analyses, you can clone the data using git:

```
git clone https://github.com/faircloth-lab/uce-probe-sets
```

Then update your `--taxon-list-config` file and provide the proper paths to the cloned data, as detailed below.

The first step of this process is to setup your `--taxon-list-config` slightly differently - by indicating taxa from external data sources using asterisks:

```
[dataset3]
genus_species1
genus_species2
genus_species3
genus_species4*
genus_species5*
```

Here, `genus_species4*` and `genus_species5*` come from an external data source.

Then, you need to pass `phyluce_assembly_get_match_counts` the location of both **your** `--locus-db` and the `--extend-locus-db`. For example:

```
# create the data matrix configuration file
phyluce_assembly_get_match_counts \
  --locus-db /path/to/uce/output/probe.matches.sqlite \
  --taxon-list-config datasets.conf \
```

```
--taxon-group 'dataset3' \  
--extend-locus-db /path/to/some/other/probe.matches.sqlite \  
--output /path/to/uce/taxon-set3/dataset3.conf
```

To keep all this extension from getting too crazy, I've limited the ability to include external data to a single set. If you have lots of data from many different enrichments, you'll need to generate a *contigs* folder containing all these various assemblies (or symlinks to them), then align the probes to these data (see *Match contigs to probes*). Once you do that, you can extend your current data set with all of these other data.

## Extracting FASTA data using the data matrix configuration file

Once we have created the data matrix configuration file containing data for our taxa of interest and those loci of interest, we need to extract the appropriate FASTA sequences from each assembly representing the taxon/OTU of interest (e.g. in  $\$ASSEMBLY/contigs$ ). This is a reasonably straightforward process that differs only slightly based on whether you are extracting a complete matrix of data, an incomplete matrix of data, and/or whether you are incorporating any external data sources.

### Complete data matrix

To generate FASTA file containing the sequence data from a complete data matrix configuration file, run:

```
phyluce_assembly_get_fastas_from_match_counts \  
--contigs /path/to/assembly/contigs/ \  
--locus-db /path/to/uce/output/probe.matches.sqlite \  
--match-count-output /path/to/uce/taxon-set1/dataset1.conf \  
--output /path/to/uce/taxon-set1/dataset1.fasta
```

### Incomplete data matrix

Similarly, to generate a FASTA file containing the sequence data from a complete data matrix configuration file, run:

```
phyluce_assembly_get_fastas_from_match_counts \  
--contigs /path/to/assembly/contigs/ \  
--locus-db /path/to/uce/output/probe.matches.sqlite \  
--match-count-output /path/to/uce/taxon-set/dataset1.conf \  
--incomplete-matrix /path/to/uce/taxon-set1/dataset1.incomplete \  
--output /path/to/uce/taxon-set1/dataset1.fasta
```

**Attention:** Note the addition of the `--incomplete-matrix` option. This creates an output file that contains the names of the **missing** loci by taxon/OTU. You can name this file anything you like. I tend to use `.incomplete` as the extension so that it is clear what this file contains.

### Incorporating outgroup/other data

When we're incorporating external data, we need to pass the name of the external database as well as the name of the external *contigs*. To generate a FASTA file containing the sequence data from a complete data matrix configuration that includes external data sources, run:

```

phyluce_assembly_get_fastas_from_match_counts \
  --contigs /path/to/assembly/contigs/ \
  --locus-db /path/to/uce/output/probe.matches.sqlite \
  --match-count-output /path/to/uce/taxon-set1/dataset1.conf \
  --incomplete-matrix /path/to/uce/taxon-set1/dataset1.incomplete \
  --extend-locus-db /path/to/some/other/probe.matches.sqlite \
  --extend-locus-contigs /path/to/some/other/contigs \
  --output /path/to/uce/taxon-set3/dataset3.fasta

```

## Aligning and trimming FASTA data

With all of that out of the way, things get much easier to deal with. Now, we need to align our data across loci, and once we're done with that, the remaining operations we can run on the data are format-conversions, QC steps, matrix trimming for completeness, and any number of other fun things.

Aligning the amount of data generated by enrichment approaches is reasonably computationally intensive - so the alignment step goes fastest if you have a multicore machine. You also have several alignment options available, although I would suggest sticking with MAFFT.

**Attention:** The alignment process, as implemented by [phyluce](#), includes trimming steps that trim ragged edges and remove alignments that become too short following trimming.

To turn trimming off and trim alignments using another approach, pass the `--no-trim` option. There are also several more options related to trimming that you can tweak. To view these, run `phyluce_align_seqcap_align --help`.

## Complete data matrix

### Alignment

To align the loci, by taxon, in the FASTA file you just created, run:

```

phyluce_align_seqcap_align \
  --fasta /path/to/uce/taxon-set1/dataset1.fasta \
  --output /path/to/uce/taxon-set1/mafft-nexus/ \
  --taxa 3 \
  --aligner mafft \
  --cores 8

```

**Attention:** If you pass more `--cores` than your machine has, you will receive an error.

**Note:** Here, we are accepting the default, output alignment format ("nexus"). To change that format to something else, pass the `--output-format` option with a choice of {fasta,nexus,phylip,clustal,emboss,stockholm}.

## Alignment stats

Once you have alignments, it's nice to get a general sense of their length and composition. You can quickly (with a multicore machine) summarize thousands of alignments by running:

```
phyluce_align_get_align_summary_data \
  --alignments /path/to/uce/taxon-set1/mafft-nexus/ \
  --cores 12
```

This will produce output that looks similar to:

```
2014-04-24 17:31:15,724 - get_align_summary_data - INFO - ===== Starting_
↳get_align_summary_data =====
2014-04-24 17:31:15,724 - get_align_summary_data - INFO - Version: git 7aec8f1
2014-04-24 17:31:15,724 - get_align_summary_data - INFO - Argument --alignments: /
↳path/to/uce/taxon-set1/mafft-nexus/
2014-04-24 17:31:15,724 - get_align_summary_data - INFO - Argument --cores: 12
2014-04-24 17:31:15,724 - get_align_summary_data - INFO - Argument --input_format:
↳nexus
2014-04-24 17:31:15,724 - get_align_summary_data - INFO - Argument --log_path: /path/
↳to/uce/taxon-set1/log
2014-04-24 17:31:15,725 - get_align_summary_data - INFO - Argument --show_taxon_
↳counts: False
2014-04-24 17:31:15,725 - get_align_summary_data - INFO - Argument --verbosity: INFO
2014-04-24 17:31:15,725 - get_align_summary_data - INFO - Getting alignment files
2014-04-24 17:31:15,729 - get_align_summary_data - INFO - Computing summary_
↳statistics using 12 cores
2014-04-24 17:31:16,653 - get_align_summary_data - INFO - -----
↳Alignment summary -----
2014-04-24 17:31:16,654 - get_align_summary_data - INFO - [Alignments] loci:      306
2014-04-24 17:31:16,654 - get_align_summary_data - INFO - [Alignments] length: 223,
↳929
2014-04-24 17:31:16,654 - get_align_summary_data - INFO - [Alignments] mean:      731.79
2014-04-24 17:31:16,654 - get_align_summary_data - INFO - [Alignments] 95% CI:    17.01
2014-04-24 17:31:16,654 - get_align_summary_data - INFO - [Alignments] min:      275
2014-04-24 17:31:16,654 - get_align_summary_data - INFO - [Alignments] max:      1,109
2014-04-24 17:31:16,655 - get_align_summary_data - INFO - -----
↳Taxon summary -----
2014-04-24 17:31:16,655 - get_align_summary_data - INFO - [Taxa] mean:          27.00
2014-04-24 17:31:16,655 - get_align_summary_data - INFO - [Taxa] 95% CI:          0.00
2014-04-24 17:31:16,656 - get_align_summary_data - INFO - [Taxa] min:           27
2014-04-24 17:31:16,656 - get_align_summary_data - INFO - [Taxa] max:           27
2014-04-24 17:31:16,656 - get_align_summary_data - INFO - ----- Missing_
↳data from trim summary -----
2014-04-24 17:31:16,656 - get_align_summary_data - INFO - [Missing] mean:        7.61
2014-04-24 17:31:16,656 - get_align_summary_data - INFO - [Missing] 95% CI:      0.24
2014-04-24 17:31:16,656 - get_align_summary_data - INFO - [Missing] min:         1.13
2014-04-24 17:31:16,657 - get_align_summary_data - INFO - [Missing] max:         15.79
2014-04-24 17:31:16,661 - get_align_summary_data - INFO - -----
↳Character count summary -----
2014-04-24 17:31:16,661 - get_align_summary_data - INFO - [All characters]      6,046,
↳083
2014-04-24 17:31:16,661 - get_align_summary_data - INFO - [Nucleotides]        4,924,
↳129
2014-04-24 17:31:16,661 - get_align_summary_data - INFO - ----- Data_
↳matrix completeness summary -----
2014-04-24 17:31:16,661 - get_align_summary_data - INFO - [Matrix 50%]          306_
↳alignments
```

```

2014-04-24 17:31:16,661 - get_align_summary_data - INFO - [Matrix 55%] 306_
↳alignments
2014-04-24 17:31:16,662 - get_align_summary_data - INFO - [Matrix 60%] 306_
↳alignments
2014-04-24 17:31:16,662 - get_align_summary_data - INFO - [Matrix 65%] 306_
↳alignments
2014-04-24 17:31:16,662 - get_align_summary_data - INFO - [Matrix 70%] 306_
↳alignments
2014-04-24 17:31:16,662 - get_align_summary_data - INFO - [Matrix 75%] 306_
↳alignments
2014-04-24 17:31:16,662 - get_align_summary_data - INFO - [Matrix 80%] 306_
↳alignments
2014-04-24 17:31:16,662 - get_align_summary_data - INFO - [Matrix 85%] 306_
↳alignments
2014-04-24 17:31:16,662 - get_align_summary_data - INFO - [Matrix 90%] 306_
↳alignments
2014-04-24 17:31:16,662 - get_align_summary_data - INFO - [Matrix 95%] 306_
↳alignments
2014-04-24 17:31:16,663 - get_align_summary_data - INFO - -----_
↳Character counts -----
2014-04-24 17:31:16,663 - get_align_summary_data - INFO - [Characters] '-' is present_
↳651,009 times
2014-04-24 17:31:16,663 - get_align_summary_data - INFO - [Characters] '?' is present_
↳470,945 times
2014-04-24 17:31:16,663 - get_align_summary_data - INFO - [Characters] 'A' is present_
↳1,386,821 times
2014-04-24 17:31:16,663 - get_align_summary_data - INFO - [Characters] 'C' is present_
↳1,089,729 times
2014-04-24 17:31:16,663 - get_align_summary_data - INFO - [Characters] 'G' is present_
↳1,094,159 times
2014-04-24 17:31:16,663 - get_align_summary_data - INFO - [Characters] 'T' is present_
↳1,353,420 times
2014-04-24 17:31:16,664 - get_align_summary_data - INFO - ===== Completed_
↳get_align_summary_data =====

```

## Locus name removal

For historical reasons, and also for users to ensure that the sequence data aligned together are from the same loci, each sequence line in the alignment file output by `seqcap_align_2` contains the `genus_species1` designator, but the `genus_species1` designator is also prepended with the locus name (e.g. `uce-1005_genus_species1`). We need to remove these if we plan to concatenate the loci (*RAxML*). More generally, it is a good idea to remove locus names from sequence lines before running any analyses. To do this, run:

```

phyluce_align_remove_locus_name_from_nexus_lines \
  --alignments /path/to/uce/taxon-set1/mafft-nexus/ \
  --output /path/to/uce/taxon-set1/mafft-nexus-clean/ \
  --taxa 3

```

## Incomplete data matrix

### Alignment

The only difference for an alignment of incomplete data is that we also pass the `--incomplete-matrix` flag, which tells the code to expect that some loci will not contain data across all taxa:

```
phyluce_align_seqcap_align \
  --fasta /path/to/uce/taxon-set2/dataset2.fasta \
  --output /path/to/uce/taxon-set2/mafft-nexus/ \
  --taxa 34 \
  --aligner mafft \
  --cores 12 \
  --incomplete-matrix
```

## Alignment stats

Once you have alignments, it's nice to get a general sense of their length and composition. You can quickly (with a multicore machine) summarize thousands of alignments by running:

```
phyluce_align_get_align_summary_data \
  --alignments /path/to/uce/taxon-set1/mafft-nexus/ \
  --cores 12
```

This will produce output that looks similar to:

```
2014-04-24 20:11:18,208 - get_align_summary_data - INFO - ===== Starting_
↳get_align_summary_data =====
2014-04-24 20:11:18,209 - get_align_summary_data - INFO - Version: git 7aec8f1
2014-04-24 20:11:18,209 - get_align_summary_data - INFO - Argument --alignments: /
↳path/to/uce/taxon-set1/mafft-nexus/
2014-04-24 20:11:18,209 - get_align_summary_data - INFO - Argument --cores: 12
2014-04-24 20:11:18,209 - get_align_summary_data - INFO - Argument --input_format:_
↳nexus
2014-04-24 20:11:18,209 - get_align_summary_data - INFO - Argument --log_path: /path/
↳to/uce/taxon-set1/log
2014-04-24 20:11:18,209 - get_align_summary_data - INFO - Argument --show_taxon_
↳counts: False
2014-04-24 20:11:18,209 - get_align_summary_data - INFO - Argument --verbosity: INFO
2014-04-24 20:11:18,210 - get_align_summary_data - INFO - Getting alignment files
2014-04-24 20:11:18,253 - get_align_summary_data - INFO - Computing summary_
↳statistics using 12 cores
2014-04-24 20:11:20,573 - get_align_summary_data - INFO - -----_
↳Alignment summary -----
2014-04-24 20:11:20,574 - get_align_summary_data - INFO - [Alignments] loci:    1,104
2014-04-24 20:11:20,574 - get_align_summary_data - INFO - [Alignments] length:  752,
↳617
2014-04-24 20:11:20,574 - get_align_summary_data - INFO - [Alignments] mean:    681.72
2014-04-24 20:11:20,574 - get_align_summary_data - INFO - [Alignments] 95% CI:  13.03
2014-04-24 20:11:20,574 - get_align_summary_data - INFO - [Alignments] min:    169
2014-04-24 20:11:20,574 - get_align_summary_data - INFO - [Alignments] max:    4,520
2014-04-24 20:11:20,576 - get_align_summary_data - INFO - -----_
↳Taxon summary -----
2014-04-24 20:11:20,576 - get_align_summary_data - INFO - [Taxa] mean:          24.29
2014-04-24 20:11:20,576 - get_align_summary_data - INFO - [Taxa] 95% CI:          0.26
2014-04-24 20:11:20,576 - get_align_summary_data - INFO - [Taxa] min:            3
2014-04-24 20:11:20,576 - get_align_summary_data - INFO - [Taxa] max:           27
2014-04-24 20:11:20,577 - get_align_summary_data - INFO - ----- Missing_
↳data from trim summary -----
2014-04-24 20:11:20,577 - get_align_summary_data - INFO - [Missing] mean:          7.97
2014-04-24 20:11:20,577 - get_align_summary_data - INFO - [Missing] 95% CI:        0.16
2014-04-24 20:11:20,578 - get_align_summary_data - INFO - [Missing] min:          0.44
2014-04-24 20:11:20,578 - get_align_summary_data - INFO - [Missing] max:         19.71
```



```

2014-04-24 20:11:20,592 - get_align_summary_data - INFO - -----
↳Character count summary -----
2014-04-24 20:11:20,592 - get_align_summary_data - INFO - [All characters]      18,
↳541,550
2014-04-24 20:11:20,592 - get_align_summary_data - INFO - [Nucleotides]      14,
↳713,956
2014-04-24 20:11:20,594 - get_align_summary_data - INFO - ----- Data
↳matrix completeness summary -----
2014-04-24 20:11:20,594 - get_align_summary_data - INFO - [Matrix 50%]      1048
↳alignments
2014-04-24 20:11:20,594 - get_align_summary_data - INFO - [Matrix 55%]      1044
↳alignments
2014-04-24 20:11:20,594 - get_align_summary_data - INFO - [Matrix 60%]      1035
↳alignments
2014-04-24 20:11:20,594 - get_align_summary_data - INFO - [Matrix 65%]      1027
↳alignments
2014-04-24 20:11:20,594 - get_align_summary_data - INFO - [Matrix 70%]      1024
↳alignments
2014-04-24 20:11:20,594 - get_align_summary_data - INFO - [Matrix 75%]      1010
↳alignments
2014-04-24 20:11:20,594 - get_align_summary_data - INFO - [Matrix 80%]      998
↳alignments
2014-04-24 20:11:20,595 - get_align_summary_data - INFO - [Matrix 85%]      994
↳alignments
2014-04-24 20:11:20,595 - get_align_summary_data - INFO - [Matrix 90%]      906
↳alignments
2014-04-24 20:11:20,595 - get_align_summary_data - INFO - [Matrix 95%]      794
↳alignments
2014-04-24 20:11:20,595 - get_align_summary_data - INFO - -----
↳Character counts -----
2014-04-24 20:11:20,595 - get_align_summary_data - INFO - [Characters] '-' is present
↳2,301,454 times
2014-04-24 20:11:20,595 - get_align_summary_data - INFO - [Characters] '?' is present
↳1,526,140 times
2014-04-24 20:11:20,595 - get_align_summary_data - INFO - [Characters] 'A' is present
↳4,092,085 times
2014-04-24 20:11:20,596 - get_align_summary_data - INFO - [Characters] 'C' is present
↳3,267,550 times
2014-04-24 20:11:20,596 - get_align_summary_data - INFO - [Characters] 'G' is present
↳3,286,742 times
2014-04-24 20:11:20,596 - get_align_summary_data - INFO - [Characters] 'T' is present
↳4,067,579 times
2014-04-24 20:11:20,596 - get_align_summary_data - INFO - ===== Completed
↳get_align_summary_data =====

```

**Note:** The alignment summary stats give you some idea of data matrix composition at varying levels of completeness in the Data matrix completeness summary section.

### Locus name removal

For historical reasons, and also for users to ensure that the sequence data aligned together are from the same loci, each sequence line in the alignment file output by seqcap\_align\_2 contains the genus\_species1 designator, but the genus\_species1 designator is also prepended with the locus name (e.g. uce-1005\_genus\_species1). We need to remove these if we plan to concatenate the loci (*RAxML*). More generally, it is a good idea to remove locus

names from sequence lines before running any analyses. To do this, run:

```
phyluce_align_remove_locus_name_from_nexus_lines \  
  --alignments /path/to/uce/taxon-set1/mafft-nexus/ \  
  --output /path/to/uce/taxon-set1/mafft-nexus-clean/ \  
  --cores 12
```

## Finalize matrix completeness

After checking the resulting alignment summary stats and checking your alignments for quality, you will generally want to cull the data set to reach your desired level of completeness. That is easily done by running the following, while inputting the set of alignments just generated using:

```
# the integer following --taxa is the number of TOTAL taxa  
phyluce_align_get_only_loci_with_min_taxa \  
  --alignments /path/to/uce/taxon-set1/mafft-nexus-clean/ \  
  --taxa 34 \  
  --percent 0.75 \  
  --output /path/to/uce/taxon-set1/mafft-nexus-min-25-taxa/ \  
  --cores 12
```

**Attention:** This program computes the floor( $\text{taxa} * \text{percent}$ ) and uses the resulting number to determine the min( $\text{taxa}$ ) allowed in an alignment of `--percent` completeness.

This will produce output that looks similar to:

```
2014-04-24 20:12:33,386 - get_only_loci_with_min_taxa - INFO - =====  
↳ Starting get_only_loci_with_min_taxa =====  
2014-04-24 20:12:33,387 - get_only_loci_with_min_taxa - INFO - Version: git 7aec8f1  
2014-04-24 20:12:33,387 - get_only_loci_with_min_taxa - INFO - Argument --alignments:   
↳ /path/to/uce/taxon-set1/mafft-nexus  
2014-04-24 20:12:33,387 - get_only_loci_with_min_taxa - INFO - Argument --cores: 12  
2014-04-24 20:12:33,387 - get_only_loci_with_min_taxa - INFO - Argument --input_  
↳ format: nexus  
2014-04-24 20:12:33,387 - get_only_loci_with_min_taxa - INFO - Argument --log_path:   
↳ None  
2014-04-24 20:12:33,387 - get_only_loci_with_min_taxa - INFO - Argument --output: /  
↳ path/to/uce/taxon-set1/mafft-nexus-min-25-taxa  
2014-04-24 20:12:33,388 - get_only_loci_with_min_taxa - INFO - Argument --percent: 0.  
↳ 75  
2014-04-24 20:12:33,388 - get_only_loci_with_min_taxa - INFO - Argument --taxa: 27  
2014-04-24 20:12:33,388 - get_only_loci_with_min_taxa - INFO - Argument --verbosity:   
↳ INFO  
2014-04-24 20:12:33,388 - get_only_loci_with_min_taxa - INFO - Getting alignment files  
2014-04-24 20:12:35,293 - get_only_loci_with_min_taxa - INFO - Copied 1010 alignments   
↳ of 1104 total containing >= 0.75 proportion of taxa (n = 20)  
2014-04-24 20:12:35,294 - get_only_loci_with_min_taxa - INFO - =====  
↳ Completed get_only_loci_with_min_taxa =====
```

## Add missing data designators

Finally, you will need to add missing data designators for taxa missing from each alignment of a given locus. This will basically allow you to generate concatenated data sets and it may reduce error messages from other programs about

files having unequal numbers of taxa. To do this, run:

```
phyluce_align_add_missing_data_designators \
  --alignments /path/to/uce/taxon-set1/mafft-nexus-min-25-taxa \
  --output /path/to/uce/taxon-set1/mafft-nexus-min-25-taxa \
  --match-count-output /path/to/uce/taxon-set/dataset1.conf \
  --incomplete-matrix /path/to/uce/taxon-set1/dataset1.incomplete \
  --log-path log \
  --cores 12
```

**Note:** Here, we're inputting the `--match-count-output` and the `--incomplete-matrix` we created earlier in the *Incomplete data matrix* and *Extracting FASTA data using the data matrix configuration file* sections.

## Operations on alignments

Many workflows for phylogenetics simply involve converting one alignment format to another or changing something about the contents of a given alignment. We use many of these manipulations in the next section (see *Preparing alignment data for analysis*), as well.

### Converting one alignment format to another

To convert one alignment type (e.g., nexus) to another (e.g., fasta), we have a relative simple bit of code to achieve that process. You can greatly speed this processing step up on a multicore machine with the `--cores` option:

```
phyluce_align_convert_one_align_to_another \
  --alignments /path/to/uce/taxon-set1/mafft-nexus \
  --output /path/to/uce/taxon-set1/mafft-fasta \
  --input-format nexus \
  --output-format fasta \
  --cores 8 \
  --log-path log
```

You can convert from/to:

1. fasta
2. nexus
3. phylip
4. clustal
5. emboss
6. stockholm

### Shortening taxon names

You can shorten taxon names (e.g. for use with strict phylip) by modifying the above command slightly to add `--shorten-names`:

```
phyluce_align_convert_one_align_to_another \
  --alignments /path/to/uce/taxon-set1/mafft-nexus \
  --output /path/to/uce/taxon-set1/mafft-fasta-shortnames \
  --input-format nexus \
```

```
--output-format fasta \  
--cores 8 \  
--shorten-names \  
--log-path log
```

### Excluding loci or taxa

You may want to exclude loci less than a certain length or having fewer than a particular number of taxa, or only containing certain taxa. You can accomplish that using:

```
phyluce_align_filter_alignments \  
  --alignments /path/to/uce/taxon-set1/mafft-nexus \  
  --output /path/to/a/new/directory \  
  --input-format nexus \  
  --containing-data-for genus_species1 genus_species2 \  
  --min-length 100 \  
  --min-taxa 5 \  
  --log-path log
```

This will filter alignments that do not contain the taxa requested, those alignments shorter than 100 bp, and those alignments having fewer than 5 taxa (taxa with only missing data are not counted).

### Extracting taxon data from alignments

Sometimes you may have alignments from which you want to extract data from a given taxon, format the alignment string as fasta, and do something with the fasta results:

```
phyluce_align_extract_taxon_fasta_from_alignments \  
  --alignments /path/to/uce/taxon-set1/mafft-nexus \  
  --taxon genus_species1 \  
  --output /path/to/output/file.fasta
```

### Preparing alignment data for analysis

Formatting data for analysis generally involves slight differences from the steps described above. There are several application-specific programs in [phyluce](#).

#### RAxML

For RAxML, you need a concatenated phylip file. This is pretty easily created if you have an input directory of nexus alignments. To create a concatenated phylip file from run:

```
phyluce_align_format_nexus_files_for_raxml \  
  --alignments /path/to/uce/taxon-set1/mafft-nexus \  
  --output /path/to/uce/taxon-set1/mafft-raxml
```

This will output a concatenated file named `mafft-raxml.phylip` in `/path/to/uce/taxon-set1/mafft-raxml`.

## PHYLIP/CloudForest

PHYLIP, PhyML, and other programs like [CloudForest](#) require input files to be in strict phylip format for analysis. Converting alignment files to this format was discussed above, and is simple a matter of:

```
phyluce_align_convert_one_align_to_another \
  --alignments /path/to/uce/taxon-set1/mafft-nexus \
  --output /path/to/uce/taxon-set1/mafft-phylip-shortnames \
  --input-format nexus \
  --output-format phylip \
  --cores 8 \
  --shorten-names \
  --log-path log
```

## MrBayes

MrBayes is a little more challenging to run. This is largely due to the fact that we usually estimate the substitution models for all loci, then we partition loci by substitution model, concatenate the data, and format an appropriate file to be input to MrBayes.

The tricky part of this process is estimating the locus-specific substitution models. Generally speaking, I do this with [CloudForest](#) now, then I strip the best-fitting substitution model from the [CloudForest](#) output, and input that file to the program that creates a nexus file for MrBayes.

First, estimate the substitution models using cloudforest (this will also give you genetrees for all loci, as a bonus). You will need your alignments in strict phylip format:

```
python cloudforest/cloudforest_mpi.py \
  /path/to/strict/phylip/alignments/ \
  /path/to/store/cloudforest/output/ \
  genetrees \
  $HOME/git/cloudforest/cloudforest/binaries/PhyML3linux64 \
  --parallelism multiprocessing \
  --cores 8
```

In the above, *genetrees* is a keyword that tells [CloudForest](#) that you mean to estimate genetrees (instead of bootstraps). Depending on the size of your dataset (and computer), this may take some time. Once this is done:

```
phyluce_genetrees_split_models_from_genetrees \
  --genetrees /path/to/cloudforest/output/genetrees.tre \
  --output /path/to/output_models.txt
```

Now, you're ready to go with formatting for MrBayes - note that we're inputting the path of the models file created above (*output\_models.txt*) on line 3:

```
phyluce_align_format_nexus_files_for_mrbayes \
  --alignments /path/to/input/nexus/ \
  --models /path/to/output_models.txt \
  --output /path/to/output/mrbayes.nexus \
  --interleave \
  --unlink
```

This should create a partitioned data file for you. The partitioning will be by model, not by locus. Should you want to fully partition by locus (which may overparameterize), then you can run:

```
phyluce_align_format_nexus_files_for_mrbayes \
  /path/to/input/nexus/ \
```

```
/path/to/output_models.txt \  
/path/to/output/mrbayes.nexus \  
--interleave \  
--unlink \  
--fully-partition
```

### CloudForest (genetree/species tree)

CloudForest is a program written by Nick Crawford and myself that helps you estimate genetrees and perform bootstrap replicates for very large datasets. Data input to CloudForest should be in strict phylip format (see *PHYLIP/CloudForest*). First, as above, run genetree analysis on your data (if you ran this above, you don't need to run it again). This will estimate the genetrees for each locus in your dataset, using its best fitting substitution model):

```
python cloudforest/cloudforest_mpi.py \  
/path/to/strict/phylip/alignments/ \  
/path/to/store/cloudforest/output/ \  
genetrees \  
$HOME/git/cloudforest/cloudforest/binaries/PhyML3linux64 \  
--parallelism multiprocessing \  
--cores 8
```

The, to generate bootstrap replicates, you can run:

```
python cloudforest/cloudforest_mpi.py \  
/path/to/strict/phylip/alignments/ \  
/path/to/store/cloudforest/output/ \  
bootstraps \  
$HOME/git/cloudforest/cloudforest/binaries/PhyML3linux64 \  
--parallelism multiprocessing \  
--cores 8 \  
--bootreps 1000 \  
--genetrees /path/to/store/cloudforest/output/genetrees.tre
```

**NOTE** that depending on your system, you may need to choose another value for the path to PhyML:

```
$HOME/git/cloudforest/cloudforest/binaries/PhyML3linux64
```

### RaXML (genetree/species tree)

We can also use RaXML to generate gene trees to turn into a species tree. To keep the taxa names similar to what I run through CloudForest, I usually input strict phylip formatted files to these runs (see *PHYLIP/CloudForest*). Once that's done, you can generate genetrees with:

```
phyluce_genetrees_run_raxml_genetrees \  
--alignments /path/to/strict/phylip/alignments/ \  
--output /path/to/store/raxml/output/ \  
--outgroup genus_species1 \  
--cores 12 \  
--threads 1
```

Number of *-cores* is the number of simultaneous trees to estimate, while *-threads* is the number of threads to use for each tree. Although somewhat counterintuitive, I've found that 1 *-thread* per locus and many loci being processed at once is the fastest route to go.

Once that's finished, you can generate bootstrap replicates for those same loci:

```
.. code-block:: bash
```

```
phyluce_genetrees_run_raxml_bootstraps -alignments /path/to/strict/phylip/alignments/ -output
/path/to/store/raxml/output/ -bootreps 100 -outgroup genus_species1 -cores 12 -threads 1
```

## Tutorial I: UCE Phylogenomics

In the following example, we are going to process raw read data from UCE enrichments performed against several divergent taxa so that you can get a feel for how a typical analysis goes. I'm also going to use several tricks that I did not cover in the *UCE Processing for Phylogenomics* section.

The taxa we are working with will be:

- Mus musculus (PE100)
- Anolis carolinensis (PE100)
- Alligator mississippiensis (PE150)
- Gallus gallus (PE250)

### Download the data

You can download the data from figshare (<http://dx.doi.org/10.6084/m9.figshare.1284521>). If you want to use the command line, you can use something like:

```
# create a project directory
mkdir uce-tutorial

# change to that directory
cd uce-tutorial

# download the data into a file names fastq.zip
wget -O fastq.zip https://ndownloader.figshare.com/articles/1284521/versions/1

# make a directory to hold the data
mkdir raw-fastq

# move the zip file into that directory
mv fastq.zip raw-fastq

# move into the directory we just created
cd raw-fastq

# unzip the fastq data
unzip fastq.zip

# delete the zip file
rm fastq.zip

# you should see 6 files in this directory now
ls -l

-rw-r--r--. 1 bcf data 152M Apr 11 14:03 Alligator_mississippiensis_GGAGCTATGG_L001_
↪R1_001.fastq.gz
```

```
-rw-r--r--. 1 bcf data 147M Apr 11 14:03 Alligator_mississippiensis_GGAGCTATGG_L001_
↪R2_001.fastq.gz
-rw-r--r--. 1 bcf data 173M Apr 11 14:03 Anolis_carolinensis_GGCGAAGGTT_L001_R1_001.
↪fastq.gz
-rw-r--r--. 1 bcf data 174M Apr 11 14:03 Anolis_carolinensis_GGCGAAGGTT_L001_R2_001.
↪fastq.gz
-rw-r--r--. 1 bcf data 67M Apr 11 14:03 Gallus_gallus_TTCTCCTTCA_L001_R1_001.fastq.gz
-rw-r--r--. 1 bcf data 73M Apr 11 14:03 Gallus_gallus_TTCTCCTTCA_L001_R2_001.fastq.gz
-rw-r--r--. 1 bcf data 122M Apr 11 14:03 Mus_musculus_CTACAACGGC_L001_R1_001.fastq.gz
-rw-r--r--. 1 bcf data 121M Apr 11 14:03 Mus_musculus_CTACAACGGC_L001_R2_001.fastq.gz
```

Alternatively, if you think of the filesystem as a tree-like structure, the directory in which we are working (*uce-tutorial*) would look like:

```
uce-tutorial
+- raw-fastq
  +- Alligator_mississippiensis_GGAGCTATGG_L001_R1_001.fastq.gz
  +- Alligator_mississippiensis_GGAGCTATGG_L001_R2_001.fastq.gz
  +- Anolis_carolinensis_GGCGAAGGTT_L001_R1_001.fastq.gz
  +- Anolis_carolinensis_GGCGAAGGTT_L001_R2_001.fastq.gz
  +- Gallus_gallus_TTCTCCTTCA_L001_R1_001.fastq.gz
  +- Gallus_gallus_TTCTCCTTCA_L001_R2_001.fastq.gz
  +- Mus_musculus_CTACAACGGC_L001_R1_001.fastq.gz
  +- Mus_musculus_CTACAACGGC_L001_R2_001.fastq.gz
```

If you do not want to use the command line, you can download the data using the figshare interface or by clicking:

<http://downloads.figshare.com/article/public/1284521>

## Count the read data

Usually, we want a count of the actual number of reads in a given sequence file for a given species. As mentioned in the *UCE Processing for Phylogenomics* section, we can do this several ways. We'll use tools from unix, because they are fast. The next line of code will count the lines in each R1 file (which should be equal to the reads in the R2 file) and divide that number by 4 to get the number of sequence reads.

```
for i in *_R1*.fastq.gz; do echo $i; gunzip -c $i | wc -l | awk '{print $1/4}'; done
```

You should see:

```
Alligator_mississippiensis_GGAGCTATGG_L001_R1_001.fastq.gz
1750000
Anolis_carolinensis_GGCGAAGGTT_L001_R1_001.fastq.gz
1874362
Gallus_gallus_TTCTCCTTCA_L001_R1_001.fastq.gz
376559
Mus_musculus_CTACAACGGC_L001_R1_001.fastq.gz
1298196
```

## Clean the read data

The data you just downloaded are actual, raw, untrimmed fastq data. This means they contain adapter contamination and low quality bases. We need to remove these - which you can do several ways. We'll use another program that I wrote (*illumiprocessor*) because it allows us to trim many different indexed adapters from individual-specific fastq



files - something that is a pain to do by hand. That said, you can certainly trim your reads however you would like. See the [illumiprocessor](#) website for instructions on installing the program.

To use this program, we will create a configuration file that we will use to inform the program about which adapters are in which READ1 and READ2 files. The data we are trimming, here, are from TruSeq v3 libraries, but the indexes are 10 nucleotides long. We will set up the trimming file with these parameters, but please see the [illumiprocessor](#) documentation for other options.

```
# this is the section where you list the adapters you used. the asterisk
# will be replaced with the appropriate index for the sample.
[adapters]
i7:AGATCGGAAGAGCACACGTCTGAACTCCAGTCAC*ATCTCGTATGCCGTCTTCTGCTTG
i5:AGATCGGAAGAGCGCTCGTGTAGGGAAAGAGTGTAGATCTCGGTGGTCGCCGTATCATT

# this is the list of indexes we used
[tag sequences]
BFIDT-166:GGAGCTATGG
BFIDT-016:GGCGAAGGTT
BFIDT-045:TTCTCCTTCA
BFIDT-011:CTACAACGGC

# this is how each index maps to each set of reads
[tag map]
Alligator_mississippiensis_GGAGCTATGG:BFIDT-166
Anolis_carolinensis_GGCGAAGGTT:BFIDT-016
Gallus_gallus_TTCTCCTTCA:BFIDT-045
Mus_musculus_CTACAACGGC:BFIDT-011

# we want to rename our read files something a bit more nice - so we will
# rename Alligator_mississippiensis_GGAGCTATGG to alligator_mississippiensis
[names]
Alligator_mississippiensis_GGAGCTATGG:alligator_mississippiensis
Anolis_carolinensis_GGCGAAGGTT:anolis_carolinensis
Gallus_gallus_TTCTCCTTCA:gallus_gallus
Mus_musculus_CTACAACGGC:mus_musculus
```

I create this file in a directory **above** the one holding my reads, so the structure looks like:

```
uce-tutorial
+- illumiprocessor.conf
+- raw-fastq
  +- Alligator_mississippiensis_GGAGCTATGG_L001_R1_001.fastq.gz
  +- Alligator_mississippiensis_GGAGCTATGG_L001_R2_001.fastq.gz
  +- Anolis_carolinensis_GGCGAAGGTT_L001_R1_001.fastq.gz
  +- Anolis_carolinensis_GGCGAAGGTT_L001_R2_001.fastq.gz
  +- Gallus_gallus_TTCTCCTTCA_L001_R1_001.fastq.gz
  +- Gallus_gallus_TTCTCCTTCA_L001_R2_001.fastq.gz
  +- Mus_musculus_CTACAACGGC_L001_R1_001.fastq.gz
  +- Mus_musculus_CTACAACGGC_L001_R2_001.fastq.gz
```

Now I run [illumiprocessor](#) against the data. Note that I am using **4 physical CPU cores** to do this work. You need to use the number of physical cores available on your machine.

```
# go to the directory containing our config file and data
cd my-analysis

# run illumiprocessor

illumiprocessor \
```

```
--input raw-fastq/ \  
--output clean-fastq \  
--config illumiprocessor.conf \  
--cores 4
```

The output should look like the following:

```
2015-04-11 14:23:57,912 - illumiprocessor - INFO - ===== Starting_  
↳illumiprocessor =====  
2015-04-11 14:23:57,913 - illumiprocessor - INFO - Version: 2.0.6  
2015-04-11 14:23:57,913 - illumiprocessor - INFO - Argument --config: illumiprocessor.  
↳conf  
2015-04-11 14:23:57,913 - illumiprocessor - INFO - Argument --cores: 4  
2015-04-11 14:23:57,913 - illumiprocessor - INFO - Argument --input: /scratch/  
↳bfaircloth-uce-tutorial/raw-fastq  
2015-04-11 14:23:57,913 - illumiprocessor - INFO - Argument --log_path: None  
2015-04-11 14:23:57,913 - illumiprocessor - INFO - Argument --min_len: 40  
2015-04-11 14:23:57,914 - illumiprocessor - INFO - Argument --no_merge: False  
2015-04-11 14:23:57,914 - illumiprocessor - INFO - Argument --output: /scratch/  
↳bfaircloth-uce-tutorial/clean-fastq  
2015-04-11 14:23:57,914 - illumiprocessor - INFO - Argument --phred: phred33  
2015-04-11 14:23:57,914 - illumiprocessor - INFO - Argument --r1_pattern: None  
2015-04-11 14:23:57,914 - illumiprocessor - INFO - Argument --r2_pattern: None  
2015-04-11 14:23:57,914 - illumiprocessor - INFO - Argument --se: False  
2015-04-11 14:23:57,914 - illumiprocessor - INFO - Argument --trimmomatic: /home/bcf/  
↳anaconda/jar/trimmomatic.jar  
2015-04-11 14:23:57,915 - illumiprocessor - INFO - Argument --verbosity: INFO  
2015-04-11 14:23:57,942 - illumiprocessor - INFO - Trimming samples with Trimmomatic  
Running....  
2015-04-11 14:25:17,714 - illumiprocessor - INFO - ===== Completed_  
↳illumiprocessor =====
```

Notice that the program has created a *log* file showing what it did, and it has also created a new directory holding the clean data that has the name *clean-fastq* (what you told it to name the directory). Within that new directory, there are taxon-specific folder for the cleaned reads. More specifically, your directory structure should look similar to the following (I've collapsed the list of raw-reads):

```
uce-tutorial  
+- clean-fastq  
| +- alligator_mississippiensis  
| +- anolis_carolinensis  
| +- gallus_gallus  
| +- mus_musculus  
+- illumiprocessor.conf  
+- illumiprocessor.log  
+- raw-fastq
```

Within each organism specific directory, there are more files and folders:

```
uce-tutorial  
+- clean-fastq  
| +- alligator_mississippiensis  
| | +- adapters.fasta  
| | +- raw-reads  
| | +- split-adapter-quality-trimmed  
| | +- stats  
| +- anolis_carolinensis  
| | +- adapters.fasta
```

```

| | +- raw-reads
| | +- split-adapter-quality-trimmed
| | +- stats
| +- gallus_gallus
| | +- adapters.fasta
| | +- raw-reads
| | +- split-adapter-quality-trimmed
| | +- stats
| +- mus_musculus
|   +- adapters.fasta
|   +- raw-reads
|   +- split-adapter-quality-trimmed
|   +- stats
+- illumiprocessor.conf
+- illumiprocessor.log
+- raw-fastq

```

And, within each of those directories nested within the species-specific directory, there are additional files or links to files:

```

uce-tutorial
+- clean-fastq
| +- alligator_mississippiensis
| | +- adapters.fasta
| | +- raw-reads
| | | +- alligator_mississippiensis-READ1.fastq.gz -> <PATH>
| | | +- alligator_mississippiensis-READ2.fastq.gz -> <PATH>
| | +- split-adapter-quality-trimmed
| | | +- alligator_mississippiensis-READ1.fastq.gz
| | | +- alligator_mississippiensis-READ2.fastq.gz
| | | +- alligator_mississippiensis-READ-singleton.fastq.gz
| | +- stats
| |   +- alligator_mississippiensis-adapter-contam.txt
| +- anolis_carolinensis
| +- gallus_gallus
| +- mus_musculus
+- illumiprocessor.conf
+- illumiprocessor.log
+- raw-fastq

```

I have collapsed the listing to show only the first taxon.

The `->` in the `raw-reads` directory above means there are [symlinks](#) to the files. I have removed the file paths and replaced them with `<PATH>` so that the figure will fit on a page.

The really important information is in the `split-adapter-quality-trimmed` directory - which now holds our reads that have had adapter-contamination and low-quality bases removed. Within this `split-adapter-quality-trimmed` directory, the `READ1` and `READ2` files hold reads that remain in a pair (the reads are in the same consecutive order in each file). The `READ-singleton` file holds `READ1` reads **OR** `READ2` reads that lost their “mate” or “paired-read” because of trimming or removal.

## Quality control

You might want to get some idea of what effect the trimming has on read counts and overall read lengths. There are certainly other (better) tools out there to do this (like [FastQC](#)), but you can get a reasonable idea of how good your reads are by running the following, which will output a CSV listing of read stats by sample:

```
# move to the directory holding our cleaned reads
cd clean-fastq/

# run this script against all directories of reads

for i in *;
do
    phyluce_assembly_get_fastq_lengths --input $i/split-adapter-quality-trimmed/ --
    ↪ csv;
done
```

The output you see should look like this:

```
# sample,reads,total bp,mean length, 95 CI length,min,max,median
All files in dir with alligator_mississippiensis-READ1.fastq.gz,3279362,294890805,89.
↪ 9232243955,0.01003996813,40,100,100.0
All files in dir with anolis_carolinensis-READ2.fastq.gz,3456457,314839345,91.
↪ 0873026917,0.00799863728974,40,100,100.0
All files in dir with gallus_gallus-READ2.fastq.gz,749026,159690692,213.197795537,0.
↪ 0588973605567,40,251,250.0
All files in dir with mus_musculus-READ-singleton.fastq.gz,2332785,211828511,90.
↪ 8049867433,0.0102813002698,40,100,100.0
```

Now, we're ready to assemble our reads.

## Assemble the data

phyluce has a lot of options for assembly - you can use [velvet](#), [abyss](#), or [trinity](#). For this tutorial, we are going to use [trinity](#), because I believe it works best for most purposes. The helper programs for the other assemblers use the same config file, too, so you can easily experiment with all of the assemblers.

To run an assembly, we need to create a another configuration file. The assembly configuration file looks like the following, assuming we want to assemble all of our data from the organisms above:

```
[samples]
alligator_mississippiensis:/scratch/bfaircloth-uce-tutorial/clean-fastq/alligator_
↪ mississippiensis/split-adapter-quality-trimmed/
anolis_carolinensis:/scratch/bfaircloth-uce-tutorial/clean-fastq/anolis_carolinensis/
↪ split-adapter-quality-trimmed/
gallus_gallus:/scratch/bfaircloth-uce-tutorial/clean-fastq/gallus_gallus/split-
↪ adapter-quality-trimmed/
mus_musculus:/scratch/bfaircloth-uce-tutorial/clean-fastq/mus_musculus/split-adapter-
↪ quality-trimmed/
```

We will save this into a file named *assembly.conf* at the top of our *uce-tutorial* directory:

```
uce-tutorial
+- assembly.conf
+- clean-fastq
+- illumiprocessor.conf
+- illumiprocessor.log
+- phyluce_assembly_assemblo_trinity.log
+- raw-fastq
+- trinity-assemblies
```

If you want to change the names on the left hand side of the colon in the config file, you can do so, but the *PATHs* on the right hand side need to point to our “clean” UCE raw reads. If you have files in multiple locations, you can use

different *PATHs* on the right-hand side.

**Attention:** Although you can easily input new *PATHs* in this file, the **structure** of the data below the *PATH* you use must be the same - meaning that the structure and naming scheme for *READ1*, *READ2*, and *READ-singleton* must be the same. Or, put another way, it must look like the following:

```
some-random-data
+- clean-fastq
  +- alligator_mississippiensis
    | +- split-adapter-quality-trimmed
    |   +- alligator_mississippiensis-READ1.fastq.gz
    |   +- alligator_mississippiensis-READ2.fastq.gz
    |   +- alligator_mississippiensis-READ-singleton.fastq.gz
  +- anolis_carolinensis
    +- split-adapter-quality-trimmed
      +- anolis_carolinensis-READ1.fastq.gz
      +- anolis_carolinensis-READ2.fastq.gz
      +- anolis_carolinensis-READ-singleton.fastq.gz
```

Now that we have that file created, copy it to our working directory, and run the *phyluce\_assembly\_assemblo\_trinity* program:

```
# make sure we are at the top-level of our uce tutorial directory
cd uce-tutorial

# run the assembly
phyluce_assembly_assemblo_trinity \
  --conf assembly.conf \
  --output trinity-assemblies \
  --clean \
  --cores 12
```

**Warning:** Note that I am using 12 physical CPU cores to do this work. You need to use the number of physical cores available on *your* machine. The *phyluce.conf* file also assumes you have **at least 8 GB** of RAM on your system, and it is better to have much more. If you use more CPU cores than you have or you specify more RAM than you have, bad things will happen.

As the assembly proceeds, you should see output similar to the following:

```
2015-04-11 18:22:41,128 - phyluce_assembly_assemblo_trinity - INFO - -----
↳--- Processing gallus_gallus ----
2015-04-11 15:30:54,183 - phyluce_assembly_assemblo_trinity - INFO - Argument --dir:
↳None
2015-04-11 15:30:54,183 - phyluce_assembly_assemblo_trinity - INFO - Argument --log_
↳path: None
2015-04-11 15:30:54,183 - phyluce_assembly_assemblo_trinity - INFO - Argument --min_
↳kmer_coverage: 2
2015-04-11 15:30:54,183 - phyluce_assembly_assemblo_trinity - INFO - Argument --
↳output: /scratch/bfaircloth-uce-tutorial/trinity-assemblies
2015-04-11 15:30:54,183 - phyluce_assembly_assemblo_trinity - INFO - Argument --
↳subfolder:
2015-04-11 15:30:54,183 - phyluce_assembly_assemblo_trinity - INFO - Argument --
↳verbosity: INFO
2015-04-11 15:30:54,184 - phyluce_assembly_assemblo_trinity - INFO - Getting input_
↳filenames and creating output directories
```

```

2015-04-11 15:30:54,186 - phyluce_assembly_assemblo_trinity - INFO - -----
↳Processing alligator_mississippiensis -----
2015-04-11 15:30:54,186 - phyluce_assembly_assemblo_trinity - INFO - Finding fastq/
↳fasta files
2015-04-11 15:30:54,189 - phyluce_assembly_assemblo_trinity - INFO - File type is
↳fastq
2015-04-11 15:30:54,190 - phyluce_assembly_assemblo_trinity - INFO - Copying raw read
↳data to /scratch/bfaircloth-uce-tutorial/trinity-assemblies/alligator_
↳mississippiensis_trinity
2015-04-11 15:30:54,561 - phyluce_assembly_assemblo_trinity - INFO - Combining
↳singleton reads with R1 data
2015-04-11 15:30:54,576 - phyluce_assembly_assemblo_trinity - INFO - Running Trinity.
↳pl for PE data
2015-04-11 16:29:19,127 - phyluce_assembly_assemblo_trinity - INFO - Removing
↳extraneous Trinity files
2015-04-11 16:29:20,957 - phyluce_assembly_assemblo_trinity - INFO - Symlinking
↳assembled contigs into /scratch/bfaircloth-uce-tutorial/trinity-assemblies/contigs
2015-04-11 16:29:20,957 - phyluce_assembly_assemblo_trinity - INFO - -----
↳ Processing anolis_carolinensis -----
...[continued]...
2015-04-11 22:30:19,558 - phyluce_assembly_assemblo_trinity - INFO - =====
↳Completed phyluce_assembly_assemblo_trinity =====

```

**Attention:** This is not a toy tutorial - the data you downloaded contain roughly 100 MB for each of READ1 and READ2, which means it's going to take some time for these samples to assemble.

Once the assembly is finished, have a look at the directory structure:

```

uce-tutorial
+- assembly.conf
+- clean-fastq
+- illumiprocessor.conf
+- illumiprocessor.log
+- phyluce_assembly_assemblo_trinity.log
+- raw-fastq
+- trinity-assemblies
  +- alligator_mississippiensis_trinity
  | +- contigs.fasta -> Trinity.fasta
  | +- Trinity.fasta
  | +- trinity.log
  +- anolis_carolinensis_trinity
  | +- contigs.fasta -> Trinity.fasta
  | +- Trinity.fasta
  | +- trinity.log
  +- contigs
  | +- alligator_mississippiensis.contigs.fasta -> ../alligator_mississippiensis_
↳trinity/Trinity.fasta
  | +- anolis_carolinensis.contigs.fasta -> ../anolis_carolinensis_trinity/
↳Trinity.fasta
  | +- gallus_gallus.contigs.fasta -> ../gallus_gallus_trinity/Trinity.fasta
  | +- mus_musculus.contigs.fasta -> ../mus_musculus_trinity/Trinity.fasta
  +- gallus_gallus_trinity
  | +- contigs.fasta -> Trinity.fasta
  | +- Trinity.fasta
  | +- trinity.log
  +- mus_musculus_trinity

```

```
+-- contigs.fasta -> Trinity.fasta
+-- Trinity.fasta
+-- trinity.log
```

Your species-specific assembly files are in the *trinity-assemblies* directory nested within species-specific directories that correspond to the name you used in the *assembly.conf* file (to the left of the colon). Each name is appended with *\_trinity* because that's what *trinity* requires. There is a *symlink* within each species-specific folder so that you now the “contigs” you assembled are in *Trinity.fasta*. *trinity.log* holds the log output from *trinity*.

There is also a *contigs* directory within this folder. The *contigs* directory is the important one, because it contains *symlinks* to all of the species- specific contigs. This means that you can treat this single folder as if it contains all of your assembled contigs.

## Assembly QC

We can get a sense of how well the assembly worked by running the following from the top of our working directory:

```
# run this script against all directories of reads

for i in trinity-assemblies/contigs/*.fasta;
do
    phyluce_assembly_get_fasta_lengths --input $i --csv;
done
```

This should output something similar to the following. I've added the header as a comment:

```
# samples, contigs, total bp, mean length, 95 CI length, min length, max length, median_
↪ legnth, contigs >1kb
alligator_mississippiensis.contigs.fasta,10587,5820479,549.776046094,3.5939422934,224,
↪11285,413.0,1182
anolis_carolinensis.contigs.fasta,2458,1067208,434.177379984,5.72662897806,224,4359,
↪319.0,34
gallus_gallus.contigs.fasta,19905,8841661,444.192966591,2.06136172068,224,9883,306.0,
↪1530
mus_musculus.contigs.fasta,2162,1126231,520.920906568,7.75103292163,224,6542,358.0,186
```

### Question: Why are my numbers slightly different than your numbers?

The process of read assembly often differs by operating system and sometimes by OS version, and some of these differences are due to libraries that underlie many of the assembly programs. Expect to see differences. You should not expect for them to be huge.

**Attention:** If you see max-contig sizes around 16KB (for vertebrates), that is commonly the entire or almost-entire mtDNA genome. You do not tend to see entire mtDNA assemblies when the input DNA was extracted from a source having few mitochondria (e.g. blood).

There are many, many other assembly QC steps you can run other than simply looking at the stats of the assembled contigs. We will not go into those here.

## Finding UCE loci

Now that we've assembled our contigs from raw reads, it's time to find those contigs which are UCE loci and move aside those that are not. The directory structure before we do this should look like the following:

```
uce-tutorial
+- assembly.conf
+- clean-fastq
+- illumiprocessor.conf
+- illumiprocessor.log
+- phyluce_assembly_assemblo_trinity.log
+- raw-fastq
+- trinity-assemblies
```

Before we locate UCE loci, you need to get the probe set used for the enrichments:

```
wget https://raw.githubusercontent.com/faircloth-lab/uce-probe-sets/master/uce-5k-
↳probe-set/uce-5k-probes.fasta
```

Now, our directory structure looks like:

```
uce-tutorial
+- assembly.conf
+- clean-fastq
+- illumiprocessor.conf
+- illumiprocessor.log
+- phyluce_assembly_assemblo_trinity.log
+- raw-fastq
+- trinity-assemblies
+- uce-5k-probes.fasta
```

Now, run the *phyluce\_assembly\_match\_contigs\_to\_probes* program:

```
phyluce_assembly_match_contigs_to_probes \
  --contigs trinity-assemblies/contigs \
  --probes uce-5k-probes.fasta \
  --output uce-search-results
```

You should see output similar to the following (also stored in *phyluce\_assembly\_assemblo\_trinity.log*):

```
2015-04-12 12:44:58,951 - phyluce_assembly_match_contigs_to_probes - INFO - =====_
↳Starting phyluce_assembly_match_contigs_to_probes =====
2015-04-12 12:44:58,951 - phyluce_assembly_match_contigs_to_probes - INFO - Version:_
↳git 2a9c49d
2015-04-12 12:44:58,951 - phyluce_assembly_match_contigs_to_probes - INFO - Argument -
↳-contigs: /scratch/uce-tutorial/trinity-assemblies/contigs
2015-04-12 12:44:58,951 - phyluce_assembly_match_contigs_to_probes - INFO - Argument -
↳-dupefile: None
2015-04-12 12:44:58,952 - phyluce_assembly_match_contigs_to_probes - INFO - Argument -
↳-keep_duplicates: None
2015-04-12 12:44:58,952 - phyluce_assembly_match_contigs_to_probes - INFO - Argument -
↳-log_path: None
2015-04-12 12:44:58,952 - phyluce_assembly_match_contigs_to_probes - INFO - Argument -
↳-min_coverage: 80
2015-04-12 12:44:58,952 - phyluce_assembly_match_contigs_to_probes - INFO - Argument -
↳-min_identity: 80
2015-04-12 12:44:58,952 - phyluce_assembly_match_contigs_to_probes - INFO - Argument -
↳-output: /scratch/uce-tutorial/uce-search-results
```



```

2015-04-12 12:44:58,952 - phyluce_assembly_match_contigs_to_probes - INFO - Argument -
↳-probes: /scratch/uce-tutorial/uce-5k-probes.fasta
2015-04-12 12:44:58,952 - phyluce_assembly_match_contigs_to_probes - INFO - Argument -
↳-regex: ^(uce-\d+)(?:_p\d+.* )
2015-04-12 12:44:58,953 - phyluce_assembly_match_contigs_to_probes - INFO - Argument -
↳-verbosity: INFO
2015-04-12 12:44:59,111 - phyluce_assembly_match_contigs_to_probes - INFO - Creating
↳the UCE-match database
2015-04-12 12:44:59,405 - phyluce_assembly_match_contigs_to_probes - INFO -
↳Processing contig data
2015-04-12 12:44:59,405 - phyluce_assembly_match_contigs_to_probes - INFO - -----
↳-----
2015-04-12 12:45:12,735 - phyluce_assembly_match_contigs_to_probes - INFO - alligator_
↳mississippiensis: 4315 (40.76%) uniques of 10587 contigs, 0 dupe probe matches, 230
↳UCE loci removed for matching multiple contigs, 40 contigs removed for matching
↳multiple UCE loci
2015-04-12 12:45:15,919 - phyluce_assembly_match_contigs_to_probes - INFO - anolis_
↳carolinensis: 703 (28.60%) uniques of 2458 contigs, 0 dupe probe matches, 138 UCE
↳loci removed for matching multiple contigs, 2 contigs removed for matching multiple
↳UCE loci
2015-04-12 12:45:33,479 - phyluce_assembly_match_contigs_to_probes - INFO - gallus_
↳gallus: 3923 (19.71%) uniques of 19905 contigs, 0 dupe probe matches, 625 UCE loci
↳removed for matching multiple contigs, 47 contigs removed for matching multiple UCE
↳loci
2015-04-12 12:45:36,604 - phyluce_assembly_match_contigs_to_probes - INFO - mus_
↳musculus: 825 (38.16%) uniques of 2162 contigs, 0 dupe probe matches, 93 UCE loci
↳removed for matching multiple contigs, 1 contigs removed for matching multiple UCE
↳loci
2015-04-12 12:45:36,604 - phyluce_assembly_match_contigs_to_probes - INFO - -----
↳-----
2015-04-12 12:45:36,605 - phyluce_assembly_match_contigs_to_probes - INFO - The LASTZ
↳alignments are in /scratch/uce-tutorial/uce-search-results
2015-04-12 12:45:36,605 - phyluce_assembly_match_contigs_to_probes - INFO - The UCE
↳match database is in /scratch/uce-tutorial/uce-search-results/probe.matches.sqlite
2015-04-12 12:45:36,605 - phyluce_assembly_match_contigs_to_probes - INFO - =====
↳Completed phyluce_assembly_match_contigs_to_probes =====

```

The header info at the top tells us exactly what version of the code we are running and keeps track of our options. The important output is:

```

alligator_mississippiensis: 4315 (40.76%) uniques of 10587 contigs, 0 dupe probe
↳matches, 230 UCE loci removed for matching multiple contigs, 40 contigs removed for
↳matching multiple UCE loci
anolis_carolinensis: 703 (28.60%) uniques of 2458 contigs, 0 dupe probe matches, 138
↳UCE loci removed for matching multiple contigs, 2 contigs removed for matching
↳multiple UCE loci
gallus_gallus: 3923 (19.71%) uniques of 19905 contigs, 0 dupe probe matches, 625 UCE
↳loci removed for matching multiple contigs, 47 contigs removed for matching
↳multiple UCE loci
mus_musculus: 825 (38.16%) uniques of 2162 contigs, 0 dupe probe matches, 93 UCE loci
↳removed for matching multiple contigs, 1 contigs removed for matching multiple UCE
↳loci

```

Which we can break down to the following (for *alligator\_mississippiensis*):

```

alligator_mississippiensis:
    4315 (40.76%) uniques of 10587 contigs
    0 dupe probe matches

```

```
230 UCE loci removed for matching multiple contigs
40 contigs removed for matching multiple UCE loci
```

These are the capture data for the *alligator\_mississippiensis* sample. We targeted 5k UCE loci in this sample and recovered roughly 4484 of those loci. Before reaching that total of 4484 loci, we removed 45 and 44 loci from the data set because they looked like duplicates (probes supposedly targeting different loci hit the same contig or two supposedly different contigs hit probes designed for a single UCE locus).

---

**Question: Why is the count of UCE loci different by sample?**

For these example data, we enriched some samples (*alligator\_mississippiensis* and *gallus\_gallus*) for 5k UCE loci, while we enriched others (*anolis\_carolinensis* and *mus\_musculus*) for 2.5k UCE loci. Additionally, the 2.5k UCE enrichments did not work very well (operator error).

---

The directory structure now looks like the following (everything collapsed but the *uce-search-results* directory):

```
uce-tutorial
+- assembly.conf
+- clean-fastq
+- illumiprocessor.conf
+- illumiprocessor.log
+- phyluce_assembly_assemblo_trinity.log
+- phyluce_assembly_match_contigs_to_probes.log
+- raw-fastq
+- trinity-assemblies
+- uce-5k-probes.fasta
+- uce-search-results
  +- alligator_mississippiensis.contigs.lastz
  +- anolis_carolinensis.contigs.lastz
  +- gallus_gallus.contigs.lastz
  +- mus_musculus.contigs.lastz
  +- probe.matches.sqlite
```

The search we just ran created *lastz* search result files for each taxon, and stored summary results of these searches in the *probe.matches.sqlite* database (see *The probe.matches.sqlite database* for more information on this database and its structure).

## Extracting UCE loci

Now that we have located UCE loci, we need to determine which taxa we want in our analysis, create a list of those taxa, and then generated a list of which UCE loci we enriched in each taxon (the “data matrix configuration file”). We will then use this list to extract FASTA data for each taxon for each UCE locus.

First, we need to decide which taxa we want in our “taxon set”. So, we create a configuration file like so:

```
[all]
alligator_mississippiensis
anolis_carolinensis
gallus_gallus
mus_musculus
```

These names need to match the assembly names we used. Here, we have just put all 4 taxa in a list that we named *all*. However, we can adjust this list in many ways (see *Creating a data matrix configuration file*).

Save this file as *taxon-set.conf* at the top level of our *uce-tutorial* directory. The directory should look like this, now:

```
uce-tutorial
+- assembly.conf
+- clean-fastq
+- illumiprocessor.conf
+- illumiprocessor.log
+- phyluce_assembly_assemblo_trinity.log
+- phyluce_assembly_match_contigs_to_probes.log
+- raw-fastq
+- taxon-set.conf
+- trinity-assemblies
+- uce-5k-probes.fasta
+- uce-search-results
```

Now that we have this file created, we do the following to create the initial list of loci for each taxon:

```
# create an output directory for this taxon set - this just keeps
# things neat
cd uce-tutorial
mkdir -p taxon-sets/all

# create the data matrix configuration file
phyluce_assembly_get_match_counts \
  --locus-db uce-search-results/probe.matches.sqlite \
  --taxon-list-config taxon-set.conf \
  --taxon-group 'all' \
  --incomplete-matrix \
  --output taxon-sets/all/all-taxa-incomplete.conf
```

The output should look like the following:

```
2015-04-12 12:56:11,835 - phyluce_assembly_get_match_counts - INFO - =====_
↳Starting phyluce_assembly_get_match_counts =====
2015-04-12 12:56:11,835 - phyluce_assembly_get_match_counts - INFO - Version: git_
↳2a9c49d
2015-04-12 12:56:11,835 - phyluce_assembly_get_match_counts - INFO - Argument --
↳extend_locus_db: None
2015-04-12 12:56:11,835 - phyluce_assembly_get_match_counts - INFO - Argument --
↳incomplete_matrix: True
2015-04-12 12:56:11,835 - phyluce_assembly_get_match_counts - INFO - Argument --keep_
↳counts: False
2015-04-12 12:56:11,835 - phyluce_assembly_get_match_counts - INFO - Argument --locus_
↳db: /scratch/uce-tutorial/uce-search-results/probe.matches.sqlite
2015-04-12 12:56:11,836 - phyluce_assembly_get_match_counts - INFO - Argument --log_
↳path: None
2015-04-12 12:56:11,836 - phyluce_assembly_get_match_counts - INFO - Argument --
↳optimize: False
2015-04-12 12:56:11,836 - phyluce_assembly_get_match_counts - INFO - Argument --
↳output: /scratch/uce-tutorial/taxon-sets/all/all-taxa-incomplete.conf
2015-04-12 12:56:11,836 - phyluce_assembly_get_match_counts - INFO - Argument --
↳random: False
2015-04-12 12:56:11,836 - phyluce_assembly_get_match_counts - INFO - Argument --
↳sample_size: 10
2015-04-12 12:56:11,836 - phyluce_assembly_get_match_counts - INFO - Argument --
↳samples: 10
2015-04-12 12:56:11,836 - phyluce_assembly_get_match_counts - INFO - Argument --
↳silent: False
2015-04-12 12:56:11,836 - phyluce_assembly_get_match_counts - INFO - Argument --taxon_
↳group: all
```

```

2015-04-12 12:56:11,837 - phyluce_assembly_get_match_counts - INFO - Argument --taxon_
↳list_config: /scratch/uce-tutorial/taxon-set.conf
2015-04-12 12:56:11,837 - phyluce_assembly_get_match_counts - INFO - Argument --
↳verbosity: INFO
2015-04-12 12:56:11,837 - phyluce_assembly_get_match_counts - INFO - There are 4 taxa_
↳in the taxon-group '[all]' in the config file taxon-set.conf
2015-04-12 12:56:11,838 - phyluce_assembly_get_match_counts - INFO - Getting UCE_
↳names from database
2015-04-12 12:56:11,844 - phyluce_assembly_get_match_counts - INFO - There are 5041_
↳total UCE loci in the database
2015-04-12 12:56:11,946 - phyluce_assembly_get_match_counts - INFO - Getting UCE_
↳matches by organism to generate a INCOMPLETE matrix
2015-04-12 12:56:11,948 - phyluce_assembly_get_match_counts - INFO - There are 4710_
↳UCE loci in an INCOMPLETE matrix
2015-04-12 12:56:11,949 - phyluce_assembly_get_match_counts - INFO - Writing the taxa_
↳and loci in the data matrix to /scratch/uce-tutorial/taxon-sets/all/all-taxa-
↳incomplete.conf
2015-04-12 12:56:11,952 - phyluce_assembly_get_match_counts - INFO - =====_
↳Completed phyluce_assembly_get_match_counts =====

```

And, our directory structure should now look like this (collapsing all but *taxon-sets*):

```

uce-tutorial
+- assembly.conf
+- clean-fastq
+- illumiprocessor.conf
+- illumiprocessor.log
+- phyluce_assembly_assemblo_trinity.log
+- phyluce_assembly_get_match_counts.log
+- phyluce_assembly_match_contigs_to_probes.log
+- taxon-set.conf
+- taxon-sets
| +- all
|   +- all-taxa-incomplete.conf
+- trinity-assemblies
+- uce-5k-probes.fasta
+- uce-search-results

```

Now, we need to extract FASTA data that correspond to the loci in *all-taxa-incomplete.conf*:

```

# change to the taxon-sets/all directory
cd taxon-sets/all

# make a log directory to hold our log files - this keeps things neat
mkdir log

# get FASTA data for taxa in our taxon set
phyluce_assembly_get_fastas_from_match_counts \
  --contigs ../../trinity-assemblies/contigs \
  --locus-db ../../uce-search-results/probe.matches.sqlite \
  --match-count-output all-taxa-incomplete.conf \
  --output all-taxa-incomplete.fasta \
  --incomplete-matrix all-taxa-incomplete.incomplete \
  --log-path log

```

The output should look something like the following:

```

2015-04-12 12:58:40,646 - phyluce_assembly_get_fastas_from_match_counts - INFO -
↳==== Starting phyluce_assembly_get_fastas_from_match_counts ====
2015-04-12 12:58:40,646 - phyluce_assembly_get_fastas_from_match_counts - INFO -
↳Version: git a6a957a
2015-04-12 12:58:40,646 - phyluce_assembly_get_fastas_from_match_counts - INFO -
↳Argument --contigs: /scratch/uce-tutorial/trinity-assemblies/contigs
2015-04-12 12:58:40,647 - phyluce_assembly_get_fastas_from_match_counts - INFO -
↳Argument --extend_locus_contigs: None
2015-04-12 12:58:40,647 - phyluce_assembly_get_fastas_from_match_counts - INFO -
↳Argument --extend_locus_db: None
2015-04-12 12:58:40,647 - phyluce_assembly_get_fastas_from_match_counts - INFO -
↳Argument --incomplete_matrix: /scratch/uce-tutorial/taxon-sets/all/all-taxa-
↳incomplete.incomplete
2015-04-12 12:58:40,647 - phyluce_assembly_get_fastas_from_match_counts - INFO -
↳Argument --locus_db: /scratch/uce-tutorial/uce-search-results/probe.matches.sqlite
2015-04-12 12:58:40,647 - phyluce_assembly_get_fastas_from_match_counts - INFO -
↳Argument --log_path: /scratch/uce-tutorial/taxon-sets/all/log
2015-04-12 12:58:40,647 - phyluce_assembly_get_fastas_from_match_counts - INFO -
↳Argument --match_count_output: /scratch/uce-tutorial/taxon-sets/all/all-taxa-
↳incomplete.conf
2015-04-12 12:58:40,647 - phyluce_assembly_get_fastas_from_match_counts - INFO -
↳Argument --output: /scratch/uce-tutorial/taxon-sets/all/all-taxa-incomplete.fasta
2015-04-12 12:58:40,647 - phyluce_assembly_get_fastas_from_match_counts - INFO -
↳Argument --verbosity: INFO
2015-04-12 12:58:40,702 - phyluce_assembly_get_fastas_from_match_counts - INFO -
↳There are 4 taxa in the match-count-config file named all-taxa-incomplete.conf
2015-04-12 12:58:40,716 - phyluce_assembly_get_fastas_from_match_counts - INFO -
↳There are 4710 UCE loci in an INCOMPLETE matrix
2015-04-12 12:58:40,717 - phyluce_assembly_get_fastas_from_match_counts - INFO -
↳-----Getting UCE loci for alligator_mississippiensis-----
2015-04-12 12:58:40,773 - phyluce_assembly_get_fastas_from_match_counts - INFO -
↳There are 4315 UCE loci for alligator_mississippiensis
2015-04-12 12:58:40,774 - phyluce_assembly_get_fastas_from_match_counts - INFO -
↳Parsing and renaming contigs for alligator_mississippiensis
2015-04-12 12:58:51,084 - phyluce_assembly_get_fastas_from_match_counts - INFO -
↳Writing missing locus information to /scratch/uce-tutorial/taxon-sets/all/all-taxa-
↳incomplete.incomplete
2015-04-12 12:58:51,086 - phyluce_assembly_get_fastas_from_match_counts - INFO -
↳-----Getting UCE loci for anolis_carolinensis-----
2015-04-12 12:58:51,120 - phyluce_assembly_get_fastas_from_match_counts - INFO -
↳There are 703 UCE loci for anolis_carolinensis
2015-04-12 12:58:51,121 - phyluce_assembly_get_fastas_from_match_counts - INFO -
↳Parsing and renaming contigs for anolis_carolinensis
2015-04-12 12:58:51,693 - phyluce_assembly_get_fastas_from_match_counts - INFO -
↳Writing missing locus information to /scratch/uce-tutorial/taxon-sets/all/all-taxa-
↳incomplete.incomplete
2015-04-12 12:58:51,701 - phyluce_assembly_get_fastas_from_match_counts - INFO -
↳-----Getting UCE loci for gallus_gallus-----
2015-04-12 12:58:51,753 - phyluce_assembly_get_fastas_from_match_counts - INFO -
↳There are 3923 UCE loci for gallus_gallus
2015-04-12 12:58:51,753 - phyluce_assembly_get_fastas_from_match_counts - INFO -
↳Parsing and renaming contigs for gallus_gallus
2015-04-12 12:59:10,959 - phyluce_assembly_get_fastas_from_match_counts - INFO -
↳Writing missing locus information to /scratch/uce-tutorial/taxon-sets/all/all-taxa-
↳incomplete.incomplete
2015-04-12 12:59:10,962 - phyluce_assembly_get_fastas_from_match_counts - INFO -
↳-----Getting UCE loci for mus_musculus-----
2015-04-12 12:59:10,997 - phyluce_assembly_get_fastas_from_match_counts - INFO -
↳There are 825 UCE loci for mus_musculus

```

```
2015-04-12 12:59:10,997 - phyluce_assembly_get_fastas_from_match_counts - INFO - ↵
↳ Parsing and renaming contigs for mus_musculus
2015-04-12 12:59:11,610 - phyluce_assembly_get_fastas_from_match_counts - INFO - ↵
↳ Writing missing locus information to /scratch/uce-tutorial/taxon-sets/all/all-taxa-
↳ incomplete.incomplete
2015-04-12 12:59:11,618 - phyluce_assembly_get_fastas_from_match_counts - INFO - ===== ↵
↳ Completed phyluce_assembly_get_fastas_from_match_counts =====
```

And, our directory structure should now look like this (collapsing all but *taxon-sets*):

```
uce-tutorial
+- assembly.conf
+- clean-fastq
+- illumiprocessor.conf
+- illumiprocessor.log
+- phyluce_assembly_assemblo_trinity.log
+- phyluce_assembly_get_match_counts.log
+- phyluce_assembly_match_contigs_to_probes.log
+- taxon-set.conf
+- taxon-sets
| +- all
|   +- all-taxa-incomplete.conf
|   +- all-taxa-incomplete.fasta
|   +- all-taxa-incomplete.incomplete
|   +- log
|   +- phyluce_assembly_get_fastas_from_match_counts.log
+- phyluce_assembly_get_fastas_from_match_counts.log
+- trinity-assemblies
+- uce-5k-probes.fasta
+- uce-search-results
```

The extracted FASTA data are in a monolithic FASTA file (all data for all organisms) named *all-taxa-incomplete.fasta*.

## Exploding the monolithic FASTA file

Lots of times we want to know individual statistics on UCE assemblies for a given taxon. We can do that by exploding the monolithic fasta file into a file of UCE loci that we have enriched by taxon, then running stats on those exploded files. To do that, run the following:

```
# explode the monolithic FASTA by taxon (you can also do by locus)
phyluce_assembly_explode_get_fastas_file \
  --input all-taxa-incomplete.fasta \
  --output-dir exploded-fastas \
  --by-taxon

# get summary stats on the FASTAS
for i in exploded-fastas/*.fasta;
do
  phyluce_assembly_get_fasta_lengths --input $i --csv;
done

# samples, contigs, total bp, mean length, 95 CI length, min length, max length, median_
↳ length, contigs >1kb
alligator-mississippiensis.unaligned.fasta, 4315, 3465679, 803.170104287, 3.80363492428,
↳ 224, 1794, 823.0, 980
anolis-carolinensis.unaligned.fasta, 703, 400214, 569.294452347, 9.16433421241, 224, 1061,
↳ 546.0, 7
```

```
gallus-gallus.unaligned.fasta,3923,3273674,834.482283966,4.26048496461,231,1864,852.0,
↪1149
mus-musculus.unaligned.fasta,825,594352,720.426666667,9.85933217965,225,1178,823.0,139
```

## Aligning UCE loci

You have lots of options when aligning UCE loci. You can align the loci and use those alignments with no trimming, you can edge-trim the alignments following some algorithm, and you can end+internally trim alignments following some algorithm. It's hard to say what is best in all situations. When taxa are “closely” related (< 30-50 MYA, perhaps), I think that edge-trimming alignments is reasonable. When the taxa you are interested in span a wider range of divergence times (> 50 MYA), you may want to think about internal trimming.

How you accomplish you edge- or internal trimming is also a decision you need to make. In *phyluce*, we implement our edge-trimming algorithm by running the alignment program “as-is” (i.e., without the *-no-trim*) option. We do internal-trimming by turning off trimming using *-no-trim*, then passing the resulting alignments (in FASTA format) to a parallel wrapper around *Gblocks*.

You also have a choice of aligner - *mafft* or *muscle* (or you can externally align UCE loci using a tool like SATé, as well).

Generally, I would use *mafft*.

## Edge trimming

Edge trimming your alignments is a relatively simple matter. You can run edge trimming, as follows:

```
# make sure we are in the correct directory
cd uce-tutorial/taxon-sets/all

# align the data
phyluce_align_seqcap_align \
  --fasta all-taxa-incomplete.fasta \
  --output mafft-nexus-edge-trimmed \
  --taxa 4 \
  --aligner mafft \
  --cores 12 \
  --incomplete-matrix \
  --log-path log
```

**Warning:** Note that I am using 12 physical CPU cores here. You need to use the number of physical cores available on *your* machine.

The output should look like this:

```
2015-04-12 13:01:31,919 - phyluce_align_seqcap_align - INFO - ===== Starting_
↪phyluce_align_seqcap_align =====
2015-04-12 13:01:31,920 - phyluce_align_seqcap_align - INFO - Version: git a6a957a
2015-04-12 13:01:31,920 - phyluce_align_seqcap_align - INFO - Argument --aligner:_
↪mafft
2015-04-12 13:01:31,920 - phyluce_align_seqcap_align - INFO - Argument --ambiguous:_
↪False
2015-04-12 13:01:31,920 - phyluce_align_seqcap_align - INFO - Argument --cores: 12
2015-04-12 13:01:31,920 - phyluce_align_seqcap_align - INFO - Argument --fasta: /
↪scratch/uce-tutorial/taxon-sets/all/all-taxa-incomplete.fasta
```

```

2015-04-12 13:01:31,920 - phyluce_align_seqcap_align - INFO - Argument --log_path: /
↳scratch/uce-tutorial/taxon-sets/all/log
2015-04-12 13:01:31,920 - phyluce_align_seqcap_align - INFO - Argument --max_
↳divergence: 0.2
2015-04-12 13:01:31,921 - phyluce_align_seqcap_align - INFO - Argument --min_length:
↳100
2015-04-12 13:01:31,921 - phyluce_align_seqcap_align - INFO - Argument --no_trim:
↳False
2015-04-12 13:01:31,921 - phyluce_align_seqcap_align - INFO - Argument --notstrict:
↳True
2015-04-12 13:01:31,921 - phyluce_align_seqcap_align - INFO - Argument --output: /
↳scratch/uce-tutorial/taxon-sets/all/mafft-nexus-edge-trimmed
2015-04-12 13:01:31,921 - phyluce_align_seqcap_align - INFO - Argument --output_
↳format: nexus
2015-04-12 13:01:31,921 - phyluce_align_seqcap_align - INFO - Argument --proportion:
↳0.65
2015-04-12 13:01:31,921 - phyluce_align_seqcap_align - INFO - Argument --taxa: 4
2015-04-12 13:01:31,921 - phyluce_align_seqcap_align - INFO - Argument --threshold: 0.
↳65
2015-04-12 13:01:31,922 - phyluce_align_seqcap_align - INFO - Argument --verbosity:
↳INFO
2015-04-12 13:01:31,922 - phyluce_align_seqcap_align - INFO - Argument --window: 20
2015-04-12 13:01:31,922 - phyluce_align_seqcap_align - INFO - Building the locus_
↳dictionary
2015-04-12 13:01:31,922 - phyluce_align_seqcap_align - INFO - Removing ALL sequences_
↳with ambiguous bases...
2015-04-12 13:01:34,051 - phyluce_align_seqcap_align - WARNING - DROPPED locus uce-
↳6169. Too few taxa (N < 3).
[many more loci dropped here]
2015-04-12 13:01:34,596 - phyluce_align_seqcap_align - INFO - Aligning with MAFFT
2015-04-12 13:01:34,598 - phyluce_align_seqcap_align - INFO - Alignment begins. 'X'_
↳indicates dropped alignments (these are reported after alignment)
.....[continued]
2015-04-12 13:02:05,847 - phyluce_align_seqcap_align - INFO - Alignment ends
2015-04-12 13:02:05,848 - phyluce_align_seqcap_align - INFO - Writing output files
2015-04-12 13:02:06,567 - phyluce_align_seqcap_align - INFO - =====
↳Completed phyluce_align_seqcap_align =====

```

The . values that you see represent loci that were aligned and successfully trimmed. Any X values that you see represent loci that were removed because trimming reduced their length to effectively nothing.

**Attention:** The number of potential alignments dropped here is abnormally large because our **sample size is so small (n=4)**.

The current directory structure should look like (I've collapsed a number of branches in the tree):

```

uce-tutorial
+- assembly.conf
...
+- taxon-sets
|   +- all
|     +- all-taxa-incomplete.conf
|     +- all-taxa-incomplete.fasta
|     +- all-taxa-incomplete.incomplete
|     +- exploded-fastas
|     +- log

```



```
|         +- mafft-nexus-edge-trimmed
|         +- uce-1008.nexus
|         +- uce-1014.nexus
|         +- uce-1039.nexus
|         ...
|         +- uce-991.nexus
...
+- uce-search-results
```

We can output summary stats for these alignments by running the following program:

```
phyluce_align_get_align_summary_data \
  --alignments mafft-nexus-edge-trimmed \
  --cores 12 \
  --log-path log
```

**Warning:** Note that I am using 12 physical CPU cores here. You need to use the number of physical cores available on *your* machine.

The output from the program should look like:

```
2015-04-12 13:40:56,410 - phyluce_align_get_align_summary_data - INFO - =====_
↳Starting phyluce_align_get_align_summary_data =====
2015-04-12 13:40:56,410 - phyluce_align_get_align_summary_data - INFO - Version: git_
↳a6a957a
2015-04-12 13:40:56,410 - phyluce_align_get_align_summary_data - INFO - Argument --
↳alignments: /scratch/uce-tutorial/taxon-sets/all/mafft-nexus-edge-trimmed
2015-04-12 13:40:56,410 - phyluce_align_get_align_summary_data - INFO - Argument --
↳cores: 12
2015-04-12 13:40:56,410 - phyluce_align_get_align_summary_data - INFO - Argument --
↳input_format: nexus
2015-04-12 13:40:56,411 - phyluce_align_get_align_summary_data - INFO - Argument --
↳log_path: /scratch/uce-tutorial/taxon-sets/all/log
2015-04-12 13:40:56,411 - phyluce_align_get_align_summary_data - INFO - Argument --
↳show_taxon_counts: False
2015-04-12 13:40:56,411 - phyluce_align_get_align_summary_data - INFO - Argument --
↳verbosity: INFO
2015-04-12 13:40:56,411 - phyluce_align_get_align_summary_data - INFO - Getting_
↳alignment files
2015-04-12 13:40:56,416 - phyluce_align_get_align_summary_data - INFO - Computing_
↳summary statistics using 12 cores
2015-04-12 13:40:56,716 - phyluce_align_get_align_summary_data - INFO - -----
↳----- Alignment summary -----
2015-04-12 13:40:56,717 - phyluce_align_get_align_summary_data - INFO - [Alignments]_
↳loci:          913
2015-04-12 13:40:56,717 - phyluce_align_get_align_summary_data - INFO - [Alignments]_
↳length:       559,055
2015-04-12 13:40:56,717 - phyluce_align_get_align_summary_data - INFO - [Alignments]_
↳mean:         612.33
2015-04-12 13:40:56,717 - phyluce_align_get_align_summary_data - INFO - [Alignments]_
↳95% CI:       13.92
2015-04-12 13:40:56,717 - phyluce_align_get_align_summary_data - INFO - [Alignments]_
↳min:          157
2015-04-12 13:40:56,717 - phyluce_align_get_align_summary_data - INFO - [Alignments]_
↳max:          1,266
2015-04-12 13:40:56,719 - phyluce_align_get_align_summary_data - INFO - -----
↳----- Taxon summary -----
```

```

2015-04-12 13:40:56,719 - phyluce_align_get_align_summary_data - INFO - [Taxa] mean: 3.38
↳ 3.38
2015-04-12 13:40:56,719 - phyluce_align_get_align_summary_data - INFO - [Taxa] 95% CI: 0.03
↳ CI: 0.03
2015-04-12 13:40:56,719 - phyluce_align_get_align_summary_data - INFO - [Taxa] min: 3
↳ 3
2015-04-12 13:40:56,719 - phyluce_align_get_align_summary_data - INFO - [Taxa] max: 4
↳ 4
2015-04-12 13:40:56,720 - phyluce_align_get_align_summary_data - INFO - -----
↳ --- Missing data from trim summary -----
2015-04-12 13:40:56,720 - phyluce_align_get_align_summary_data - INFO - [Missing] mean: 9.61
↳ mean: 9.61
2015-04-12 13:40:56,720 - phyluce_align_get_align_summary_data - INFO - [Missing] 95% CI: 0.42
↳ CI: 0.42
2015-04-12 13:40:56,720 - phyluce_align_get_align_summary_data - INFO - [Missing] min: 0.00
↳ min: 0.00
2015-04-12 13:40:56,720 - phyluce_align_get_align_summary_data - INFO - [Missing] max: 27.98
↳ max: 27.98
2015-04-12 13:40:56,732 - phyluce_align_get_align_summary_data - INFO - -----
↳ ----- Character count summary -----
2015-04-12 13:40:56,732 - phyluce_align_get_align_summary_data - INFO - [All characters] 1,852,854
↳ characters] 1,852,854
2015-04-12 13:40:56,732 - phyluce_align_get_align_summary_data - INFO - [Nucleotides] 1,600,577
↳ 1,600,577
2015-04-12 13:40:56,733 - phyluce_align_get_align_summary_data - INFO - -----
↳ -- Data matrix completeness summary -----
2015-04-12 13:40:56,733 - phyluce_align_get_align_summary_data - INFO - [Matrix 50%] 913 alignments
↳ 913 alignments
2015-04-12 13:40:56,733 - phyluce_align_get_align_summary_data - INFO - [Matrix 55%] 913 alignments
↳ 913 alignments
2015-04-12 13:40:56,733 - phyluce_align_get_align_summary_data - INFO - [Matrix 60%] 913 alignments
↳ 913 alignments
2015-04-12 13:40:56,733 - phyluce_align_get_align_summary_data - INFO - [Matrix 65%] 913 alignments
↳ 913 alignments
2015-04-12 13:40:56,734 - phyluce_align_get_align_summary_data - INFO - [Matrix 70%] 913 alignments
↳ 913 alignments
2015-04-12 13:40:56,734 - phyluce_align_get_align_summary_data - INFO - [Matrix 75%] 913 alignments
↳ 913 alignments
2015-04-12 13:40:56,734 - phyluce_align_get_align_summary_data - INFO - [Matrix 80%] 346 alignments
↳ 346 alignments
2015-04-12 13:40:56,734 - phyluce_align_get_align_summary_data - INFO - [Matrix 85%] 346 alignments
↳ 346 alignments
2015-04-12 13:40:56,734 - phyluce_align_get_align_summary_data - INFO - [Matrix 90%] 346 alignments
↳ 346 alignments
2015-04-12 13:40:56,734 - phyluce_align_get_align_summary_data - INFO - [Matrix 95%] 346 alignments
↳ 346 alignments
2015-04-12 13:40:56,735 - phyluce_align_get_align_summary_data - INFO - -----
↳ ----- Character counts -----
2015-04-12 13:40:56,735 - phyluce_align_get_align_summary_data - INFO - [Characters] '-' is present 72,168 times
↳ '-' is present 72,168 times
2015-04-12 13:40:56,735 - phyluce_align_get_align_summary_data - INFO - [Characters] '?' is present 180,109 times
↳ '?' is present 180,109 times
2015-04-12 13:40:56,735 - phyluce_align_get_align_summary_data - INFO - [Characters] 'A' is present 494,411 times
↳ 'A' is present 494,411 times
2015-04-12 13:40:56,735 - phyluce_align_get_align_summary_data - INFO - [Characters] 'C' is present 308,796 times
↳ 'C' is present 308,796 times
2015-04-12 13:40:56,735 - phyluce_align_get_align_summary_data - INFO - [Characters] 'G' is present 295,171 times
↳ 'G' is present 295,171 times

```

```
2015-04-12 13:40:56,735 - phyluce_align_get_align_summary_data - INFO - [Characters]
↳ 'T' is present 502,199 times
2015-04-12 13:40:56,736 - phyluce_align_get_align_summary_data - INFO - =====_
↳ Completed phyluce_align_get_align_summary_data =====
```

**Attention:** Note that there are only 2 sets of counts in the *Data matrix completeness* section because (1) we dropped all loci having fewer than 3 taxa and (2) that only leaves two remaining options.

The most important data here are the number of loci we have and the number of loci in data matrices of different completeness. The locus length stats are also reasonably important, but they can also be misleading because edge-trimming does not remove internal gaps that often inflate the length of alignments.

### Internal trimming

Now, let's do the same thing, but run internal trimming on the resulting alignments. We will do that by turning off trimming `-no-trim` and outputting FASTA formatted alignments with `-output-format fasta`.

```
# make sure we are in the correct directory
cd uce-tutorial/taxon-sets/all

# align the data - turn off trimming and output FASTA
phyluce_align_seqcap_align \
  --fasta all-taxa-incomplete.fasta \
  --output mafft-nexus-internal-trimmed \
  --taxa 4 \
  --aligner mafft \
  --cores 12 \
  --incomplete-matrix \
  --output-format fasta \
  --no-trim \
  --log-path log
```

**Attention:** The number of UCE loci dropped here is abnormally large because our **sample size is so small (n=4)**.

**Warning:** Note that I am using 12 physical CPU cores here. You need to use the number of physical cores available on *your* machine.

The output from the program should be the roughly the same as what we saw before. The current directory structure should look like (I've collapsed a number of branches in the tree):

```
uce-tutorial
+- assembly.conf
...
+- taxon-sets
|   +- all
|     +- all-taxa-incomplete.conf
|     +- all-taxa-incomplete.fasta
|     +- all-taxa-incomplete.incomplete
|     +- exploded-fastas
|     +- log
```

```
|         +- mafft-nexus-edge-trimmed
|         +- mafft-nexus-internal-trimmed
|             +- uce-1008.nexus
|             +- uce-1014.nexus
|             +- uce-1039.nexus
|             ...
|             +- uce-991.nexus
...
+- uce-search-results
```

Now, we are going to trim these loci using **Gblocks**:

```
# run gblocks trimming on the alignments
phyluce_align_get_gblocks_trimmed_alignments_from_untrimmed \
    --alignments mafft-nexus-internal-trimmed \
    --output mafft-nexus-internal-trimmed-gblocks \
    --cores 12 \
    --log log
```

The output should look like this:

```
2015-04-12 13:51:52,450 - phyluce_align_get_gblocks_trimmed_alignments_from_untrimmed_
↳ INFO - Starting phyluce_align_get_gblocks_trimmed_alignments_from_untrimmed
2015-04-12 13:51:52,450 - phyluce_align_get_gblocks_trimmed_alignments_from_untrimmed_
↳ INFO - Version: git a6a957a
2015-04-12 13:51:52,451 - phyluce_align_get_gblocks_trimmed_alignments_from_untrimmed_
↳ INFO - Argument --alignments: /scratch/uce-tutorial/taxon-sets/all/mafft-nexus-
↳ internal-trimmed
2015-04-12 13:51:52,451 - phyluce_align_get_gblocks_trimmed_alignments_from_untrimmed_
↳ INFO - Argument --b1: 0.5
2015-04-12 13:51:52,451 - phyluce_align_get_gblocks_trimmed_alignments_from_untrimmed_
↳ INFO - Argument --b2: 0.85
2015-04-12 13:51:52,451 - phyluce_align_get_gblocks_trimmed_alignments_from_untrimmed_
↳ INFO - Argument --b3: 8
2015-04-12 13:51:52,451 - phyluce_align_get_gblocks_trimmed_alignments_from_untrimmed_
↳ INFO - Argument --b4: 10
2015-04-12 13:51:52,451 - phyluce_align_get_gblocks_trimmed_alignments_from_untrimmed_
↳ INFO - Argument --cores: 12
2015-04-12 13:51:52,451 - phyluce_align_get_gblocks_trimmed_alignments_from_untrimmed_
↳ INFO - Argument --input_format: fasta
2015-04-12 13:51:52,451 - phyluce_align_get_gblocks_trimmed_alignments_from_untrimmed_
↳ INFO - Argument --log_path: /scratch/uce-tutorial/taxon-sets/all/log
2015-04-12 13:51:52,452 - phyluce_align_get_gblocks_trimmed_alignments_from_untrimmed_
↳ INFO - Argument --output: /scratch/uce-tutorial/taxon-sets/all/mafft-nexus-
↳ internal-trimmed-gblocks
2015-04-12 13:51:52,452 - phyluce_align_get_gblocks_trimmed_alignments_from_untrimmed_
↳ INFO - Argument --output_format: nexus
2015-04-12 13:51:52,452 - phyluce_align_get_gblocks_trimmed_alignments_from_untrimmed_
↳ INFO - Argument --verbosity: INFO
2015-04-12 13:51:52,452 - phyluce_align_get_gblocks_trimmed_alignments_from_untrimmed_
↳ INFO - Starting phyluce_align_get_gblocks_trimmed_alignments_from_untrimmed
2015-04-12 13:51:52,452 - phyluce_align_get_gblocks_trimmed_alignments_from_untrimmed_
↳ INFO - Getting aligned sequences for trimming
2015-04-12 13:51:52,460 - phyluce_align_get_gblocks_trimmed_alignments_from_untrimmed_
↳ INFO - Alignment trimming begins.
.....[continued]
2015-04-12 13:51:53,063 - phyluce_align_get_gblocks_trimmed_alignments_from_untrimmed_
↳ INFO - Alignment trimming ends
```

```
2015-04-12 13:51:53,063 - phyluce_align_get_gblocks_trimmed_alignments_from_untrimmed_
↳ INFO - Writing output files
2015-04-12 13:51:53,725 - phyluce_align_get_gblocks_trimmed_alignments_from_untrimmed_
↳ INFO - Completed phyluce_align_get_gblocks_trimmed_alignments_from_untrimmed
```

The . values that you see represent loci that were aligned and successfully trimmed. Andy X values that you see represent loci that were aligned and trimmed so much that there was nothing left.

The current directory structure should look like (I've collapsed a number of branches in the tree):

```
uce-tutorial
+- assembly.conf
...
+- taxon-sets
|   +- all
|     +- all-taxa-incomplete.conf
|     +- all-taxa-incomplete.fasta
|     +- all-taxa-incomplete.incomplete
|     +- exploded-fastas
|     +- log
|     +- mafft-nexus-edge-trimmed
|     +- mafft-nexus-internal-trimmed
|     +- mafft-nexus-internal-trimmed-gblocks
|       +- uce-1008.nexus
|       +- uce-1014.nexus
|       +- uce-1039.nexus
|       ...
|       +- uce-991.nexus
...
+- uce-search-results
```

We can output summary stats for these alignments by running the following program:

```
phyluce_align_get_align_summary_data \
  --alignments mafft-nexus-internal-trimmed-gblocks \
  --cores 12 \
  --log-path log
```

**Warning:** Note that I am using 12 physical CPU cores here. You need to use the number of physical cores available on *your* machine.

The output from the program should look like:

```
2015-04-12 13:54:48,251 - phyluce_align_get_align_summary_data - INFO - =====_
↳ Starting phyluce_align_get_align_summary_data =====
2015-04-12 13:54:48,251 - phyluce_align_get_align_summary_data - INFO - Version: git_
↳ a6a957a
2015-04-12 13:54:48,251 - phyluce_align_get_align_summary_data - INFO - Argument --
↳ alignments: /scratch/uce-tutorial/taxon-sets/all/mafft-nexus-internal-trimmed-
↳ gblocks
2015-04-12 13:54:48,251 - phyluce_align_get_align_summary_data - INFO - Argument --
↳ cores: 12
2015-04-12 13:54:48,251 - phyluce_align_get_align_summary_data - INFO - Argument --
↳ input_format: nexus
2015-04-12 13:54:48,251 - phyluce_align_get_align_summary_data - INFO - Argument --
↳ log_path: /scratch/uce-tutorial/taxon-sets/all/log
```

```

2015-04-12 13:54:48,251 - phyluce_align_get_align_summary_data - INFO - Argument --
↳show_taxon_counts: False
2015-04-12 13:54:48,252 - phyluce_align_get_align_summary_data - INFO - Argument --
↳verbosity: INFO
2015-04-12 13:54:48,252 - phyluce_align_get_align_summary_data - INFO - Getting_
↳alignment files
2015-04-12 13:54:48,257 - phyluce_align_get_align_summary_data - INFO - Computing_
↳summary statistics using 12 cores
2015-04-12 13:54:48,523 - phyluce_align_get_align_summary_data - INFO - -----
↳----- Alignment summary -----
2015-04-12 13:54:48,523 - phyluce_align_get_align_summary_data - INFO - [Alignments]_
↳loci:          913
2015-04-12 13:54:48,523 - phyluce_align_get_align_summary_data - INFO - [Alignments]_
↳length:       499,929
2015-04-12 13:54:48,523 - phyluce_align_get_align_summary_data - INFO - [Alignments]_
↳mean:         547.57
2015-04-12 13:54:48,523 - phyluce_align_get_align_summary_data - INFO - [Alignments]_
↳95% CI:       14.10
2015-04-12 13:54:48,523 - phyluce_align_get_align_summary_data - INFO - [Alignments]_
↳min:          101
2015-04-12 13:54:48,524 - phyluce_align_get_align_summary_data - INFO - [Alignments]_
↳max:           1,180
2015-04-12 13:54:48,525 - phyluce_align_get_align_summary_data - INFO - -----
↳----- Taxon summary -----
2015-04-12 13:54:48,525 - phyluce_align_get_align_summary_data - INFO - [Taxa] mean:  _
↳          3.38
2015-04-12 13:54:48,525 - phyluce_align_get_align_summary_data - INFO - [Taxa] 95%_
↳CI:         0.03
2015-04-12 13:54:48,525 - phyluce_align_get_align_summary_data - INFO - [Taxa] min:  _
↳          3
2015-04-12 13:54:48,526 - phyluce_align_get_align_summary_data - INFO - [Taxa] max:  _
↳          4
2015-04-12 13:54:48,526 - phyluce_align_get_align_summary_data - INFO - -----
↳--- Missing data from trim summary -----
2015-04-12 13:54:48,526 - phyluce_align_get_align_summary_data - INFO - [Missing]_
↳mean:        0.00
2015-04-12 13:54:48,527 - phyluce_align_get_align_summary_data - INFO - [Missing] 95%_
↳CI:          0.00
2015-04-12 13:54:48,527 - phyluce_align_get_align_summary_data - INFO - [Missing]_
↳min:         0.00
2015-04-12 13:54:48,527 - phyluce_align_get_align_summary_data - INFO - [Missing]_
↳max:         0.00
2015-04-12 13:54:48,537 - phyluce_align_get_align_summary_data - INFO - -----
↳----- Character count summary -----
2015-04-12 13:54:48,538 - phyluce_align_get_align_summary_data - INFO - [All_
↳characters]   1,717,433
2015-04-12 13:54:48,538 - phyluce_align_get_align_summary_data - INFO - [Nucleotides]_
↳          1,621,491
2015-04-12 13:54:48,538 - phyluce_align_get_align_summary_data - INFO - -----
↳--- Data matrix completeness summary -----
2015-04-12 13:54:48,538 - phyluce_align_get_align_summary_data - INFO - [Matrix 50%] _
↳          913 alignments
2015-04-12 13:54:48,538 - phyluce_align_get_align_summary_data - INFO - [Matrix 55%] _
↳          913 alignments
2015-04-12 13:54:48,538 - phyluce_align_get_align_summary_data - INFO - [Matrix 60%] _
↳          913 alignments
2015-04-12 13:54:48,539 - phyluce_align_get_align_summary_data - INFO - [Matrix 65%] _
↳          913 alignments

```

```

2015-04-12 13:54:48,539 - phyluce_align_get_align_summary_data - INFO - [Matrix 70%] ↵
↪ 913 alignments
2015-04-12 13:54:48,539 - phyluce_align_get_align_summary_data - INFO - [Matrix 75%] ↵
↪ 913 alignments
2015-04-12 13:54:48,539 - phyluce_align_get_align_summary_data - INFO - [Matrix 80%] ↵
↪ 346 alignments
2015-04-12 13:54:48,539 - phyluce_align_get_align_summary_data - INFO - [Matrix 85%] ↵
↪ 346 alignments
2015-04-12 13:54:48,539 - phyluce_align_get_align_summary_data - INFO - [Matrix 90%] ↵
↪ 346 alignments
2015-04-12 13:54:48,539 - phyluce_align_get_align_summary_data - INFO - [Matrix 95%] ↵
↪ 346 alignments
2015-04-12 13:54:48,540 - phyluce_align_get_align_summary_data - INFO - -----
↪ ----- Character counts -----
2015-04-12 13:54:48,540 - phyluce_align_get_align_summary_data - INFO - [Characters]
↪ '-' is present 95,942 times
2015-04-12 13:54:48,540 - phyluce_align_get_align_summary_data - INFO - [Characters]
↪ 'A' is present 498,643 times
2015-04-12 13:54:48,540 - phyluce_align_get_align_summary_data - INFO - [Characters]
↪ 'C' is present 314,200 times
2015-04-12 13:54:48,540 - phyluce_align_get_align_summary_data - INFO - [Characters]
↪ 'G' is present 300,494 times
2015-04-12 13:54:48,540 - phyluce_align_get_align_summary_data - INFO - [Characters]
↪ 'T' is present 508,154 times
2015-04-12 13:54:48,540 - phyluce_align_get_align_summary_data - INFO - ===== ↵
↪ Completed phyluce_align_get_align_summary_data =====

```

## Alignment cleaning

If you look at the alignments we currently have, you will notice that each alignment contains the locus name along with the taxon name. This is not what we want downstream, but it does enable us to ensure the correct data went into each alignment. So, we need to clean our alignments. For the remainder of this tutorial, we will work with the `Gblocks` trimmed alignments, so we will clean those alignments:

```

# make sure we are in the correct directory
cd uce-tutorial/taxon-sets/all

# align the data - turn off trimming and output FASTA
phyluce_align_remove_locus_name_from_nexus_lines \
  --alignments mafft-nexus-internal-trimmed-gblocks \
  --output mafft-nexus-internal-trimmed-gblocks-clean \
  --cores 12 \
  --log-path log

```

The output should be similar to:

```

2015-04-12 14:08:19,925 - phyluce_align_remove_locus_name_from_nexus_lines - INFO - ↵
↪=== Starting phyluce_align_remove_locus_name_from_nexus_lines ===
2015-04-12 14:08:19,925 - phyluce_align_remove_locus_name_from_nexus_lines - INFO - ↵
↪Version: git a6a957a
2015-04-12 14:08:19,925 - phyluce_align_remove_locus_name_from_nexus_lines - INFO - ↵
↪Argument --alignments: /scratch/uce-tutorial/taxon-sets/all/mafft-nexus-internal-
↪trimmed-gblocks
2015-04-12 14:08:19,925 - phyluce_align_remove_locus_name_from_nexus_lines - INFO - ↵
↪Argument --cores: 1
2015-04-12 14:08:19,926 - phyluce_align_remove_locus_name_from_nexus_lines - INFO - ↵
↪Argument --input_format: nexus

```

```
2015-04-12 14:08:19,926 - phyluce_align_remove_locus_name_from_nexus_lines - INFO -
↳Argument --log_path: None
2015-04-12 14:08:19,926 - phyluce_align_remove_locus_name_from_nexus_lines - INFO -
↳Argument --output: /scratch/uce-tutorial/taxon-sets/all/mafft-nexus-internal-
↳trimmed-gblocks-clean
2015-04-12 14:08:19,926 - phyluce_align_remove_locus_name_from_nexus_lines - INFO -
↳Argument --output_format: nexus
2015-04-12 14:08:19,926 - phyluce_align_remove_locus_name_from_nexus_lines - INFO -
↳Argument --taxa: None
2015-04-12 14:08:19,926 - phyluce_align_remove_locus_name_from_nexus_lines - INFO -
↳Argument --verbosity: INFO
2015-04-12 14:08:19,926 - phyluce_align_remove_locus_name_from_nexus_lines - INFO -
↳Getting alignment files
Running.....[continued]
2015-04-12 14:08:22,001 - phyluce_align_remove_locus_name_from_nexus_lines - INFO -
↳Taxon names in alignments: alligator_mississippiensis, anolis_carolinensis, gallus_
↳gallus, mus_musculus
2015-04-12 14:08:22,001 - phyluce_align_remove_locus_name_from_nexus_lines - INFO -
↳=== Completed phyluce_align_remove_locus_name_from_nexus_lines ==
```

The current directory structure should look like (I've collapsed a number of branches in the tree):

```
uce-tutorial
+- assembly.conf
...
+- taxon-sets
|   +- all
|       +- all-taxa-incomplete.conf
|       +- all-taxa-incomplete.fasta
|       +- all-taxa-incomplete.incomplete
|       +- exploded-fastas
|       +- log
|       +- mafft-nexus-edge-trimmed
|       +- mafft-nexus-internal-trimmed
|       +- mafft-nexus-internal-trimmed-gblocks
|       +- mafft-nexus-internal-trimmed-gblocks-clean
|           +- uce-1008.nexus
|           +- uce-1014.nexus
|           +- uce-1039.nexus
|           ...
|           +- uce-991.nexus
...
+- uce-search-results
```

Now, if you look at the alignments, you will see that the locus names are removed. We're ready to generate our final data matrices.

## Final data matrices

For the most part, I analyze 75% and 95% complete matrices, where “completeness” for the 75% matrix means that, in a study of 100 taxa (total), all alignments will contain at least 75 of these 100 taxa. Similarly, for the 95% matrix, in a study of 100 taxa, all alignments will contain 95 of these 100 taxa.

**Attention:** Notice that this metric for completeness does not pay attention to which taxa are in which alignments - so the 75%, above, does **not** mean that a given taxon will have data in 75 of 100 alignments.



To create a 75% data matrix, run the following. Notice that the integer following `--taxa` is the **total** number of organisms in the study.

```
# make sure we are in the correct directory
cd uce-tutorial/taxon-sets/all

# the integer following --taxa is the number of TOTAL taxa
# and I use "75p" to denote the 75% complete matrix
phyluce_align_get_only_loci_with_min_taxa \
  --alignments mafft-nexus-internal-trimmed-gblocks-clean \
  --taxa 4 \
  --percent 0.75 \
  --output mafft-nexus-internal-trimmed-gblocks-clean-75p \
  --cores 12 \
  --log-path log
```

The output should look like the following:

```
2015-04-12 14:17:37,589 - phyluce_align_get_only_loci_with_min_taxa - INFO - =====_
↳Starting phyluce_align_get_only_loci_with_min_taxa =====
2015-04-12 14:17:37,589 - phyluce_align_get_only_loci_with_min_taxa - INFO - Version:_
↳git a6a957a
2015-04-12 14:17:37,589 - phyluce_align_get_only_loci_with_min_taxa - INFO - Argument_
↳--alignments: /scratch/uce-tutorial/taxon-sets/all/mafft-nexus-internal-trimmed-
↳gblocks-clean
2015-04-12 14:17:37,589 - phyluce_align_get_only_loci_with_min_taxa - INFO - Argument_
↳--cores: 12
2015-04-12 14:17:37,590 - phyluce_align_get_only_loci_with_min_taxa - INFO - Argument_
↳--input_format: nexus
2015-04-12 14:17:37,590 - phyluce_align_get_only_loci_with_min_taxa - INFO - Argument_
↳--log_path: /scratch/uce-tutorial/taxon-sets/all/log
2015-04-12 14:17:37,590 - phyluce_align_get_only_loci_with_min_taxa - INFO - Argument_
↳--output: /scratch/uce-tutorial/taxon-sets/all/mafft-nexus-internal-trimmed-gblocks-
↳clean-75p
2015-04-12 14:17:37,590 - phyluce_align_get_only_loci_with_min_taxa - INFO - Argument_
↳--percent: 0.75
2015-04-12 14:17:37,590 - phyluce_align_get_only_loci_with_min_taxa - INFO - Argument_
↳--taxa: 4
2015-04-12 14:17:37,590 - phyluce_align_get_only_loci_with_min_taxa - INFO - Argument_
↳--verbosity: INFO
2015-04-12 14:17:37,590 - phyluce_align_get_only_loci_with_min_taxa - INFO - Getting_
↳alignment files
2015-04-12 14:17:37,780 - phyluce_align_get_only_loci_with_min_taxa - INFO - Copied_
↳913 alignments of 913 total containing ≥ 0.75 proportion of taxa (n = 3)
2015-04-12 14:17:37,780 - phyluce_align_get_only_loci_with_min_taxa - INFO - =====_
↳Completed phyluce_align_get_only_loci_with_min_taxa =====
```

The current directory structure should look like (I've collapsed a number of branches in the tree):

```
uce-tutorial
+- assembly.conf
...
+- taxon-sets
|   +- all
|       +- all-taxa-incomplete.conf
|       +- all-taxa-incomplete.fasta
|       +- all-taxa-incomplete.incomplete
|       +- exploded-fastas
|       +- log
```

```
|      +- mafft-nexus-edge-trimmed
|      +- mafft-nexus-internal-trimmed
|      +- mafft-nexus-internal-trimmed-gblocks
|      +- mafft-nexus-internal-trimmed-gblocks-clean
|      +- mafft-nexus-internal-trimmed-gblocks-clean-75p
|          +- uce-1008.nexus
|          +- uce-1014.nexus
|          +- uce-1039.nexus
|          ...
|      +- uce-991.nexus
...
+- uce-search-results
```

## Preparing data for RAxML and ExaML

Now that we have our *75p* data matrix completed, we can generate input files for subsequent phylogenetic analysis. For the most part, I favor [ExaBayes](#), [RAxML](#), and [ExaML](#) (usually in that order). Formatting our *75p* data into a phylip file for these programs is rather easy. To do that, run:

```
# make sure we are in the correct directory
cd uce-tutorial/taxon-sets/all

# build the concatenated data matrix
phyluce_align_format_nexus_files_for_raxml \
  --alignments mafft-nexus-internal-trimmed-gblocks-clean-75p \
  --output mafft-nexus-internal-trimmed-gblocks-clean-75p-raxml \
  --charsets \
  --log-path log
```

The output from this program will look like:

```
2015-04-12 14:40:52,276 - phyluce_align_format_nexus_files_for_raxml - INFO - =====_
↳Starting phyluce_align_format_nexus_files_for_raxml =====
2015-04-12 14:40:52,276 - phyluce_align_format_nexus_files_for_raxml - INFO - _
↳Version: git a6a957a
2015-04-12 14:40:52,276 - phyluce_align_format_nexus_files_for_raxml - INFO - _
↳Argument --alignments: /scratch/uce-tutorial/taxon-sets/all/mafft-nexus-internal-
↳trimmed-gblocks-clean-75p
2015-04-12 14:40:52,277 - phyluce_align_format_nexus_files_for_raxml - INFO - _
↳Argument --charsets: True
2015-04-12 14:40:52,277 - phyluce_align_format_nexus_files_for_raxml - INFO - _
↳Argument --log_path: /scratch/uce-tutorial/taxon-sets/all/log
2015-04-12 14:40:52,277 - phyluce_align_format_nexus_files_for_raxml - INFO - _
↳Argument --nexus: False
2015-04-12 14:40:52,277 - phyluce_align_format_nexus_files_for_raxml - INFO - _
↳Argument --output: /scratch/uce-tutorial/taxon-sets/all/mafft-nexus-internal-
↳trimmed-gblocks-clean-75p-raxml
2015-04-12 14:40:52,277 - phyluce_align_format_nexus_files_for_raxml - INFO - _
↳Argument --verbosity: INFO
2015-04-12 14:40:52,277 - phyluce_align_format_nexus_files_for_raxml - INFO - Reading_
↳input alignments in NEXUS format
2015-04-12 14:40:53,291 - phyluce_align_format_nexus_files_for_raxml - INFO - _
↳Concatenating files
2015-04-12 14:40:54,242 - phyluce_align_format_nexus_files_for_raxml - INFO - Writing_
↳charsets to /scratch/uce-tutorial/taxon-sets/all/mafft-nexus-internal-trimmed-
↳gblocks-clean-75p-raxml/mafft-nexus-internal-trimmed-gblocks-clean-75p.charsets
```

```
2015-04-12 14:40:54,242 - phyluce_align_format_nexus_files_for_raxml - INFO - Writing
↳ concatenated PHYLIP alignment to /scratch/uce-tutorial/taxon-sets/all/mafft-nexus-
↳ internal-trimmed-gblocks-clean-75p-raxml/mafft-nexus-internal-trimmed-gblocks-clean-
↳ 75p.phylip
2015-04-12 14:40:54,245 - phyluce_align_format_nexus_files_for_raxml - INFO - =====
↳ Completed phyluce_align_format_nexus_files_for_raxml =====
```

**Attention:** Notice that using the `-charsets` flag will output the charsets as well as the concatenated PHYLIP file. You generally want these and the cost for them is low. I would always run this option - even if you later do not use them.

The current directory structure should look like (I've collapsed a number of branches in the tree):

```
uce-tutorial
+- assembly.conf
...
+- taxon-sets
|   +- all
|     +- all-taxa-incomplete.conf
|     +- all-taxa-incomplete.fasta
|     +- all-taxa-incomplete.incomplete
|     +- exploded-fastas
|     +- log
|     +- mafft-nexus-edge-trimmed
|     +- mafft-nexus-internal-trimmed
|     +- mafft-nexus-internal-trimmed-gblocks
|     +- mafft-nexus-internal-trimmed-gblocks-clean
|     +- mafft-nexus-internal-trimmed-gblocks-clean-75p
|     +- mafft-nexus-internal-trimmed-gblocks-clean-75p-raxml
|         +- mafft-nexus-internal-trimmed-gblocks-clean-75p.charsets
|         +- mafft-nexus-internal-trimmed-gblocks-clean-75p.phylip
...
+- uce-search-results
```

Usually, I will create a separate directory to hold the alignment and input/output files for ExaBayes:

```
cp -R mafft-nexus-internal-trimmed-gblocks-clean-75p-raxml mafft-nexus-internal-
↳ trimmed-gblocks-clean-75p-exbayes
```

## RAxML

The above data are ready to analysis in **RAxML**. I usually run **RAxML** in two pieces - first running the “best” tree search, then running the bootstrap replicates, and I almost always run these analyses on the HPC.

```
# make sure we are in the correct directory
cd uce-tutorial/taxon-sets/all/mafft-nexus-internal-trimmed-gblocks-clean-75p-raxml

# get two random numbers
for i in 1 2; do echo $RANDOM; done

# that output the following
# 19877
# 7175
```

```
# run the search for the "best" ML tree
raxmlHPC-PTHREADS-SSE3 \
  -m GTRGAMMA \
  -N 20 \
  -p 19877 \
  -n best \
  -s mafft-nexus-internal-trimmed-gblocks-clean-75p.phylip \
  -T 12

# analyze bootstrap data sets using the autoMRE function of RAxML
raxmlHPC-PTHREADS-SSE3 \
  -m GTRGAMMA \
  -N autoMRE \
  -p 19877 \
  -b 7175 \
  -n bootreps \
  -s mafft-nexus-internal-trimmed-gblocks-clean-75p.phylip \
  -T 12
```

**Warning:** Note that I am using 12 physical CPU cores here (*-T 12*). You need to use the number of physical cores available on *your* machine.

Once those are finished running, we need to reconcile the “best” tree with the ML bootstrap replicates:

```
# reconcile the "best" ML tree with the bootreps
raxmlHPC-SSE3 \
  -m GTRGAMMA \
  -f b \
  -t RAxML_bestTree.best \
  -z RAxML_bootstrap.bootreps
```

## ExaBayes

I have not included [ExaBayes](#) in the conda package for [phyluce](#) because it generally requires MPI and should really be run on an HPC or a large local machine. If you would like to run it on your machine, you’ll need to get it installed on your own.

To run data in [ExaBayes](#), you need the PHYLIP file you just created as well as 2 other files in order to the the program - one giving partition information and another giving configuration options. Create these.

### aln.part

```
DNA, p1=1-499929
```

### config.nexus

```
#NEXUS

begin run;
  numruns 4
```

```
numgen 1e6
end;
```

## Run ExaBayes

Now, we can run ExaBayes:

```
# make sure we are in the correct directory
cd uce-tutorial/taxon-sets/all/mafft-nexus-internal-trimmed-gblocks-clean-75p-exabayes

# get a random number
echo $RANDOM

# this outputs
5341

# ensure the files we created above are in this directory, then
mpirun -np 12 exabayes -f mafft-nexus-internal-trimmed-gblocks-clean-75p.phylip -n_
↪run1 -q aln.part -s 5341 -c config.nexus -R 4
```

**Warning:** Note that I am using 12 physical CPU cores here (*-T 12*). You need to use the number of physical cores available on *your* machine.

Once those are finished running, we can summarize the posterior using the *consense* and *postProcParams* program:

```
consense -f ExaBayes_topologies.run1.phylip.* -n some_descriptive_name_here
postProcParam -f ExaBayes_parameters.run1.phylip.* -n some_descriptive_name_here
```

## Tutorial III: Harvesting UCE Loci From Genomes

In many cases, genomic data exist for some (or many) taxa, and you want to “harvest” those loci from the genome(s) available to you for inclusion in a study. This tutorial is meant to explain how to do this. For this example, we’ll download two genome sequences from web repositories, locate, and extract UCE loci from these genomes for subsequent analysis.

The taxa we are working with will be:

- *Gallus gallus*
- *Alligator mississippiensis*

### Starting directory structure

To keep things clear, we’re going to assume you are working in some directory, which I’ll call *uce-genome*. We’ll be working from the top of this directory in the steps below:

```
uce-genome
```

## Download the data

You can download genome assemblies from any number of sources - some better than others. Here, we're going to download one genome assembly (chicken; galGal4) from the [UCSC Genome Browser](#) and another (alligator) from NCBI. We're using two different sources so you can see some of the differences in the process... and what you might need to do in order to "clean up" a given genome sequence. We'll start with the easy one.

### chicken (galGal4)

The [UCSC Genome Browser](#) is a great resource for lots of data that is easy to find and easy to use. In particular, we're interested in the [UCSC Genome Browser Downloads](#) area, where you can find genome sequence, genome annotations, etc. for many model (and non-model) taxa. In this case, we want the *galGal4* genome sequence, which is the 4th "official" assembly of the chicken genome sequence (AKA *GCA\_000002315.2*). You can find this by navigating from the [UCSC Genome Browser Downloads](#) to the [Chicken](#) section of the page. Under the *galGal4* heading, click on [Full data set](#). This will take you to the data download page for *galGal4* where there is a listing of all the data you can download for the *galGal4* assembly. We're interested in the *galGal4.2bit* file, which is a compressed representation of the genome in *2bit* format. You can either click on this file name to download the file, or navigate to your *uce-genome* folder and:

```
$ cd uce-genome
$ mkdir galGal4
$ cd galGal4
wget http://hgdownload.soe.ucsc.edu/goldenPath/galGal4/bigZips/galGal4.2bit
```

This put a 2bit file in our *uce-genome* directory, so that our directory structure looks like:

```
uce-genome
+- galGal4
   +- galGal4.2bit
```

There are various utilities for dealing with *2bit* files that you can download as part of the [Kent Source Archive](#), a set of programs for dealing with genome-scale data. Probably the most important of these are *faToTwoBit*, which we use below; *twoBitInfo*, which gives us information on a given *2bit* file; and *twoBitToFa* which converts a *2bit* file back to FASTA format. If we run *twoBitInfo* against *galGal4.2bit*, we see something like:

```
$ cd uce-genome/galGal4
$ twoBitInfo galGal4.2bit sizes.tab

$ head -n 5 sizes.tab
chr1      195276750
chr10     19911089
chr10_AADN03010416_random  11080
chr10_AADN03010420_random  8415
chr10_JH375184_random     3009
```

which shows us the scaffolds/contigs and their sizes.

### alligator (allMis1)

Although you can get the alligator genomes from [UCSC](#), we'll download it from [NCBI](#), so that you can see the differences in the process/data. Explaining [NCBI](#) is beyond the scope of what we're doing here, so we'll just navigate to the [NCBI alligator genome assembly page](#). If you click the [Assembly](#) link on the right hand side of the page (under "Related information"), this will take you to the [suite of assemblies](#) for alligator.

We'll download the *Algmis\_Hirise\_1.0* assembly, which is an improvement on the original assembly. To do this, click on the *Algmis\_Hirise\_1.0* link, and look for the [Download the GenBank assembly](#) link on the top right corner of the next page. This will take you to an FTP page, where you want to download `GCA_001541155.1_Algmis_Hirise_1.0_genomic.fna.gz`. You can do this by clicking on the link or by:

```
$ cd uce-genome
$ mkdir allMis2
$ cd allMis2
wget ftp://ftp.ncbi.nlm.nih.gov/genomes/all/GCA_001541155.1_Algmis_Hirise_1.0/GCA_
↪001541155.1_Algmis_Hirise_1.0_genomic.fna.gz
```

This put a gzipped fasta file in our *uce-genome/allMis2* directory, so that our directory structure looks like:

```
uce-genome
+- allMis2
|   +- GCA_001541155.1_Algmis_Hirise_1.0_genomic.fna.gz
+- galGal4
   +- galGal4.2bit
   +- sizes.tab
```

Now, we need to unzip this and have a look at the file:

```
$ cd uce-genome/allMis2
$ gunzip GCA_001541155.1_Algmis_Hirise_1.0_genomic.fna.gz

# take a look at the contents of this file:
$ head -n 1 GCA_001541155.1_Algmis_Hirise_1.0_genomic.fna
>LPUV01000001.1 Alligator mississippiensis ScZkoYb_3522, whole genome shotgun sequence
```

Note that the header has a lot of text in it, and this text is not always good. We're going to convert this file to *2bit* format using *faToTwoBit* from the [Kent Source Archive](#), which will remove everything following the first space in the header line:

```
$ faToTwoBit GCA_001541155.1_Algmis_Hirise_1.0_genomic.fna allMis2.2bit
$ twoBitInfo allMis2.2bit sizes.tab
$ head -n 5 sizes.tab
LPUV01000001.1 1279
LPUV01000002.1 4883
LPUV01000003.1 1105
LPUV01000004.1 22955
LPUV01000005.1 1137523
```

Now that we've converted the fasta file to *2bit* format, we can delete the fasta file:

```
$ rm GCA_001541155.1_Algmis_Hirise_1.0_genomic.fna
```

And your directory structure should look like:

```
uce-genome
+- allMis2
|   +- allMis2.2bit
|   +- sizes.tab
+- galGal4
   +- galGal4.2bit
   +- sizes.tab
```

**Attention:** It's easiest to use *2bit* files of each genome you want to search for UCE loci.

## Finding UCE loci

Now that we've downloaded our genome assemblies, it's time to find those contigs/scaffolds in the assembly that contain UCE loci.

### Get the UCE probes

To do that, we need to get the UCE probe set we want to search for into our directory. Here, we'll search for the UCE 5k probe set. However, you can use whichever probe set you like - you just need to know the path to this file.

```
> cd uce-genome
# download the probe set
wget https://raw.githubusercontent.com/faircloth-lab/uce-probe-sets/master/uce-5k-
probe-set/uce-5k-probes.fasta
```

Now, our directory structure looks like:

```
uce-genome
+- allMis2
|   +- allMis2.2bit
|   +- sizes.tab
+- galGal4
|   +- galGal4.2bit
|   +- sizes.tab
+- uce-5k-probes.fasta
```

### Align the probes to the genomes

We need to align the probes to the genome sequences. The program that we are going to run assumes that each *2bit* genome file is in its own directory (which we have ensured, above).

**Attention:** If you use some other organizational structure, you still need to ensure that each genome sequence is in its own directory. So, if you have some directory *genomes*, the genomes must be organized within that folder like:

```
genomes
+- allMis2
|   +- allMis2.2bit
+- galGal4
   +- galGal4.2bit
```

Now, we need to think of some name for the database to create (here *tutorial3.sqlite*), the name of an output in which to store the *lastz* search results, the path to the genome sequences, the name of the genome sequences, the path to the probe file, and a number of compute cores to use:

```
> cd uce-genome
# run the search
> phyluce_probe_run_multiple_lastzs_sqlite \
```



```
--db tutorial3.sqlite \  
--output tutorial3-genome-lastz \  
--scaffoldlist galGal4 allMis2 \  
--genome-base-path ./ \  
--probefile /nfs/data1/uce-probe-sets/uce-5k-probe-set/uce-5k-probes.fasta \  
--cores 12
```

The program will create some output that looks like:

```
Running against galGal4.2bit  
Running with the --huge option. Chunking files into 10000000 bp...  
/tmp/tmpXup1D.lastz  
  
< .. snip .. >  
  
Cleaning up the chunked files...  
Cleaning /home/bcf/tmp/uce-genome/tutorial3-genome-lastz/uce-5k-probes.fasta_v_  
→galGal4.lastz  
Creating galGal4 table  
Inserting data to galGal4 table  
  
Running against allMis2.2bit  
Running with the --huge option. Chunking files into 10000000 bp...  
Running the targets against 109 queries...  
/tmp/tmpuuHOpS.fasta  
  
< .. snip .. >  
  
Cleaning up the chunked files...  
Cleaning /home/bcf/tmp/uce-genome/tutorial3-genome-lastz/uce-5k-probes.fasta_v_  
→allMis2.lastz  
Creating allMis2 table  
Inserting data to allMis2 table
```

The directory structure should now look like this:

```
uce-genome  
+- allMis2  
| +- allMis2.2bit  
| +- sizes.tab  
+- galGal4  
| +- galGal4.2bit  
| +- sizes.tab  
+- tutorial3-genome-lastz  
| +- uce-5k-probes.fasta_v_allMis2.lastz.clean  
| +- uce-5k-probes.fasta_v_galGal4.lastz.clean  
+- tutorial3.sqlite  
+- uce-5k-probes.fasta
```

### Extracting FASTA sequence matching UCE loci from genome sequences

Once you have run the search, you can extract the loci identified during the search from each respective genome sequence plus some sequence flanking each locus (at a distance you can specify). Before you do that, however, you need to create a configuration file that tells the program where to find each genome. In the case above, that file would look like so (the full path to my working directory is */home/bcf/tmp/uce-genome*):

```
[scaffolds]
galGal4:/home/bcf/tmp/uce-genome/galGal4/galGal4.2bit
allMis2:/home/bcf/tmp/uce-genome/allMis2/allMis2.2bit
```

**Attention:** Be careful to make sure that your capitalization is consistent with your file names.

With that file created as *genomes.conf*, our directory structure looks like:

```
uce-genome
+- allMis2
| +- allMis2.2bit
| +- sizes.tab
+- galGal4
| +- galGal4.2bit
| +- sizes.tab
+- genomes.conf
+- tutorial3-genome-lastz
| +- uce-5k-probes.fasta_v_allMis2.lastz.clean
| +- uce-5k-probes.fasta_v_galGal4.lastz.clean
+- tutorial3.sqlite
+- uce-5k-probes.fasta
```

Now, we can extract FASTA data from each genome for each UCE locus. To do this, we need to input the path to the *lastz* files from above, the path to the *conf* file we just created, the amount of flanking sequence (to each side) that we would like to slice, a name pattern, matching the *lastz* files that we would like to use, and the name of the output directory we want to create:

```
phyluce_probe_slice_sequence_from_genomes \
  --lastz tutorial3-genome-lastz \
  --conf genomes.conf \
  --flank 500 \
  --name-pattern "uce-5k-probes.fasta_v_{}.lastz.clean" \
  --output tutorial-genome-fasta
```

You should see output similar to:

```
2016-06-01 16:02:18,907 - Phyluce - INFO - ===== Starting Phyluce:
↳Slice Sequence =====
2016-06-01 16:02:18,908 - Phyluce - INFO - ----- Working on galGal4
↳genome -----
2016-06-01 16:02:18,908 - Phyluce - INFO - Reading galGal4 genome
2016-06-01 16:02:28,036 - Phyluce - INFO - galGal4: 4966 uces, 35 dupes, 4931 non-
↳dupes, 2 orient drop, 40 length drop, 4889 written
2016-06-01 16:02:28,036 - Phyluce - INFO - ----- Working on allMis2
↳genome -----
2016-06-01 16:02:28,037 - Phyluce - INFO - Reading allMis2 genome
2016-06-01 16:02:37,230 - Phyluce - INFO - allMis2: 4830 uces, 7 dupes, 4823 non-
↳dupes, 2 orient drop, 3 length drop, 4818 written
```

And, your directory structure should look like:

```
uce-genome
+- allMis2
| +- allMis2.2bit
| +- sizes.tab
```

```

+- galGal4
|   +- galGal4.2bit
|   +- sizes.tab
+- genomes.conf
+- tutorial3-genome-lastz
|   +- uce-5k-probes.fasta_v_allMis2.lastz.clean
|   +- uce-5k-probes.fasta_v_galGal4.lastz.clean
+- tutorial3.sqlite
+- tutorial-genome-fasta
|   +- allmis2.fasta
|   +- galgal4.fasta
+- uce-5k-probes.fasta

```

The UCE contig sequence are in each respective file within *tutorial-genome-fasta*.

### Using the extracted sequences in downstream analyses

The easiest way for you to use the extracted sequences is to symlink them into an appropriate *contigs* folder that resulted from a [PHYLUCE](#) assembly process and then proceed with the *Extracting UCE loci* procedure.

For more information on the structure of this folder, look at the *Assemble the data* section of *Tutorial I: UCE Phylogenomics* for more information.

**Attention:** Although we have already extracted the UCE loci from each genome sequence and even though it seems redundant to go back through the *Extracting UCE loci* process, this is the best path forward.

## Tutorial IV: Identifying UCE Loci and Designing Baits To Target Them

The first few tutorials have given you a feel for how to perform phylogenetic/phylogeographic analyses using existing probe sets, new data, and genomes. However, what if you are working on a set of organisms without a probe set targeting conserved loci? How do you identify those loci and design baits to target them? This Tutorial shows you how to do that.

In the examples below, we'll follow a UCE identification and probe design workflow using data from Coleoptera (beetles). Although you can follow the entire tutorial from beginning to end, I've also made the BAM files containing mapped reads available, which lets you skip the computationally expensive step of performing read simulation and alignment.

### Starting directory structure

To keep things clear, we're going to assume you are working in some directory, which I'll call *uce-coleoptera*. We'll be working from the top of this directory in the steps below:

```
uce-coleoptera
```

### Data download and preparation

## Download the genomes

**Attention:** You do not necessarily need to do this as part of this tutorial for UCE identification and probe design - If you only want to follow the steps for locus identification (skipping probe design and in-silico testing), you can simply download the prepared FASTQ/BAM files from figshare.

Make a directory to hold the genome sequences:

```
> mkdir genomes
```

Now, get the genome sequences:

```
# Anoplophora glabripennis (Asian longhorned beetle)
> wget ftp://ftp.ncbi.nlm.nih.gov/genomes/all/GCA/000/390/285/GCA_000390285.1_Agla_1.
↪0/GCA_000390285.1_Agla_1.0_genomic.fna.gz

# Agrilus planipennis (emerald ash borer)
> wget ftp://ftp.ncbi.nlm.nih.gov/genomes/all/GCA/000/699/045/GCA_000699045.1_Apla_1.
↪0/GCA_000699045.1_Apla_1.0_genomic.fna.gz

# Leptinotarsa decemlineata (Colorado potato beetle)
> wget ftp://ftp.ncbi.nlm.nih.gov/genomes/all/GCA/000/500/325/GCA_000500325.1_Ldec_1.
↪5/GCA_000500325.1_Ldec_1.5_genomic.fna.gz

# Onthophagus taurus (beetles)
> wget ftp://ftp.ncbi.nlm.nih.gov/genomes/all/GCA/000/648/695/GCA_000648695.1_Otau_1.
↪0/GCA_000648695.1_Otau_1.0_genomic.fna.gz

# Dendroctonus ponderosae (mountain pine beetle)
> wget ftp://ftp.ncbi.nlm.nih.gov/genomes/all/GCA/000/355/655/GCA_000355655.1_
↪DendPond_male_1.0/GCA_000355655.1_DendPond_male_1.0_genomic.fna.gz

# Tribolium castaneum (red flour beetle)
> wget ftp://ftp.ncbi.nlm.nih.gov/genomes/all/GCA/000/002/335/GCA_000002335.2_Tcas_3.
↪0/GCA_000002335.2_Tcas_3.0_genomic.fna.gz

# Mengenilla moldrzyki (twisted-wing parasites) [Outgroup]
> wget ftp://ftp.ncbi.nlm.nih.gov/genomes/all/GCA/000/281/935/GCA_000281935.1_Memo_1.
↪0/GCA_000281935.1_Memo_1.0_genomic.fna.gz
```

You need to unzip all of the genome sequences

```
> gunzip *
```

Finally, your directory structure should look something like:

```
uce-coleoptera
+- genomes
  +- GCA_000002335.2_Tcas_3.0_genomic.fna
  +- GCA_000281935.1_Memo_1.0_genomic.fna
  +- GCA_000355655.1_DendPond_male_1.0_genomic.fna
  +- GCA_000390285.1_Agla_1.0_genomic.fna
  +- GCA_000500325.1_Ldec_1.5_genomic.fna
  +- GCA_000648695.1_Otau_1.0_genomic.fna
  +- GCA_000699045.1_Apla_1.0_genomic.fna
```

## Cleanup the genome sequences

**Attention:** You do not necessarily need to do this as part of this tutorial for UCE identification and probe design - If you only want to follow the steps for locus identification (skipping probe design and in-silico testing), you can simply download the prepared FASTQ/BAM files from [figshare](#).

When you get genome sequences from NCBI, the FASTA headers of most scaffold/contigs contain a lot of extra cruft that can cause problems with some of the steps in the UCE identification and probe design workflow. I usually remove the extra stuff, maintaining only the accession number information for each contig / scaffold. To do that, I use a little python script that looks like the following:

```
from Bio import SeqIO
with open("Name_of_Genome_File.fna", "rU") as infile:
with open("outfileName.fasta", "w") as outf:
    for seq in SeqIO.parse(infile, 'fasta'):
        seq.name = ""
        seq.description = ""
        outf.write(seq.format('fasta'))
```

For these genome assemblies, the important bits are in the FASTA header just before the space in the name. The code above basically keeps this information before the space and discards the remaining FASTA header information.

In the case of the genomes above, here are the commands I ran (note also that this creates an output file with a different name from the input file):

```
from Bio import SeqIO

# Anoplophora glabripennis (Asian longhorned beetle)
with open("GCA_000390285.1_Agla_1.0_genomic.fna", "rU") as infile:
    with open("anoGla1.fasta", "w") as outf:
        for seq in SeqIO.parse(infile, 'fasta'):
            seq.name = ""
            seq.description = ""
            outf.write(seq.format('fasta'))

# Agrilus planipennis (emerald ash borer)
with open("GCA_000699045.1_Apla_1.0_genomic.fna", "rU") as infile:
    with open("agrPlal.fasta", "w") as outf:
        for seq in SeqIO.parse(infile, 'fasta'):
            seq.name = ""
            seq.description = ""
            outf.write(seq.format('fasta'))

# Dendroctonus ponderosae (mountain pine beetle)
with open("GCA_000355655.1_DendPond_male_1.0_genomic.fna", "rU") as infile:
    with open("denPon1.fasta", "w") as outf:
        for seq in SeqIO.parse(infile, 'fasta'):
            seq.name = ""
            seq.description = ""
            outf.write(seq.format('fasta'))

# Leptinotarsa decemlineata (Colorado potato beetle)
with open("GCA_000500325.1_Ldec_1.5_genomic.fna", "rU") as infile:
    with open("lepDec1.fasta", "w") as outf:
        for seq in SeqIO.parse(infile, 'fasta'):
            seq.name = ""
```

```
seq.description = ""
outf.write(seq.format('fasta'))

# Mengenilla moldrzyki (twisted-wing parasites) [Outgroup]
with open("GCA_000281935.1_Memo_1.0_genomic.fna", "rU") as infile:
    with open("menMoll.fasta", "w") as outf:
        for seq in SeqIO.parse(infile, 'fasta'):
            seq.name = ""
            seq.description = ""
            outf.write(seq.format('fasta'))

# Onthophagus taurus (beetles)
with open("GCA_000648695.1_Otau_1.0_genomic.fna", "rU") as infile:
    with open("ontTaul.fasta", "w") as outf:
        for seq in SeqIO.parse(infile, 'fasta'):
            seq.name = ""
            seq.description = ""
            outf.write(seq.format('fasta'))

# Tribolium castaneum (red flour beetle)
with open("GCA_000002335.2_Tcas_3.0_genomic.fna", "rU") as infile:
    with open("triCas1.fasta", "w") as outf:
        for seq in SeqIO.parse(infile, 'fasta'):
            seq.name = ""
            seq.description = ""
            outf.write(seq.format('fasta'))
```

Now, you can remove the original files downloaded from NCBI:

```
rm *.fna
```

And, your directory structure should look something like:

```
uce-coleoptera
+- genomes
  +- anoGlal.fasta
  +- agrPlal.fasta
  +- denPonl.fasta
  +- lepDecl.fasta
  +- menMoll.fasta
  +- ontTaul.fasta
  +- triCas1.fasta
```

### Put genomes in their own directories

Because of some historical reasons (and how I organize our lab data), each genome sequence needs to be in its own directory. We can do that pretty easily by running:

```
> cd uce-coleoptera
> cd genomes
> for critter in *; do mkdir ${critter%.*}; mv $critter ${critter%.*}; done
```

Now the directory structure looks like:

```
uce-coleoptera
+- genomes
```

```

+- agrPlal
| +- agrPlal.fasta
+- anoGla1
| +- anoGla1.fasta
+- denPon1
| +- denPon1.fasta
+- lepDecl
| +- lepDecl.fasta
+- menMoll
| +- menMoll.fasta
+- ontTau1
| +- ontTau1.fasta
+- triCas1
  +- triCas1.fasta

```

### Convert genomes to 2bit format

Later during the workflow, we're going to need to have our genomes in *2bit* format, which is a binary format for genomic data that is easier and faster to work with than FASTA format. We'll use a BASH script to convert all of the sequences, above, to *2bit* format:

```

> cd uce-coleoptera
> cd genomes
> for critter in *; do faToTwoBit $critter/$critter.fasta $critter/${critter%.*}.2bit;
↪ done

```

Now the directory structure looks like:

```

uce-coleoptera
+- genomes
  +- agrPlal
    | +- agrPlal.2bit
    | +- agrPlal.fasta
  +- anoGla1
    | +- anoGla1.2bit
    | +- anoGla1.fasta
  +- denPon1
    | +- denPon1.2bit
    | +- denPon1.fasta
  +- lepDecl
    | +- lepDecl.2bit
    | +- lepDecl.fasta
  +- menMoll
    | +- menMoll.2bit
    | +- menMoll.fasta
  +- ontTau1
    | +- ontTau1.2bit
    | +- ontTau1.fasta
  +- triCas1
    +- triCas1.2bit
    +- triCas1.fasta

```

## Simulate reads from genomes

**Attention:** You do not necessarily need to take this step as part of the tutorial - you can simply download prepared, simulated FASTQ files from [figshare](#).

In order to locate UCE loci across a selection of different genomes, we're going to align reads from each taxon, above, to a reference genome sequence (the "base" genome sequence) using a permissive raw read aligner. You can use reads from low-coverage, genome scans or you can use reads simulated from particular genomes. In this tutorial, we're going to use this latter approach and simulate reads (without sequencing error) from the genomes that we will align to the base genome. To accomplish this, we'll use *art*, which is a robust read simulator that is reasonably flexible.

Because we're using simulated reads to locate UCE loci, we want to turn off the built-in feature of *art* that adds some sequencing error to simulated reads. This results in a general form of the call to *art* that looks like:

```
> art_illumina \
  --paired \
  --in ~/path/to/input/genome.fasta \
  --out prefix-of-output-file \
  --len 100 \
  --fcov 2 \
  --mflen 200 \
  --sdev 150 \
  -ir 0.0 -ir2 0.0 -dr 0.0 -dr2 0.0 -qs 100 -qs2 100 -na
```

This simulates reads from the *-in genome.fasta* that are 100 bp in length, cover the genome randomly to roughly 2X, have an insert size of 200 bp, and have an inserts size standard deviation of 150. The last line turns off the simulation of sequencing error in each of these reads and also turns off the creation of an alignment file showing where the reads came from in the genome sequence.

In the case of the Coleoptera genomes you downloaded, here are the commands I ran to simulate the read data that we will use in the next step (note also that this creates an output file with a different name from the input file). First, create a directory to hold the simulated read data:

```
> cd uce-coleoptera
> mkdir reads
> cd reads
```

Then, simulate the reads using *art*:

```
> art_illumina \
  --paired \
  --in ../genomes/agrPla1/agrPla1.fasta \
  --out agrPla1-pe100-reads \
  --len 100 --fcov 2 --mflen 200 --sdev 150 -ir 0.0 -ir2 0.0 -dr 0.0 -dr2 0.0 -qs_
↪100 -qs2 100 -na

> art_illumina \
  --paired \
  --in ../genomes/anoGla1/anoGla1.fasta \
  --out anoGla1-pe100-reads \
  --len 100 --fcov 2 --mflen 200 --sdev 150 -ir 0.0 -ir2 0.0 -dr 0.0 -dr2 0.0 -qs_
↪100 -qs2 100 -na

> art_illumina \
  --paired \
  --in ../genomes/denPon1/denPon1.fasta \
```



```

--out denPon1-pe100-reads \
--len 100 --fcov 2 --mflen 200 --sdev 150 -ir 0.0 -ir2 0.0 -dr 0.0 -dr2 0.0 -qs_
↪100 -qs2 100 -na

> art_illumina \
  --paired \
  --in ../genomes/lepDec1/lepDec1.fasta \
  --out lepDec1-pe100-reads \
  --len 100 --fcov 2 --mflen 200 --sdev 150 -ir 0.0 -ir2 0.0 -dr 0.0 -dr2 0.0 -qs_
↪100 -qs2 100 -na

> art_illumina \
  --paired \
  --in ../genomes/ontTaul/ontTaul.fasta \
  --out ontTaul-pe100-reads \
  --len 100 --fcov 2 --mflen 200 --sdev 150 -ir 0.0 -ir2 0.0 -dr 0.0 -dr2 0.0 -qs_
↪100 -qs2 100 -na

```

Now, you should see 2 read files for each taxon. We are going to “break” the pairs by merging the read information together, then we are going to zip the resulting file that contains the *R1* and *R2* data. We can accomplish this pretty easily using a quick BASH script:

**Attention:** Note that we’ve dropped *menMol1* from the read simulation process. This is largely because it is an outgroup to beetles. We’ll use it later, when we’re performing *in silico* tests of the UCE bait set.

```

for critter in agrPlal1 anoGla1 denPon1 lepDec1 ontTaul;
do
  echo "working on $critter";
  touch $critter-pe100-reads.fq;
  cat $critter-pe100-reads1.fq > $critter-pe100-reads.fq;
  cat $critter-pe100-reads2.fq >> $critter-pe100-reads.fq;
  rm $critter-pe100-reads1.fq;
  rm $critter-pe100-reads2.fq;
  gzip $critter-pe100-reads.fq;
done;

```

If we take a look at our directory structure, it should look like:

```

uce-coleoptera
+- genomes (collapsed)
+- reads
  +- anoGla1-pe100-reads.fq.gz
  +- agrPlal1-pe100-reads.fq.gz
  +- denPon1-pe100-reads.fq.gz
  +- lepDec1-pe100-reads.fq.gz
  +- ontTaul-pe100-reads.fq.gz

```

## Read alignment to the base genome

**Attention:** You do not necessarily need to run this step as part of the tutorial - you can simply download the prepared, [BAM files from figshare](#).

Because we also provide the BAM files created below, you can choose to just start with the BAM files in the *Conserved locus identification* section.

## Prepare the base genome

Now that we have read data representing each of our exemplar taxa, we need to align these reads to the “base” genome sequence, in this case the genome sequence of *Tribolium castaneum* (aka *triCas1*). We selected this assembly as the “base” genome because of its age (i.e., better-assembled) and level of annotation.

We will perform the read alignments to *triCas1* using the permissive read aligner, *stampy*, which works well when aligning sequences to a divergent reference sequence. However, before running the alignments, we need to prepare the base genome. And, before we do that, let’s create a directory to work in:

```
> cd uce-coleoptera
> mkdir base
> cd base
```

Now, let’s copy the base genome to this directory (for simplicity):

```
> cp ../genomes/triCas1.fasta ./
```

If we take a look at our directory structure, it now looks like:

```
uce-coleoptera
+- base
| +- triCas1.fasta
+- genomes
| +- agrPla1
| | +- agrPla1.2bit
| | +- agrPla1.fasta
| +- anoGla1
| | +- anoGla1.2bit
| | +- anoGla1.fasta
| +- denPon1
| | +- denPon1.2bit
| | +- denPon1.fasta
| +- lepDecl
| | +- lepDecl.2bit
| | +- lepDecl.fasta
| +- menMoll
| | +- menMoll.2bit
| | +- menMoll.fasta
| +- ontTau1
| | +- ontTau1.2bit
| | +- ontTau1.fasta
| +- triCas1
| | +- triCas1.2bit
| | +- triCas1.fasta
+- reads
  +- anoGla1-pe100-reads
  +- agrPla1-pe100-reads
  +- denPon1-pe100-reads
  +- lepDecl-pe100-reads
  +- ontTau1-pe100-reads
```

Now, we need to run the commands to prepare the *triCas1* genome for use with *stampy*:

```
> cd uce-coleoptera
> cd base
> stampy.py --species="tribolium-castaneum" --assembly="triCas1" -G triCas1 triCas1.
↳fasta
> stampy.py -g triCas1 -H triCas1
```

If we look at our directory structure, it should look like:

```
uce-coleoptera
+- base
| +- triCas1.fasta
| +- triCas1.sthash
| +- triCas1.stidx
+- genomes (collapsed)
+- reads
  +- anoGlal-pe100-reads
  +- agrPlal-pe100-reads
  +- denPon1-pe100-reads
  +- lepDecl-pe100-reads
  +- ontTaul-pe100-reads
```

### Align reads to the base genome

**Attention:** You do not necessarily need to run this step as part of the tutorial - you can simply download the prepared, [BAM files from figshare](#).

Because we also provide the BAM files created below, you can choose to just start with the BAM files in the *Conserved locus identification* section.

Now that we've prepared our base genome, we need to perform the actual alignment of the simulated reads to the base genome. And, before we do that, let's create a directory to hold the resulting alignments:

```
> cd uce-coleoptera
> mkdir alignments
```

Our resulting directory structure should now look like:

```
uce-coleoptera
+- alignments
+- base
| +- triCas1.fasta
| +- triCas1.sthash
| +- triCas1.stidx
+- genomes (collapsed)
+- reads
  +- anoGlal-pe100-reads
  +- agrPlal-pe100-reads
  +- denPon1-pe100-reads
  +- lepDecl-pe100-reads
  +- ontTaul-pe100-reads
```

Now, we need to perform the alignments on a taxon-by-taxon basis (and/or you can run these in parallel using HPC). To do this easily (and on a local computer) we can use a BASH script to run the alignments serially:

**Warning:** Note that I am using 16 physical CPU cores (*\$cores*) to do this work. You need to use the number of physical cores available on *your* machine.

```
export cores=16
export base=triCas1
export base_dir=$HOME/uce-coleoptera/alignments
for critter in agrPlal anoGlal denPonl lepDecl ontTaul;
do
    export reads=$critter-pe100-reads.fq.gz;
    mkdir -p $base_dir/$critter;
    cd $base_dir/$critter;
    stampy.py --maxbasequal 93 -g ../../base/$base -h ../../base/$base
    --substitutionrate=0.05 -t$cores --insertsize=400 -M
    ../../reads/$reads | samtools view -Sb - > $critter-to-$base.bam;
done;
```

This code basically loops over each exemplar genomes, aligns the reads to the base genome sequence, and converts the resulting output to **BAM** format (which is a binary, compressed version of **SAM** format).

When the alignments have completed, your directory structure should look something like:

```
uce-coleoptera
+- alignments
| +- anoGlal
| | +- anoGlal-to-triCas1.bam
| +- agrPlal
| | +- agrPlal-to-triCas1.bam
| +- denPonl
| | +- denPonl-to-triCas1.bam
| +- lepDecl
| | +- lepDecl-to-triCas1.bam
| +- ontTaul
| | +- ontTaul-to-triCas1.bam
+- base
| +- triCas1.fasta
| +- triCas1.sthash
| +- triCas1.stidx
+- genomes (collapsed)
+- reads
  +- anoGlal-pe100-reads
  +- agrPlal-pe100-reads
  +- denPonl-pe100-reads
  +- lepDecl-pe100-reads
  +- ontTaul-pe100-reads
```

Now, these **BAM** files are pretty large because they contain all mapped as well as all unmapped reads. We want to remove those unmapped reads, which will also reduce file-size. We can do that using **samtools view** and a **BASH** script. We're also going to create a directory named *all* and symlink all of the reduced **BAM** files to this directory.

**Warning:** The script, as-written, removes the **BAM** files containing both mapped and unmapped reads. If you don't want to do this, remove the `rm $critter/$critter-to-triCas1.bam;` line.

```
> cd uce-coleoptera
> cd alignments
```

```
> mkdir all
> for critter in agrPlal anoGla1 denPon1 lepDec1 ontTaul;
do
    samtools view -h -F 4 -b $critter/$critter-to-triCas1.bam > $critter/$critter-
↳to-triCas1-MAPPING.bam;
    rm $critter/$critter-to-triCas1.bam;
    ln -s ../$critter/$critter-to-triCas1-MAPPING.bam all/$critter-to-triCas1-
↳MAPPING.bam;
done;
```

Now, your directory structure should look something like:

```
uce-coleoptera
+- alignments
| +- anoGla1
| | +- anoGla1-to-triCas1-MAPPING.bam
| +- all
| | +- agrPlal-to-triCas1-MAPPING.bam -> ../agrPlal/agrPlal-to-triCas1-MAPPING.
↳bam
| | +- anoGla1-to-triCas1-MAPPING.bam -> ../anoGla1/anoGla1-to-triCas1-MAPPING.
↳bam
| | +- denPon1-to-triCas1-MAPPING.bam -> ../denPon1/denPon1-to-triCas1-MAPPING.
↳bam
| | +- lepDec1-to-triCas1-MAPPING.bam -> ../lepDec1/lepDec1-to-triCas1-MAPPING.
↳bam
| | +- ontTaul-to-triCas1-MAPPING.bam -> ../ontTaul/ontTaul-to-triCas1-MAPPING.
↳bam
| +- agrPlal
| | +- agrPlal-to-triCas1-MAPPING.bam
| +- denPon1
| | +- denPon1-to-triCas1-MAPPING.bam
| +- lepDec1
| | +- lepDec1-to-triCas1-MAPPING.bam
| +- ontTaul
| | +- ontTaul-to-triCas1-MAPPING.bam
+- base
| +- triCas1.fasta
| +- triCas1.sthash
| +- triCas1.stidx
+- genomes (collapsed)
+- reads
  +- anoGla1-pe100-reads
  +- agrPlal-pe100-reads
  +- denPon1-pe100-reads
  +- lepDec1-pe100-reads
  +- ontTaul-pe100-reads
```

These \*-MAPPING.bam files are available from [figshare](#).

### What it all means

Basically, because we've now mapped simulated (or actual) sequence data from several exemplar taxa to a closely-related "base" genome sequence, we've essentially identified putatively orthologous loci shared between the exemplar taxa and the base taxon. These conserved regions are where the simulated (or actual) sequence data mapped with a sequence divergence of < 5%.

Now, we need to filter this large number of conserved regions to remove things like repetitive regions, but also to find

which loci are shared among *all* exemplar taxa - not simply a single exemplar and the base taxon.

## Conserved locus identification

You can download the [alignment data](#) generated in the steps above from [figshare](#) for use in subsequent steps, rather than generating them yourself (thus, saving you time).

**Attention:** If you are starting the tutorial at this position after downloading the *\*-MAPPING.bam* files from [figshare](#), then you will need to create a directory to work in named *uce-coleoptera* and then place all of the *\*-MAPPING.bam* files in a subdirectory of *uce-coleoptera* names *alignments/all*. Your resulting directory structure should look like:

```
uce-coleoptera
+- alignments
| +- all
|   +- agrPlal-to-triCas1-MAPPING.bam -> ../agrPlal/agrPlal-to-triCas1-MAPPING.
|   ↪bam
|   +- anoGla1-to-triCas1-MAPPING.bam -> ../anoGla1/anoGla1-to-triCas1-MAPPING.
|   ↪bam
|   +- denPon1-to-triCas1-MAPPING.bam -> ../denPon1/denPon1-to-triCas1-MAPPING.
|   ↪bam
|   +- lepDec1-to-triCas1-MAPPING.bam -> ../lepDec1/lepDec1-to-triCas1-MAPPING.
|   ↪bam
|   +- ontTau1-to-triCas1-MAPPING.bam -> ../ontTau1/ontTau1-to-triCas1-MAPPING.
|   ↪bam
+- genomes (collapsed)
```

If you want to go beyond conserved locus identification and design probes from the target taxa, you will also need to download the appropriate genomes. See the [Data download and preparation](#) section.

## Convert BAMS to BEDS

In the steps above, we have generated BAM files representing reads that [stampy](#) has mapped to the base genome. Those reads that map align to putatively conserved sequence regions (that we need to filter), and these alignments of mapping reads should reside in our *alignments/all* directory.

To begin the filtering process, we're going to convert each [BAM](#) file to [BED](#) format, which is an interval-based format that is easy and fast to manipulate with a software suite known as [bedtools](#). But before we do that, we are doing to create a *bed* directory to hold all of these [BED](#) format files.

```
> cd uce-coleoptera
# make a directory to hold the BED data
> mkdir bed
> cd bed
# now, convert our *-MAPPING.bam files to BED format
> for i in ../alignments/all/*.bam; do echo $i; bedtools bamtobed -i $i -bed12 >_
↪ `basename $i`.bed; done
```

Your directory structure should look something like:

```
uce-coleoptera
+- alignments
| +- all
```

```
|     +- agrPlal-to-triCas1-MAPPING.bam
|     +- anoGlal-to-triCas1-MAPPING.bam
|     +- denPon1-to-triCas1-MAPPING.bam
|     +- lepDecl-to-triCas1-MAPPING.bam
|     +- ontTaul-to-triCas1-MAPPING.bam
+- bed
|     +- agrPlal-to-triCas1-MAPPING.bam.bed
|     +- anoGlal-to-triCas1-MAPPING.bam.bed
|     +- denPon1-to-triCas1-MAPPING.bam.bed
|     +- lepDecl-to-triCas1-MAPPING.bam.bed
|     +- ontTaul-to-triCas1-MAPPING.bam.bed
+- genomes (collapsed)
```

### Sort the converted BEDs

Before moving forward with the *merge* command, below, we need to sort the resulting BED files, which orders each lines of data in the *BED* file by chromosome/scaffold/contig and position along that chromosome/scaffold/contig. Again, we can do this with some bash scripting:

```
> cd uce-coleoptera
> cd bed
> for i in *.bed; do echo $i; bedtools sort -i $i > ${i%.*}.sort.bed; done
```

Your directory structure should look something like the following (note that I have collapsed the directory listing for *all*):

```
uce-coleoptera
+- alignments
|   +- all (collapsed)
+- bed
|   +- agrPlal-to-triCas1-MAPPING.bam.bed
|   +- agrPlal-to-triCas1-MAPPING.bam.sort.bed
|   +- anoGlal-to-triCas1-MAPPING.bam.bed
|   +- anoGlal-to-triCas1-MAPPING.bam.sort.bed
|   +- denPon1-to-triCas1-MAPPING.bam.bed
|   +- denPon1-to-triCas1-MAPPING.bam.sort.bed
|   +- lepDecl-to-triCas1-MAPPING.bam.bed
|   +- lepDecl-to-triCas1-MAPPING.bam.sort.bed
|   +- ontTaul-to-triCas1-MAPPING.bam.bed
|   +- ontTaul-to-triCas1-MAPPING.bam.sort.bed
+- genomes (collapsed)
```

### Merge overlapping or nearly-overlapping intervals

Because there may be small gaps between proximate regions of conservation (which may result because we're using data that are either from low-coverage, simulated reads or low-coverage *actual* reads) we need to merge together putative regions of conservation that are proximate. Luckily *bedtools* has a handy tool to do that - the *merge* function.

```
> cd uce-coleoptera
> cd bed
> for i in *.bam.sort.bed; do echo $i; bedtools merge -i $i > ${i%.*}.merge.bed; done
```

Your directory structure should look something like the following (note that I have collapsed the directory listing for *all*):

```
uce-coleoptera
+- alignments
| +- all (collapsed)
+- bed
| +- agrPlal-to-triCas1-MAPPING.bam.bed
| +- agrPlal-to-triCas1-MAPPING.bam.sort.bed
| +- agrPlal-to-triCas1-MAPPING.bam.sort.merge.bed
| +- anoGlal-to-triCas1-MAPPING.bam.bed
| +- anoGlal-to-triCas1-MAPPING.bam.sort.bed
| +- anoGlal-to-triCas1-MAPPING.bam.sort.merge.bed
| +- denPon1-to-triCas1-MAPPING.bam.bed
| +- denPon1-to-triCas1-MAPPING.bam.sort.bed
| +- denPon1-to-triCas1-MAPPING.bam.sort.merge.bed
| +- lepDecl-to-triCas1-MAPPING.bam.bed
| +- lepDecl-to-triCas1-MAPPING.bam.sort.bed
| +- lepDecl-to-triCas1-MAPPING.bam.sort.merge.bed
| +- ontTaul-to-triCas1-MAPPING.bam.bed
| +- ontTaul-to-triCas1-MAPPING.bam.sort.bed
| +- ontTaul-to-triCas1-MAPPING.bam.sort.merge.bed
+- genomes (collapsed)
```

To get some idea of the total number of merged, putatively conserved regions in each exemplar taxon that are shared with the base genome, we can simply loop over the files and count the number of lines in each:

```
> cd uce-coleoptera
> cd bed

> for i in *.bam.sort.merge.bed; do wc -l $i; done
19810 agrPlal-to-triCas1-MAPPING.bam.sort.merge.bed
48350 anoGlal-to-triCas1-MAPPING.bam.sort.merge.bed
21390 denPon1-to-triCas1-MAPPING.bam.sort.merge.bed
33144 lepDecl-to-triCas1-MAPPING.bam.sort.merge.bed
25188 ontTaul-to-triCas1-MAPPING.bam.sort.merge.bed
```

## Remove repetitive intervals

At this point, we've mapped reads to the base genome, kept those regions where reads mapped, converted that to a BED, and merged intervals that are very close to one another.

What we have not done is remove any putatively conserved intervals shared between exemplar taxa and the base genome that are repetitive regions. To do this, we're going to run a python program over all of the BED files for each exemplar taxon. This program will look at the intervals shared between the exemplar taxon and the base genome and it will remove intervals where the base genome is shorter than 80 bp and where > 25 % of the base genome is masked.

```
> cd uce-coleoptera
> cd bed

> for i in *.sort.merge.bed;
do
    phyluce_probe_strip_masked_loci_from_set \
        --bed $i \
        --twobit ../genomes/triCas1/triCas1.2bit \
        --output ${i%.*}.strip.bed \
        --filter-mask 0.25 \
        --min-length 80
done;
```



```

Screened 19810 sequences from agrPlal-to-triCas1-MAPPING.bam.sort.merge.bed.
↳Filtered 3113 with > 25.0% masked bases or > 0 N-bases or < 80 length. Kept 16697.
Screened 48350 sequences from anoGlal-to-triCas1-MAPPING.bam.sort.merge.bed.
↳Filtered 13226 with > 25.0% masked bases or > 0 N-bases or < 80 length. Kept 35124.
Screened 21390 sequences from denPon1-to-triCas1-MAPPING.bam.sort.merge.bed.
↳Filtered 3008 with > 25.0% masked bases or > 0 N-bases or < 80 length. Kept 18382.
Screened 33144 sequences from lepDecl-to-triCas1-MAPPING.bam.sort.merge.bed.
↳Filtered 6585 with > 25.0% masked bases or > 0 N-bases or < 80 length. Kept 26559.
Screened 9637 sequences from menMoll-to-triCas1-MAPPING.bam.sort.merge.bed. Filtered
↳4379 with > 25.0% masked bases or > 0 N-bases or < 80 length. Kept 5258.
Screened 25188 sequences from ontTaul-to-triCas1-MAPPING.bam.sort.merge.bed.
↳Filtered 6505 with > 25.0% masked bases or > 0 N-bases or < 80 length. Kept 18683.

```

When this finishes, your directory structure should look like:

```

uce-coleoptera
+- alignments
| +- all (collapsed)
+- bed
| +- agrPlal-to-triCas1-MAPPING.bam.bed
| +- agrPlal-to-triCas1-MAPPING.bam.sort.bed
| +- agrPlal-to-triCas1-MAPPING.bam.sort.merge.bed
| +- agrPlal-to-triCas1-MAPPING.bam.sort.merge.strip.bed
| +- anoGlal-to-triCas1-MAPPING.bam.bed
| +- anoGlal-to-triCas1-MAPPING.bam.sort.bed
| +- anoGlal-to-triCas1-MAPPING.bam.sort.merge.bed
| +- anoGlal-to-triCas1-MAPPING.bam.sort.merge.strip.bed
| +- denPon1-to-triCas1-MAPPING.bam.bed
| +- denPon1-to-triCas1-MAPPING.bam.sort.bed
| +- denPon1-to-triCas1-MAPPING.bam.sort.merge.bed
| +- denPon1-to-triCas1-MAPPING.bam.sort.merge.strip.bed
| +- lepDecl-to-triCas1-MAPPING.bam.bed
| +- lepDecl-to-triCas1-MAPPING.bam.sort.bed
| +- lepDecl-to-triCas1-MAPPING.bam.sort.merge.bed
| +- lepDecl-to-triCas1-MAPPING.bam.sort.merge.strip.bed
| +- menMoll-to-triCas1-MAPPING.bam.bed
| +- menMoll-to-triCas1-MAPPING.bam.sort.bed
| +- menMoll-to-triCas1-MAPPING.bam.sort.merge.bed
| +- menMoll-to-triCas1-MAPPING.bam.sort.merge.strip.bed
| +- ontTaul-to-triCas1-MAPPING.bam.bed
| +- ontTaul-to-triCas1-MAPPING.bam.sort.bed
| +- ontTaul-to-triCas1-MAPPING.bam.sort.merge.bed
| +- ontTaul-to-triCas1-MAPPING.bam.sort.merge.strip.bed
+- genomes (collapsed)

```

### Determining locus presence in multiple genomes

Up to this point, we've been processing each file on a taxon-by-taxon basis, where each taxon had data aligned to the base genome. Now, we need to determine which loci are conserved *across* taxa. To do that, we first need to prepare a configuration file (named *bed-files.conf*) that gives the paths to each of our *\*.bam.sort.merge.strip.bed* files. That file needs to be in configuration file format, like so:

```

[beds]
agrPlal:agrPlal-to-triCas1-MAPPING.bam.sort.merge.strip.bed
anoGlal:anoGlal-to-triCas1-MAPPING.bam.sort.merge.strip.bed

```

```
denPon1:denPon1-to-triCas1-MAPPING.bam.sort.merge.strip.bed
lepDecl:lepDecl-to-triCas1-MAPPING.bam.sort.merge.strip.bed
ontTaul:ontTaul-to-triCas1-MAPPING.bam.sort.merge.strip.bed
```

The *[beds]* line is the “header” line, and that is followed by each taxon name (on the left) and the name of the BED file we want to process (on the right). You should place this file in the *bed* directory. If you place it elsewhere, you’ll need to use full paths on the right hand side.

Your directory structure should now look like (note new *bed-files.conf*)

```
uce-coleoptera
+- alignments
| +- all (collapsed)
+- bed
| +- agrPlal-to-triCas1-MAPPING.bam.bed
| +- agrPlal-to-triCas1-MAPPING.bam.sort.bed
| +- agrPlal-to-triCas1-MAPPING.bam.sort.merge.bed
| +- agrPlal-to-triCas1-MAPPING.bam.sort.merge.strip.bed
| +- anogla1-to-triCas1-MAPPING.bam.bed
| +- anogla1-to-triCas1-MAPPING.bam.sort.bed
| +- anogla1-to-triCas1-MAPPING.bam.sort.merge.bed
| +- anogla1-to-triCas1-MAPPING.bam.sort.merge.strip.bed
| +- bed-files.conf
| +- denPon1-to-triCas1-MAPPING.bam.bed
| +- denPon1-to-triCas1-MAPPING.bam.sort.bed
| +- denPon1-to-triCas1-MAPPING.bam.sort.merge.bed
| +- denPon1-to-triCas1-MAPPING.bam.sort.merge.strip.bed
| +- lepDecl-to-triCas1-MAPPING.bam.bed
| +- lepDecl-to-triCas1-MAPPING.bam.sort.bed
| +- lepDecl-to-triCas1-MAPPING.bam.sort.merge.bed
| +- lepDecl-to-triCas1-MAPPING.bam.sort.merge.strip.bed
| +- menMoll-to-triCas1-MAPPING.bam.bed
| +- menMoll-to-triCas1-MAPPING.bam.sort.bed
| +- menMoll-to-triCas1-MAPPING.bam.sort.merge.bed
| +- menMoll-to-triCas1-MAPPING.bam.sort.merge.strip.bed
| +- ontTaul-to-triCas1-MAPPING.bam.bed
| +- ontTaul-to-triCas1-MAPPING.bam.sort.bed
| +- ontTaul-to-triCas1-MAPPING.bam.sort.merge.bed
| +- ontTaul-to-triCas1-MAPPING.bam.sort.merge.strip.bed
+- genomes (collapsed)
```

Now, we’re going to run the following program, that creates a record of which alignment intervals are shared among taxa. We need to pass the location of the *bed-files.conf* to this program, along with the name of our base genome, and a name for the output database that will be created:

```
> phyluce_probe_get_multi_merge_table \
  --conf bed-files.conf \
  --base-taxon triCas1 \
  --output coleoptera-to-triCas1.sqlite

agrplal.....
anogla1.....
denpon1.....
lepdecl.....
onttaul.....
Creating database
Inserting results
```

The program shows the results of inserting data for each exemplar taxon that we've selected. If we take a look at the table contents (see *The probe.matches.sqlite database* for more instructions on *sqlite* databases), we see something like the following:

```
sqlite> select * from triCas1 limit 10;
```

uce	chromo	start	stop	agrpla1	anogla1	denpon1	
↪lepdec1	onttau1						↪
↪-----	-----	-----	-----	-----	-----	-----	↪
1	GG695505.1	532	632	1	1	0	0↪
↪	1						
2	GG695547.1	826	926	0	1	0	0↪
↪	0						
3	GG695547.1	1121	1221	0	1	0	0↪
↪	0						
4	GG695547.1	1293	1393	0	1	0	0↪
↪	0						
5	GG694821.1	1002	1102	0	0	0	1↪
↪	0						
6	GG695519.1	73	193	0	1	1	0↪
↪	0						
7	GG695519.1	222	380	1	0	1	0↪
↪	1						
8	GG695519.1	925	1129	1	1	1	0↪
↪	1						
9	DS497688.1	17907	18022	0	0	1	0↪
↪	0						
10	DS497688.1	19840	19934	0	1	0	0↪
↪	0						

The first row of this table (which is limited to 10 rows of results by the query although it is 60699 rows long) shows that for *triCas1* contig *GG695505.1*, *agrpla1*, *anogla1*, and *onttau1* have reads that overlap at position 532 to 632.

### Determining shared, conserved, loci

Now that we have our table of results, we can run a quick query (using a *Python* program) against the table to look at results, more generally. The following code queries the database and writes out the number of loci shared by the base taxon (*triCas1*) and 1, 2, 3, 4, and 5 (all) of the exemplar taxa that we've aligned to the base genome. You need to give the program the path to the database created above and the name of the base taxon:

```
> python phyluce_probe_query_multi_merge_table \
    --db coleoptera-to-triCas1.sqlite \
    --base-taxon triCas1
```

```
Loci shared by triCas1 + 0 taxa: 60,699.0
Loci shared by triCas1 + 1 taxa: 60,699.0
Loci shared by triCas1 + 2 taxa: 32,431.0
Loci shared by triCas1 + 3 taxa: 15,834.0
Loci shared by triCas1 + 4 taxa: 6,471.0
Loci shared by triCas1 + 5 taxa: 1,822.0
```

The output from this program basically says that, if we are interested in only those loci found in all exemplar taxa that align to the base genome, there are 1,822 of those. Similarly, if we're willing to be a little less strict about things, there are 6,471 conserved loci that are shared by *triCas1* and 4 of the exemplar taxa.

**Question: How conservative should I be?**

Basically, the question boils down to “Should I select only the set of loci shared by all exemplars and the base genome or should I be more liberal?”. It’s also a hard question to answer. In most cases, I’m pretty happy selecting  $n-1$  or  $n-2$  where  $n$  is the total number of exemplar taxa. In the example below, however, we’ve selected  $n$  as the “ideal”. This is largely because we have so little information about coleopteran genomes - so we want to be pretty darn sure these loci are found in most/all of them.

---

Now that we have a general sense of the number of conserved loci in each class of sharing across exemplars (e.g. 5 (all), 4, 3, 2, 1), we need to extract those loci that fall within one of these classes. In this case (and as noted in the box, above), we’re going to output only those conserved loci that we’ve identified as being shared between the base genome and *all* exemplars. We do that with a slightly different Python script. This script takes the database name, the base genome, the count of exemplar taxa shared across, and the name of an output file as input. The output file will be BED formatted.

```
> phyluce_probe_query_multi_merge_table \  
    --db coleoptera-to-triCas1.sqlite \  
    --base-taxon triCas1 \  
    --output triCas1+5.bed \  
    --specific-counts 5  
  
Counter({'anogla1': 1822, 'lepdeci1': 1822, 'agrpla1': 1822, 'denpon1': 1822, 'onttaul  
↳ ': 1822})
```

## Conserved locus validation

### Extract FASTA sequence from base genome for temp bait design

Now that we’ve identified conserved sequences shared among the base genome and the exemplar taxa, we need to start designing baits to capture these loci. The first step in this process is to extract FASTA sequences from the base genome that correspond to the loci we’ve identified. We do that with a Python script that takes as input the BED file we created, above, the 2bit-formatted base genomes, a length of sequence we want to extract (160 bp), and an output FASTA filename.

---

#### Question: Why buffer to 160 bp?

We are extracting FASTA regions of 160 bp because that allows us to place 2 baits right in the center of this region at 3x tiling density which means that standard 120 bp baits will overlap by 40 bp and have 80 bp to each side (total length 160 bp).

---

To run the code, we use:

```
> phyluce_probe_get_genome_sequences_from_bed \  
    --bed triCas1+5.bed \  
    --twobit ../genomes/triCas1/triCas1.2bit \  
    --buffer-to 160 \  
    --output triCas1+5.fasta  
  
Screened 1822 sequences. Filtered 7 < 160 bp or with > 25.0% masked bases or > 0 N-  
↳ bases. Kept 1815.
```

That should produce a fasta file whose contents look similar to:

```
>slice_0 |DS497688.1:249724-250111  
AAAATCAAAGTCGAATACAAAGCGAATCTAAGACTTTCTATCCTGAAGAGATCAGTTCC
```



```
> phyluce_probe_easy_lastz \  
  --target triCas1+5.temp.probes \  
  --query triCas1+5.temp.probes \  
  --identity 50 --coverage 50 \  
  --output triCas1+5.temp.probes-TO-SELF-PROBES.lastz  
  
Started:  Fri Jun 03, 2016  13:57:54  
Ended:    Fri Jun 03, 2016  13:57:55  
Time for execution:  0.0284410158793 minutes
```

Now that we've run the alignments, we need to screen them alignments and remove the duplicate baits from the bait set. This program takes as input the `lastz` results from above and the temp-probe file, as well as the probe-prefix that we used during probe design, above. The results are written to a file that is equivalent to the probe file name + *DUPE-SCREENED*, so in this case the output file is named *triCas1+5.temp-DUPE-SCREENED.probes*.

```
> phyluce_probe_remove_duplicate_hits_from_probes_using_lastz \  
  --fasta triCas1+5.temp.probes \  
  --lastz triCas1+5.temp.probes-TO-SELF-PROBES.lastz \  
  --probe-prefix=uce-  
  
Parsing lastz file...  
Screening results...  
Screened 3601 fasta sequences.  Filtered 292 duplicates. Kept 3019.
```

### Align baits against exemplar genomes

Now that we have a duplicate-free (or putatively duplicate free) set of temporary baits designed from conserved loci in the base genome, we're going to use some in-silico alignments to see if we can locate these loci in the several exemplar genomes.

**Attention:** For the following analyses, you need genome assemblies for each of the exemplar taxa, formatted as *2bit* files.

We'll use the results of these alignments to design a bait set that includes baits designed from the base genome, but also from the exemplar taxa. This should allow our bait set to work more consistently across broad groups of organisms.

In terms of directory structure, things should look pretty similar to the following:

```
uce-coleoptera  
+- alignments  
| +- all (collapsed)  
+- bed  
| +- agrPlal-to-triCas1-MAPPING.bam.bed  
| +- agrPlal-to-triCas1-MAPPING.bam.sort.bed  
| +- agrPlal-to-triCas1-MAPPING.bam.sort.merge.bed  
| +- agrPlal-to-triCas1-MAPPING.bam.sort.merge.strip.bed  
| +- anoGlal-to-triCas1-MAPPING.bam.bed  
| +- anoGlal-to-triCas1-MAPPING.bam.sort.bed  
| +- anoGlal-to-triCas1-MAPPING.bam.sort.merge.bed  
| +- anoGlal-to-triCas1-MAPPING.bam.sort.merge.strip.bed  
| +- bed-files.conf  
| +- coleoptera-to-triCas1.sqlite  
| +- denPon1-to-triCas1-MAPPING.bam.bed  
| +- denPon1-to-triCas1-MAPPING.bam.sort.bed
```

```

| +- denPon1-to-triCas1-MAPPING.bam.sort.merge.bed
| +- denPon1-to-triCas1-MAPPING.bam.sort.merge.strip.bed
| +- lepDecl-to-triCas1-MAPPING.bam.bed
| +- lepDecl-to-triCas1-MAPPING.bam.sort.bed
| +- lepDecl-to-triCas1-MAPPING.bam.sort.merge.bed
| +- lepDecl-to-triCas1-MAPPING.bam.sort.merge.strip.bed
| +- menMoll-to-triCas1-MAPPING.bam.bed
| +- menMoll-to-triCas1-MAPPING.bam.sort.bed
| +- menMoll-to-triCas1-MAPPING.bam.sort.merge.bed
| +- menMoll-to-triCas1-MAPPING.bam.sort.merge.strip.bed
| +- ontTaul-to-triCas1-MAPPING.bam.bed
| +- ontTaul-to-triCas1-MAPPING.bam.sort.bed
| +- ontTaul-to-triCas1-MAPPING.bam.sort.merge.bed
| +- ontTaul-to-triCas1-MAPPING.bam.sort.merge.strip.bed
| +- triCas1+5.bed
| +- triCas1+5.bed.missing.matrix
| +- triCas1+5.fasta
+- genomes
  +- agrPla1
  | +- agrPla1.2bit
  | +- agrPla1.fasta
  +- anoGla1
  | +- anoGla1.2bit
  | +- anoGla1.fasta
  +- denPon1
  | +- denPon1.2bit
  | +- denPon1.fasta
  +- lepDecl
  | +- lepDecl.2bit
  | +- lepDecl.fasta
  +- menMoll
  | +- menMoll.2bit
  | +- menMoll.fasta
  +- ontTaul
  | +- ontTaul.2bit
  | +- ontTaul.fasta
  +- triCas1
  | +- triCas1.2bit
  | +- triCas1.fasta

```

Note that we have all the genomes in their directory, in both FASTA and *2bit* formats. We're also have a new genome sequence in here - that of *menMoll* (*Mengenilla moldrzyki* [twisted-wing parasites]), which represents the outgroup to Coleoptera. We're adding this taxon because it helps us bridge the base of the tree - e.g. the divergence between the outgroup and the exemplar taxa that we're using to design probes.

---

#### Question: What exemplar taxa should I use for bait design?

This is a really hard question to answer. In old, divergent groups with few genomic resources, the answer is usually "all the species" with genomic data. Basically, you want to include exemplars that make the divergence among baits targeting the same loci something >20% or so. That said, even this number is a bit of a guess - no one has systematically tested how "sticky" baits are when they are used to enrich loci across divergent groups. We know they are pretty sticky and in certain cases can enrich loci as much as 35%-40% divergent from the bait sequence. Generally speaking, I try to include exemplar taxa during probe design that bridge the known diversity of a given group... again, in many cases this is hard (or impossible) to know given current data. So, you may have to take a bit of a guess.

---

So, assuming that you have the appropriate *2bit* files in *uce-coleoptera/genomes*, we are going to align the temporary

probes that we've designed to the exemplar genomes, and we're going to run these and subsequent bait design steps in a new directory, named *probe-design*. So:

```
> cd uce-coleoptera
> mkdir probe-design
```

Now, you're directory structure should look like:

```
uce-coleoptera
+- alignments
| +- all (collapsed)
+- bed
| +- agrPla1-to-triCas1-MAPPING.bam.bed
| +- agrPla1-to-triCas1-MAPPING.bam.sort.bed
| +- agrPla1-to-triCas1-MAPPING.bam.sort.merge.bed
| +- agrPla1-to-triCas1-MAPPING.bam.sort.merge.strip.bed
| +- anoGla1-to-triCas1-MAPPING.bam.bed
| +- anoGla1-to-triCas1-MAPPING.bam.sort.bed
| +- anoGla1-to-triCas1-MAPPING.bam.sort.merge.bed
| +- anoGla1-to-triCas1-MAPPING.bam.sort.merge.strip.bed
| +- bed-files.conf
| +- coleoptera-to-triCas1.sqlite
| +- denPon1-to-triCas1-MAPPING.bam.bed
| +- denPon1-to-triCas1-MAPPING.bam.sort.bed
| +- denPon1-to-triCas1-MAPPING.bam.sort.merge.bed
| +- denPon1-to-triCas1-MAPPING.bam.sort.merge.strip.bed
| +- lepDecl-to-triCas1-MAPPING.bam.bed
| +- lepDecl-to-triCas1-MAPPING.bam.sort.bed
| +- lepDecl-to-triCas1-MAPPING.bam.sort.merge.bed
| +- lepDecl-to-triCas1-MAPPING.bam.sort.merge.strip.bed
| +- menMoll-to-triCas1-MAPPING.bam.bed
| +- menMoll-to-triCas1-MAPPING.bam.sort.bed
| +- menMoll-to-triCas1-MAPPING.bam.sort.merge.bed
| +- menMoll-to-triCas1-MAPPING.bam.sort.merge.strip.bed
| +- ontTaul-to-triCas1-MAPPING.bam.bed
| +- ontTaul-to-triCas1-MAPPING.bam.sort.bed
| +- ontTaul-to-triCas1-MAPPING.bam.sort.merge.bed
| +- ontTaul-to-triCas1-MAPPING.bam.sort.merge.strip.bed
| +- triCas1+5.bed
| +- triCas1+5.bed.missing.matrix
| +- triCas1+5.fasta
+- genomes (collapsed)
+- probe-design
```

We need to align the temporary probe sequences to each genome, which we can do using the following code, which takes as input our temporary probe file, the list of genomes we want to align the probes against, the path to the genomes, the minimum sequence identity to accept a match (on the low end of the spectrum for this step), and number of compute cores to use, and the name of an output database to create and the output directory in which to store the lastz results.

**Warning:** Note that I am using 16 physical CPU cores (*-cores*) to do this work. You need to use the number of physical cores available on *your* machine.

```
> phyluce_probe_run_multiple_lastzs_sqlite \
  --probefile ../bed/triCas1+5.temp-DUPE-SCREENED.probes \
  --scaffoldlist agrPla1 anoGla1 denPon1 lepDecl1 ontTaul1 triCas1 menMoll1 \
```



```

--genome-base-path ../genomes \
--identity 50 \
--cores 16 \
--db triCas1+5+menMoll.sqlite \
--output coleoptera-genome-lastz

Running against agrPla1.2bit
Running with the --huge option. Chunking files into 10000000 bp

< ... snip ... >

Cleaning up the chunked files...
Cleaning /nfs/data1/working/bfaircloth-insects/coleoptera/temp/probe-design/
↪coleoptera-genome-lastz/triCas1+5.temp-DUPE-SCREENED.probes_v_menMoll.lastz
Creating menMoll table
Inserting data to menMoll table

```

### Extract sequence around conserved loci from exemplar genomes

Based on the alignments of the temporary probe set to the exemplar genomes, we need to extract FASTA data from each of the exemplar sequences so that we can design baits targeting the conserved loci in each. This is pretty similar to what we did earlier for the temporary probe set, except that now we're running the extraction across all the exemplar taxa.

Before we begin, we need to make a configuration file with all the genome locations in it (again, as before):

```

[scaffolds]
menMoll:/path/to/uce-coleoptera/genomes/menMoll/menMoll.2bit
agrPla1:/path/to/uce-coleoptera/genomes/agrPla1/agrPla1.2bit
anoGla1:/path/to/uce-coleoptera/genomes/anoGla1/anoGla1.2bit
denPon1:/path/to/uce-coleoptera/genomes/denPon1/denPon1.2bit
lepDecl:/path/to/uce-coleoptera/genomes/lepDecl/lepDecl.2bit
ontTau1:/path/to/uce-coleoptera/genomes/ontTau1/ontTau1.2bit
triCas1:/path/to/uce-coleoptera/genomes/triCas1/triCas1.2bit

```

Using the configuration file, we need to extract the FASTA sequence that we need from each exemplar taxon. Here, we're buffering each locus to 180 bp to give us a little more room to work with during the probe design step. The program takes our config file as input, along with the folder of `lastz` results created above. The `--name-pattern` argument allows us to match files into the `-lastz` directory, `-probes` is how we buffer the sequence, and we pass the name of an output directory to `-output`:

```

> phyluce_probe_slice_sequence_from_genomes \
  --conf coleoptera-genome.conf \
  --lastz coleoptera-genome-lastz \
  --probes 180 \
  --name-pattern "triCas1+5.temp-DUPE-SCREENED.probes_v_{}.lastz.clean" \
  --output coleoptera-genome-fasta

2016-06-03 15:07:16,642 - Phyluce - INFO - ===== Starting Phyluce:↵
↪Slice Sequence =====
2016-06-03 15:07:16,644 - Phyluce - INFO - ----- Working on menMoll_↵
↪genome -----
2016-06-03 15:07:16,645 - Phyluce - INFO - Reading menMoll genome
2016-06-03 15:07:20,221 - Phyluce - INFO - menMoll: 884 uces, 139 dupes, 745 non-↵
↪dupes, 0 orient drop, 2 length drop, 738 written
2016-06-03 15:07:20,222 - Phyluce - INFO - ----- Working on agrPla1_↵
↪genome -----

```

```

2016-06-03 15:07:20,223 - Phyluce - INFO - Reading agrPla1 genome
2016-06-03 15:07:24,759 - Phyluce - INFO - agrPla1: 1410 uces, 184 dupes, 1226 non-
↳dupes, 7 orient drop, 63 length drop, 1156 written
2016-06-03 15:07:24,760 - Phyluce - INFO - ----- Working on anoGla1_
↳genome -----
2016-06-03 15:07:24,761 - Phyluce - INFO - Reading anoGla1 genome
2016-06-03 15:07:29,926 - Phyluce - INFO - anoGla1: 1474 uces, 224 dupes, 1250 non-
↳dupes, 6 orient drop, 35 length drop, 1209 written
2016-06-03 15:07:29,926 - Phyluce - INFO - ----- Working on denPon1_
↳genome -----
2016-06-03 15:07:29,929 - Phyluce - INFO - Reading denPon1 genome
2016-06-03 15:07:34,472 - Phyluce - INFO - denPon1: 1361 uces, 305 dupes, 1056 non-
↳dupes, 6 orient drop, 30 length drop, 1020 written
2016-06-03 15:07:34,472 - Phyluce - INFO - ----- Working on lepDecl_
↳genome -----
2016-06-03 15:07:34,473 - Phyluce - INFO - Reading lepDecl genome
2016-06-03 15:07:40,020 - Phyluce - INFO - lepDecl: 1436 uces, 259 dupes, 1177 non-
↳dupes, 10 orient drop, 28 length drop, 1139 written
2016-06-03 15:07:40,021 - Phyluce - INFO - ----- Working on ontTaul_
↳genome -----
2016-06-03 15:07:40,022 - Phyluce - INFO - Reading ontTaul genome
2016-06-03 15:07:44,350 - Phyluce - INFO - ontTaul: 1361 uces, 206 dupes, 1155 non-
↳dupes, 9 orient drop, 41 length drop, 1105 written
2016-06-03 15:07:44,350 - Phyluce - INFO - ----- Working on triCas1_
↳genome -----
2016-06-03 15:07:44,351 - Phyluce - INFO - Reading triCas1 genome
2016-06-03 15:07:49,499 - Phyluce - INFO - triCas1: 1513 uces, 199 dupes, 1314 non-
↳dupes, 26 orient drop, 46 length drop, 1242 written

```

If we look at our directory structure, it looks something like:

```

uce-coleoptera
+- alignments (collapsed)
+- bed (collapsed)
+- genomes (collapsed)
+- probe-design
  +- coleoptera-genome.conf
  +- coleoptera-genome-fasta
  | +- agrpla1.fasta
  | +- anogla1.fasta
  | +- denpon1.fasta
  | +- lepdecl.fasta
  | +- menmoll.fasta
  | +- onttaul.fasta
  | +- tricasl.fasta
  +- coleoptera-genome-lastz
  | +- triCas1+5.temp-DUPE-SCREENED.probes_v_agrPla1.lastz.clean
  | +- triCas1+5.temp-DUPE-SCREENED.probes_v_anogla1.lastz.clean
  | +- triCas1+5.temp-DUPE-SCREENED.probes_v_denPon1.lastz.clean
  | +- triCas1+5.temp-DUPE-SCREENED.probes_v_lepDecl.lastz.clean
  | +- triCas1+5.temp-DUPE-SCREENED.probes_v_menMoll.lastz.clean
  | +- triCas1+5.temp-DUPE-SCREENED.probes_v_ontTaul.lastz.clean
  | +- triCas1+5.temp-DUPE-SCREENED.probes_v_tricasl.lastz.clean
  +- triCas1+5+menMoll.sqlite

```

The FASTA files we just created are in *uce-coleoptera/coleoptera-genome-fasta*. The output from the program that you see basically shows you how many UCE loci we extracted from each of the exemplar genomes. As expected, the lowest number we located and extracted are from the *menMoll* (outgroup) genome.

If we have a look in one of these FASTA files, it looks like:

```
> less probe-design/coleoptera-genome-fasta/agrplal.fasta

>slice_0|contig:KL218988.1|slice:301686-301866|uce:uce-500|match:301721-
↳301831|orient:+|probes:1
TCGAACTTCTGGTGCTTGTCCACCCTTGATGTCCGCACCGAATTCCTTCACGAGACTGTTT
CTAAAACCTTTGGACAAGATGATTGGAGACACAGAAACAGACGAACATCAGAAAACGTGTAT
ATCTTCACTTTCTAGAGCATTACATACAAACCTTATTACCAGATGTACTCAGCAGCAGCTTT
>slice_1|contig:KL218988.1|slice:319624-319804|uce:uce-501|match:319637-
↳319791|orient:+|probes:2
atgattttttcaaAGGTTACAGCGAAGTCCTCGATTCTACAGCAGATCGAAGAACTAGGA
GAAGAGACTGGCCTGGTGTGCTGTATTTGTCGCGAGGGATACAAGTATCAACCTGCCAAG
GTATTGGGAATTTATACGTTTACAAAGAGGTGCAACGTGGACGAGTTCGAAGCAAACCA
>slice_2|contig:KL219144.1|slice:184423-184603|uce:uce-503|match:184453-
↳184573|orient:+|probes:1
agttttaataatcttACCTAAAGAACTAAAATGAAGAAGCATTTCGTCTGCTCGTAAGT
CTTGACAAATGACATCATCAGCGTAGACACATATTATAGCAAGGCATATAAATAGGTGAA
AGTAGTCTGTTAGATAATTTGCCCAACAAGCTTCCCACAGTCTAAGGGCAACACCTTCGG
>slice_3|contig:KL219144.1|slice:237138-237318|uce:uce-504|match:237150-
↳237306|orient:+|probes:2
CTGTGCAAGAGTGAGTGCCATTGATGCAACACTTGAGCGAGATGATCTAAACCTCCATGG
TGAAAATGAAGAATTTTATATTGAGATTCCCTCGAAGCAACAACCACCTGCCCTGATGTG
CAGCTTGAGTCGTTAAAGAAAAGCCTTAAAGATCTCATTGGCTTAGATCAACGCTGAAC
```

### Find which loci we detect consistently

As before, we want to determine which loci we are detecting consistently across all of the exemplar taxa when doing these in-silico searches. To do that, we'll run another bit of python code. Here, we're working in the *uce-coleoptera/coleoptera-genome-lastz* directory. This program will create a relational database that houses detections of loci in the exemplar taxa. It takes, as input, the folder of FASTAs we just created, the base genome taxon, and a name to use as the output database:

```
> phyluce_probe_get_multi_fasta_table \
  --fastas coleoptera-genome-fasta \
  --output multifastas.sqlite \
  --base-taxon triCas1

menmoll.
agrplal..
anogla1..
denpon1..
lepdec1..
onttau1..
tricas1..
Creating database
Inserting results
```

If we take a look at the table contents in the database (see *The probe.matches.sqlite database* for more instructions on *sqlite* databases), we see something like the following:

locus	menmoll	agrplal	anogla1	denpon1	lepdec1	onttau1	↳
↳tricas1							
↳-----							
uce-500	0	1	1	0	1	1	1
uce-501	1	1	1	1	1	1	1

uce-503	1	1	1	1	0	1	1
uce-504	1	1	1	1	1	1	1
uce-505	1	0	1	0	0	1	1
uce-506	1	1	1	1	1	1	1
uce-507	0	1	1	1	1	1	1
uce-508	1	1	1	1	1	0	1
uce-509	1	0	1	1	1	1	1
uce-967	0	0	0	1	1	1	0

Which shows our detection of conserved loci in each of the exemplar taxa when we search for them using the temporary probes that we designed from the base genome. As before, we can get some idea of the distribution of hits among exemplar taxa (e.g., are loci detected in “all”, n-1 taxa, n-2 taxa, etc.).

```
> phyluce_probe_query_multi_fasta_table \
  --db multifastas.sqlite \
  --base-taxon triCas1

Loci shared by 0 taxa: 1,437.0
Loci shared by 1 taxa: 1,437.0
Loci shared by 2 taxa: 1,355.0
Loci shared by 3 taxa: 1,303.0
Loci shared by 4 taxa: 1,209.0
Loci shared by 5 taxa: 1,099.0
Loci shared by 6 taxa: 820.0
Loci shared by 7 taxa: 386.0
```

Again, we’ve got to make a decision here about how conservative we want to be regarding baits that hit all/some taxa. We only get 386 loci that we detect in all exemplars (including the base genome and the *menMoll* outgroup). That seems too strict (particularly because this total includes *menMoll*, which is really divergent from our taxa of interest). We also have to keep in mind that we can randomly fail to detect loci that are actually present, either by chance or do to sequence divergences that are >50% (the value we used in our search). In the end, I settled on loci we detected in  $\geq 4$  exemplar taxa. So we need to extract those from the database and store them in *triCas1+5-back-to-4.conf*:

```
phyluce_probe_query_multi_fasta_table \
  --db multifastas.sqlite \
  --base-taxon triCas1 \
  --output triCas1+5-back-to-4.conf \
  --specific-counts 4

Counter({'trिकास1': 1160, 'anoglal': 1140, 'agrplal': 1091, 'lepdec1': 1080, 'onttaul
↪': 1043, 'denpon1': 969, 'menmoll': 658})
Total loci = 1209
```

The values above show the number of loci detected in each exemplar taxon and the total number of loci we’ll be targeting with the bait set we’re about to design.

Your directory should look something like the following:

```
uce-coleoptera
+- alignments (collapsed)
+- bed (collapsed)
+- genomes (collapsed)
+- probe-design
  +- coleoptera-genome.conf
  +- coleoptera-genome-fasta
  | +- agrplal.fasta
  | +- anoglal.fasta
  | +- denpon1.fasta
```



## Remove duplicates from our bait set

As before, we need to check out bait set for duplicate loci. This time, the search is going to take longer, because of the larger number of baits. We'll align all the probes to themselves, then read in the alignments, and filter the probe list to remove putative duplicates.

Align probes to themselves at low stringency to identify duplicates:

```
phyluce_probe_easy_lastz \
  --target coleoptera-v1-master-probe-list.fasta \
  --query coleoptera-v1-master-probe-list.fasta \
  --identity 50 \
  --coverage 50 \
  --output coleoptera-v1-master-probe-list-TO-SELF-PROBES.lastz
```

```
Started: Fri Jun 03, 2016 15:50:52
Ended: Fri Jun 03, 2016 15:51:11
Time for execution: 0.322272149722 minutes
```

Now, screen the alignments and filter our master probe list to remove duplicates:

```
phyluce_probe_remove_duplicate_hits_from_probes_using_lastz \
  --fasta coleoptera-v1-master-probe-list.fasta \
  --lastz coleoptera-v1-master-probe-list-TO-SELF-PROBES.lastz \
  --probe-prefix=uce-
```

```
Parsing lastz file...
Screening results...
Screened 14112 fasta sequences. Filtered 37 duplicates. Kept 13674.
```

The master probe list that has been filtered of putatively duplicate loci is now located in *coleoptera-v1-master-probe-list-DUPE-SCREENED.fasta*.

Your directory should look something like the following:

```
uce-coleoptera
+- alignments (collapsed)
+- bed (collapsed)
+- genomes (collapsed)
+- probe-design
  +- coleoptera-genome.conf
  +- coleoptera-genome-fasta
  | +- agrpl1a1.fasta
  | +- anogl1a1.fasta
  | +- denpon1.fasta
  | +- lepdec1.fasta
  | +- menmoll1.fasta
  | +- onttau1.fasta
  | +- tricas1.fasta
+- coleoptera-genome-lastz
  | +- triCas1+5.temp-DUPE-SCREENED.probes_v_agrPl1a1.lastz.clean
  | +- triCas1+5.temp-DUPE-SCREENED.probes_v_anoGl1a1.lastz.clean
  | +- triCas1+5.temp-DUPE-SCREENED.probes_v_denPon1.lastz.clean
  | +- triCas1+5.temp-DUPE-SCREENED.probes_v_lepDec1.lastz.clean
  | +- triCas1+5.temp-DUPE-SCREENED.probes_v_menMoll1.lastz.clean
  | +- triCas1+5.temp-DUPE-SCREENED.probes_v_ontTau1.lastz.clean
  | +- triCas1+5.temp-DUPE-SCREENED.probes_v_triCas1.lastz.clean
+- coleoptera-v1-master-probe-list-DUPE-SCREENED.fasta
+- coleoptera-v1-master-probe-list.fasta
```

```

+- coleoptera-v1-master-probe-list-TO-SELF-PROBES.lastz
+- multifastas.sqlite
+- triCas1+5-back-to-4.conf
+- triCas1+5-back-to-4.conf.missing.matrix
+- triCas1+5+menMol1.sqlite

```

## The master bait list

What we've created, above, is the master bait list that contains baits targeting the conserved locus we identified. Because we've designed probes from multiple exemplar taxa, the number of overall baits is high (as high as 14 baits targeting each conserved locus). This bait set is ready for synthesis and subsequent enrichment of these conserved loci shared among coleoptera.

## Subsetting the master probe list

Sometimes we might not want to synthesize *all* of the baits for *all* of the loci. For instance, we might be enriching loci from species that are nested within the clade defined by (('Anoplophora glabripennis (Asian longhorned beetle)':'Leptinotarsa decemlineata (Colorado potato beetle)')'Dendroctonus ponderosae (mountain pine beetle)'), and because we're only working with these species, we might want to drop the baits targeting UCE loci in *Agrilus planipennis* (emerald ash borer), *Tribolium castaneum* (red flour beetle), *Onthophagus taurus* (taurus scarab), and *Mengenilla moldrzyki* (Strepsiptera). This is actually pretty easy to do - we just need to subset the baits to include those taxa that we *do* want. Given the example, above, we can run:

```

> phyluce_probe_get_subsets_of_tiled_probes \
  --probes coleoptera-v1-master-probe-list-DUPE-SCREENED.fasta \
  --taxa anogla1 lepdeci1 denpon1 \
  --output coleoptera-v1-master-probe-list-DUPE-SCREENED-SUBSET-CLADE_1.fasta

All probes = 13674
--- Probes by taxon ---
anogla1 2189
menmoll 1236
lepdeci1 2083
denpon1 1867
onttaul 1970
agrpla1 2086
tricas1 2243
--- Post filtering ---
Conserved locus count = 1169
Probe Count = 6139

```

## In-silico test of the bait design

Now that we've designed our baits, it's always good to run a sanity check on the data - if we use the baits to collect data from a selection of available genomes (or other genetic data), can we reconstruct a phylogeny that is sane, given what we know about the specific taxa?

How we do that is outlined below. Many of the steps we've run before, so I'm not going to explain these quite as much as I have previously. Several other of the steps that we're going to run are also outlined in *Tutorial I: UCE Phylogenomics*.

First, we need to make a directory to hold our in-silico test results:

```
> cd uce-coleoptera
> mkdir probe-design-test
> cd probe-design-test
```

Now, our directory tree should look something like:

```
uce-coleoptera
+- alignments (collapsed)
+- bed (collapsed)
+- genomes (collapsed)
+- probe-design (collapsed)
+- probe-design-test
```

### Align our bait set to the extant genome sequences

Here (and if you have them), you may want to include *all* the genomic data you have access to - particularly if you removed some taxa because they were closely related to other taxa and designing probes from these closely related groups was redundant. To run these alignments (you've seen this before):

```
phyluce_probe_run_multiple_lastzs_sqlite \
  --db triCas1+5+strepsiptera-test.sqlite \
  --output coleoptera-genome-lastz \
  --probefile ../probe-design/coleoptera-v1-master-probe-list-DUPE-SCREENED.fasta \
  --scaffoldlist agrPlal anoGlal denPon1 lepDec1 ontTau1 triCas1 menMoll \
  --genome-base-path ../genomes \
  --identity 50 \
  --cores 16
```

**Warning:** Note that I am using 16 physical CPU cores (*-cores*) to do this work. You need to use the number of physical cores available on *your* machine.

Now, we need to extract fasta data for each of these loci. This is effectively the same as what we've done before, but notice the use of *-flank* in place of *-probe*. This tells the program that we want to extract larger chunks of sequence, in this base 400 bp to the each side of a given locus (if possible):

```
> phyluce_probe_slice_sequence_from_genomes \
  --conf coleoptera-genome.conf \
  --lastz coleoptera-genome-lastz \
  --output coleoptera-genome-fasta \
  --flank 400 \
  --name-pattern "coleoptera-v1-master-probe-list-DUPE-SCREENED.fasta_v_{}.lastz.
↳clean"

2016-06-03 16:36:47,712 - Phyluce - INFO - ===== Starting Phyluce:↳
↳Slice Sequence =====
2016-06-03 16:36:47,714 - Phyluce - INFO - ----- Working on menMoll↳
↳genome -----
2016-06-03 16:36:47,715 - Phyluce - INFO - Reading menMoll genome
2016-06-03 16:36:57,447 - Phyluce - INFO - menMoll: 766 uces, 73 dupes, 693 non-dupes,
↳ 0 orient drop, 2 length drop, 691 written
2016-06-03 16:36:57,447 - Phyluce - INFO - ----- Working on agrPlal↳
↳genome -----
2016-06-03 16:36:57,448 - Phyluce - INFO - Reading agrPlal genome
2016-06-03 16:37:12,538 - Phyluce - INFO - agrPlal: 1151 uces, 81 dupes, 1070 non-
↳dupes, 0 orient drop, 34 length drop, 1036 written
```



```

2016-06-03 16:37:12,538 - Phyluce - INFO - ----- Working on anoGla1_
↳genome -----
2016-06-03 16:37:12,539 - Phyluce - INFO - Reading anoGla1 genome
2016-06-03 16:37:29,320 - Phyluce - INFO - anoGla1: 1167 uces, 103 dupes, 1064 non-
↳dupes, 0 orient drop, 17 length drop, 1047 written
2016-06-03 16:37:29,321 - Phyluce - INFO - ----- Working on denPon1_
↳genome -----
2016-06-03 16:37:29,322 - Phyluce - INFO - Reading denPon1 genome
2016-06-03 16:37:44,958 - Phyluce - INFO - denPon1: 1126 uces, 174 dupes, 952 non-
↳dupes, 1 orient drop, 16 length drop, 935 written
2016-06-03 16:37:44,959 - Phyluce - INFO - ----- Working on lepDecl_
↳genome -----
2016-06-03 16:37:44,959 - Phyluce - INFO - Reading lepDecl genome
2016-06-03 16:38:01,794 - Phyluce - INFO - lepDecl: 1156 uces, 142 dupes, 1014 non-
↳dupes, 6 orient drop, 22 length drop, 986 written
2016-06-03 16:38:01,794 - Phyluce - INFO - ----- Working on ontTaul_
↳genome -----
2016-06-03 16:38:01,796 - Phyluce - INFO - Reading ontTaul genome
2016-06-03 16:38:16,611 - Phyluce - INFO - ontTaul: 1134 uces, 100 dupes, 1034 non-
↳dupes, 3 orient drop, 14 length drop, 1017 written
2016-06-03 16:38:16,612 - Phyluce - INFO - ----- Working on triCas1_
↳genome -----
2016-06-03 16:38:16,613 - Phyluce - INFO - Reading triCas1 genome
2016-06-03 16:38:32,786 - Phyluce - INFO - triCas1: 1172 uces, 70 dupes, 1102 non-
↳dupes, 13 orient drop, 13 length drop, 1076 written

```

### Match contigs to baits

In the step above, we essentially extracted FASTA data for each taxon, and wrote those out into individual FASTA files. These are the equivalent of the assembled contigs that we use in the standard [phyluce](#) pipeline, so now, we're going to use that workflow. Note that the filtering in `phyluce_assembly_match_contigs_to_probes` is more strict than what we used above to identify contigs.

```

> phyluce_assembly_match_contigs_to_probes \
  --contigs coleoptera-genome-fasta \
  --probes ../probe-design/coleoptera-v1-master-probe-list-DUPE-SCREENED.fasta \
  --output in-silico-lastz \
  --min_coverage 67 \
  --log-path log

2016-06-03 16:40:36,888 - phyluce_assembly_match_contigs_to_probes - INFO - agrplal:
↳903 (87.16%) uniques of 1036 contigs, 0 dupe probe matches, 116 UCE loci removed
↳for matching multiple contigs, 117 contigs removed for matching multiple UCE loci
2016-06-03 16:41:06,688 - phyluce_assembly_match_contigs_to_probes - INFO - anogla1:
↳927 (88.54%) uniques of 1047 contigs, 0 dupe probe matches, 111 UCE loci removed
↳for matching multiple contigs, 116 contigs removed for matching multiple UCE loci
2016-06-03 16:41:28,524 - phyluce_assembly_match_contigs_to_probes - INFO - denpon1:
↳819 (87.59%) uniques of 935 contigs, 0 dupe probe matches, 85 UCE loci removed for
↳matching multiple contigs, 89 contigs removed for matching multiple UCE loci
2016-06-03 16:41:54,879 - phyluce_assembly_match_contigs_to_probes - INFO - lepdecl:
↳900 (91.28%) uniques of 986 contigs, 0 dupe probe matches, 80 UCE loci removed for
↳matching multiple contigs, 81 contigs removed for matching multiple UCE loci
2016-06-03 16:42:05,300 - phyluce_assembly_match_contigs_to_probes - INFO - menmoll:
↳527 (76.27%) uniques of 691 contigs, 0 dupe probe matches, 58 UCE loci removed for
↳matching multiple contigs, 58 contigs removed for matching multiple UCE loci
2016-06-03 16:42:29,353 - phyluce_assembly_match_contigs_to_probes - INFO - onttaul:
↳854 (83.97%) uniques of 1017 contigs, 0 dupe probe matches, 127 UCE loci removed
↳for matching multiple contigs, 130 contigs removed for matching multiple UCE loci

```

```
2016-06-03 16:43:01,303 - phyluce_assembly_match_contigs_to_probes - INFO - tricas1:
↳934 (86.80%) uniques of 1076 contigs, 0 dupe probe matches, 140 UCE loci removed
↳for matching multiple contigs, 141 contigs removed for matching multiple UCE loci
```

## Get match counts and extract FASTA information

Now, we need to get the count of matches that we recovered to UCE loci in the probe set, and extract all of the “good” loci to a monolithic FASTA (see *Tutorial I: UCE Phylogenomics* if this is not making sense):

```
phyluce_assembly_get_match_counts \
  --locus-db in-silico-lastz/probe.matches.sqlite \
  --taxon-list-config in-silico-coleoptera-taxon-sets.conf \
  --taxon-group 'all' \
  --output taxon-sets/insilico-incomplete/insilico-incomplete.conf \
  --log-path log \
  --incomplete-matrix

2016-06-03 16:50:02,610 - phyluce_assembly_get_match_counts - INFO - There are 7 taxa
↳in the taxon-group '[all]' in the config file in-silico-coleoptera-taxon-sets.conf
2016-06-03 16:50:02,610 - phyluce_assembly_get_match_counts - INFO - Getting UCE
↳names from database
2016-06-03 16:50:02,617 - phyluce_assembly_get_match_counts - INFO - There are 1172
↳total UCE loci in the database
2016-06-03 16:50:02,708 - phyluce_assembly_get_match_counts - INFO - Getting UCE
↳matches by organism to generate a INCOMPLETE matrix
2016-06-03 16:50:02,709 - phyluce_assembly_get_match_counts - INFO - There are 1093
↳UCE loci in an INCOMPLETE matrix
2016-06-03 16:50:02,709 - phyluce_assembly_get_match_counts - INFO - Writing the taxa
↳and loci in the data matrix to /nfs/data/working/bfaircloth-insects/coleoptera/
↳triCas1+5+strepsiptera-test/taxon-sets/insilico-incomplete/insilico-incomplete.conf
```

Now, extract the FASTA information for each locus into a monolithic FASTA file:

```
phyluce_assembly_get_fastas_from_match_counts \
  --contigs ../../coleoptera-genome-fasta \
  --locus-db ../../in-silico-lastz/probe.matches.sqlite \
  --match-count-output insilico-incomplete.conf \
  --output insilico-incomplete.fasta \
  --incomplete-matrix insilico-incomplete.incomplete \
  --log-path log
```

## Align the conserved locus data

Now, we need to align the sequence data for each conserved locus in our data set. We’ll do this using standard *phyluce* tools (*mafft*). First, change into the working directory:

```
cd taxon-sets/insilico-incomplete
```

Now, align the sequences:

```
phyluce_align_seqcap_align \
  --fasta insilico-incomplete.fasta \
  --output mafft \
  --taxa 7 \
  --incomplete-matrix \
```

```
--cores 12 \  
--no-trim \  
--output-format fasta \  
--log-path log
```

**Warning:** Note that I am using 12 physical CPU cores (*-cores*) to do this work. You need to use the number of physical cores available on *your* machine.

### Trim the conserved locus alignments

Still following the standard *phyluce* workflow, trim the resulting alignments:

```
phyluce_align_get_gblocks_trimmed_alignments_from_untrimmed \  
--alignments mafft \  
--output mafft-gblocks \  
--b1 0.5 \  
--b4 8 \  
--cores 12 \  
--log log
```

**Warning:** Note that I am using 12 physical CPU cores (*-cores*) to do this work. You need to use the number of physical cores available on *your* machine.

### Remove the locus names from each alignment

And, remove the locus names from each of the resulting alignments:

```
phyluce_align_remove_locus_name_from_nexus_lines \  
--alignments mafft-gblocks \  
--output mafft-gblocks-clean \  
--cores 12 \  
--log-path log
```

**Warning:** Note that I am using 12 physical CPU cores (*-cores*) to do this work. You need to use the number of physical cores available on *your* machine.

### Get stats across the aligned loci

Compute stats across the alignments:

```
python ~/git/phyluce/bin/align/phyluce_align_get_align_summary_data \  
--alignments mafft-gblocks-clean \  
--cores 12 \  
--log-path log
```

```
2016-06-03 16:57:11,675 - phyluce_align_get_align_summary_data - INFO - =====  
↪Starting phyluce_align_get_align_summary_data =====  
2016-06-03 16:57:11,675 - phyluce_align_get_align_summary_data - INFO - Version: git_  
↪6ab3a4b
```

```

2016-06-03 16:57:11,675 - phyluce_align_get_align_summary_data - INFO - Argument --
↳alignments: triCas1+5+strepsiptera-test/taxon-sets/insilico-incomplete/mafft-
↳gblocks-clean
2016-06-03 16:57:11,675 - phyluce_align_get_align_summary_data - INFO - Argument --
↳cores: 12
2016-06-03 16:57:11,675 - phyluce_align_get_align_summary_data - INFO - Argument --
↳input_format: nexus
2016-06-03 16:57:11,675 - phyluce_align_get_align_summary_data - INFO - Argument --
↳log_path: triCas1+5+strepsiptera-test/taxon-sets/insilico-incomplete/log
2016-06-03 16:57:11,676 - phyluce_align_get_align_summary_data - INFO - Argument --
↳output: None
2016-06-03 16:57:11,676 - phyluce_align_get_align_summary_data - INFO - Argument --
↳show_taxon_counts: False
2016-06-03 16:57:11,676 - phyluce_align_get_align_summary_data - INFO - Argument --
↳verbosity: INFO
2016-06-03 16:57:11,676 - phyluce_align_get_align_summary_data - INFO - Getting_
↳alignment files
2016-06-03 16:57:11,710 - phyluce_align_get_align_summary_data - INFO - Computing_
↳summary statistics using 12 cores
2016-06-03 16:57:15,381 - phyluce_align_get_align_summary_data - INFO - -----
↳----- Alignment summary -----
2016-06-03 16:57:15,382 - phyluce_align_get_align_summary_data - INFO - [Alignments]_
↳loci: 994
2016-06-03 16:57:15,382 - phyluce_align_get_align_summary_data - INFO - [Alignments]_
↳length: 644,447
2016-06-03 16:57:15,382 - phyluce_align_get_align_summary_data - INFO - [Alignments]_
↳mean: 648.34
2016-06-03 16:57:15,383 - phyluce_align_get_align_summary_data - INFO - [Alignments]_
↳95% CI: 9.39
2016-06-03 16:57:15,383 - phyluce_align_get_align_summary_data - INFO - [Alignments]_
↳min: 240
2016-06-03 16:57:15,383 - phyluce_align_get_align_summary_data - INFO - [Alignments]_
↳max: 1,444
2016-06-03 16:57:15,383 - phyluce_align_get_align_summary_data - INFO - -----
↳----- Informative Sites summary -----
2016-06-03 16:57:15,384 - phyluce_align_get_align_summary_data - INFO - [Sites] loci:_
↳ 994
2016-06-03 16:57:15,384 - phyluce_align_get_align_summary_data - INFO - [Sites]_
↳total: 169,024
2016-06-03 16:57:15,384 - phyluce_align_get_align_summary_data - INFO - [Sites] mean:_
↳ 170.04
2016-06-03 16:57:15,384 - phyluce_align_get_align_summary_data - INFO - [Sites] 95%_
↳CI: 4.48
2016-06-03 16:57:15,384 - phyluce_align_get_align_summary_data - INFO - [Sites] min:_
↳ 0
2016-06-03 16:57:15,384 - phyluce_align_get_align_summary_data - INFO - [Sites] max:_
↳ 390
2016-06-03 16:57:15,386 - phyluce_align_get_align_summary_data - INFO - -----
↳----- Taxon summary -----
2016-06-03 16:57:15,386 - phyluce_align_get_align_summary_data - INFO - [Taxa] mean:_
↳ 5.76
2016-06-03 16:57:15,386 - phyluce_align_get_align_summary_data - INFO - [Taxa] 95%_
↳CI: 0.07
2016-06-03 16:57:15,386 - phyluce_align_get_align_summary_data - INFO - [Taxa] min:_
↳ 3
2016-06-03 16:57:15,386 - phyluce_align_get_align_summary_data - INFO - [Taxa] max:_
↳ 7
2016-06-03 16:57:15,387 - phyluce_align_get_align_summary_data - INFO - -----
↳----- Missing data from trim summary -----

```

```

2016-06-03 16:57:15,387 - phyluce_align_get_align_summary_data - INFO - [Missing]_
↳mean: 0.00
2016-06-03 16:57:15,387 - phyluce_align_get_align_summary_data - INFO - [Missing] 95%_
↳CI: 0.00
2016-06-03 16:57:15,387 - phyluce_align_get_align_summary_data - INFO - [Missing]_
↳min: 0.00
2016-06-03 16:57:15,387 - phyluce_align_get_align_summary_data - INFO - [Missing]_
↳max: 0.00
2016-06-03 16:57:15,399 - phyluce_align_get_align_summary_data - INFO - -----
↳----- Character count summary -----
2016-06-03 16:57:15,399 - phyluce_align_get_align_summary_data - INFO - [All_
↳characters] 3,655,040
2016-06-03 16:57:15,399 - phyluce_align_get_align_summary_data - INFO - [Nucleotides]_
↳ 3,518,743
2016-06-03 16:57:15,400 - phyluce_align_get_align_summary_data - INFO - -----
↳-- Data matrix completeness summary -----
2016-06-03 16:57:15,400 - phyluce_align_get_align_summary_data - INFO - [Matrix 50%] _
↳ 946 alignments
2016-06-03 16:57:15,400 - phyluce_align_get_align_summary_data - INFO - [Matrix 55%] _
↳ 946 alignments
2016-06-03 16:57:15,400 - phyluce_align_get_align_summary_data - INFO - [Matrix 60%] _
↳ 865 alignments
2016-06-03 16:57:15,400 - phyluce_align_get_align_summary_data - INFO - [Matrix 65%] _
↳ 865 alignments
2016-06-03 16:57:15,401 - phyluce_align_get_align_summary_data - INFO - [Matrix 70%] _
↳ 865 alignments
2016-06-03 16:57:15,401 - phyluce_align_get_align_summary_data - INFO - [Matrix 75%] _
↳ 657 alignments
2016-06-03 16:57:15,401 - phyluce_align_get_align_summary_data - INFO - [Matrix 80%] _
↳ 657 alignments
2016-06-03 16:57:15,401 - phyluce_align_get_align_summary_data - INFO - [Matrix 85%] _
↳ 657 alignments
2016-06-03 16:57:15,401 - phyluce_align_get_align_summary_data - INFO - [Matrix 90%] _
↳ 275 alignments
2016-06-03 16:57:15,401 - phyluce_align_get_align_summary_data - INFO - [Matrix 95%] _
↳ 275 alignments
2016-06-03 16:57:15,401 - phyluce_align_get_align_summary_data - INFO - -----
↳----- Character counts -----
2016-06-03 16:57:15,399 - phyluce_align_get_align_summary_data - INFO - [All_
↳characters] 3,655,040
2016-06-03 16:57:15,399 - phyluce_align_get_align_summary_data - INFO - [Nucleotides]_
↳ 3,518,743
2016-06-03 16:57:15,400 - phyluce_align_get_align_summary_data - INFO - -----
↳-- Data matrix completeness summary -----
2016-06-03 16:57:15,400 - phyluce_align_get_align_summary_data - INFO - [Matrix 50%] _
↳ 946 alignments
2016-06-03 16:57:15,400 - phyluce_align_get_align_summary_data - INFO - [Matrix 55%] _
↳ 946 alignments
2016-06-03 16:57:15,400 - phyluce_align_get_align_summary_data - INFO - [Matrix 60%] _
↳ 865 alignments
2016-06-03 16:57:15,400 - phyluce_align_get_align_summary_data - INFO - [Matrix 65%] _
↳ 865 alignments
2016-06-03 16:57:15,401 - phyluce_align_get_align_summary_data - INFO - [Matrix 70%] _
↳ 865 alignments
2016-06-03 16:57:15,401 - phyluce_align_get_align_summary_data - INFO - [Matrix 75%] _
↳ 657 alignments
2016-06-03 16:57:15,401 - phyluce_align_get_align_summary_data - INFO - [Matrix 80%] _
↳ 657 alignments

```

```
2016-06-03 16:57:15,401 - phyluce_align_get_align_summary_data - INFO - [Matrix 85%] └
↪ 657 alignments
2016-06-03 16:57:15,401 - phyluce_align_get_align_summary_data - INFO - [Matrix 90%] └
↪ 275 alignments
2016-06-03 16:57:15,401 - phyluce_align_get_align_summary_data - INFO - [Matrix 95%] └
↪ 275 alignments
2016-06-03 16:57:15,401 - phyluce_align_get_align_summary_data - INFO - -----
↪----- Character counts -----
2016-06-03 16:57:15,401 - phyluce_align_get_align_summary_data - INFO - [Characters]
↪'-' is present 136,297 times
2016-06-03 16:57:15,402 - phyluce_align_get_align_summary_data - INFO - [Characters]
↪'A' is present 1,047,965 times
2016-06-03 16:57:15,402 - phyluce_align_get_align_summary_data - INFO - [Characters]
↪'C' is present 708,469 times
2016-06-03 16:57:15,402 - phyluce_align_get_align_summary_data - INFO - [Characters]
↪'G' is present 706,567 times
2016-06-03 16:57:15,402 - phyluce_align_get_align_summary_data - INFO - [Characters]
↪'T' is present 1,055,742 times
2016-06-03 16:57:15,402 - phyluce_align_get_align_summary_data - INFO - =====└
↪Completed phyluce_align_get_align_summary_data =====
```

**Warning:** Note that I am using 12 physical CPU cores (`-cores`) to do this work. You need to use the number of physical cores available on *your* machine.

### Generate an incomplete matrix

Now, given the alignments that we have, let's generate a 70% complete matrix

```
phyluce_align_get_only_loci_with_min_taxa \  
  --alignments mafft-gblocks-clean \  
  --taxa 7 \  
  --output mafft-gblocks-70p \  
  --percent 0.75 \  
  --cores 12 \  
  --log log  
  
2016-06-03 16:58:19,687 - phyluce_align_get_only_loci_with_min_taxa - INFO - Copied└  
↪865 alignments of 994 total containing  $\geq 0.75$  proportion of taxa (n = 5)
```

**Warning:** Note that I am using 12 physical CPU cores (`-cores`) to do this work. You need to use the number of physical cores available on *your* machine.

### Prep raxml files, run raxml ML searches, and reconcile best tree w/ bootreps

Setup the PHYLIP-formatted files for raxml:

```
phyluce_align_format_nexus_files_for_raxml \  
  --alignments mafft-gblocks-70p \  
  --output mafft-gblocks-70p-raxml \  
  --log-path log --charsets
```

Now, run raxml against this phylip file

```
raxmlHPC-PTHREADS-SSE3 -m GTRGAMMA -N 20 -p 772374015 -n BEST -s mafft-gblocks-70p.
↪phylip -o menmoll -T 10
raxml/raxmlHPC-PTHREADS-SSE3 -m GTRGAMMA -N autoMRE -p 772374015 -b 444353738 -n
↪bootrep -s mafft-gblocks-70p.phylip -o menmoll -T 10
```

**Warning:** Note that I am using 10 physical CPU cores (*-cores*) to do this work. You need to use the number of physical cores available on *your* machine.

Now, reconcile the best ML tree w/ the bootreps:

```
raxmlHPC-SSE3 -f b \
  -m GTRGAMMA \
  -t RAxML_bestTree.BEST \
  -z RAxML_bootstrap.bootrep \
  -n FINAL -o menmoll
```

And rename the tips. To do this, setup a config file with the old and new names, like:

```
[all]
agrplal:Agrilus planipennis (emerald ash borer)
anogla1:Anoplophora glabripennis (Asian longhorned beetle)
denpon1:Dendroctonus ponderosae (mountain pine beetle)
lepdec1:Leptinotarsa decemlineata (Colorado potato beetle)
menmoll:Mengenilla moldrzyki (Strepsiptera)
onttaul:Onthophagus taurus (taurus scarab)
tricas1:Tribolium castaneum (red flour beetle)
```

And rename the tips:

```
phyluce_genetrees_rename_tree_leaves \
  --order left:right \
  --input-format newick \
  --output-format newick \
  --config rename.conf \
  --section all \
  --input RAxML_bipartitions.FINAL \
  --output RAxML_bipartitions.NAME.FINAL.tre
```





### Citing

If you use the [phyluce](#) code in any form, please cite the following manuscript:

If you are processing UCE data that you have collected by targeted enrichment using our probes/protocols, please cite the following manuscripts, which describes the first use of the general approach:

Then, if you are working in groups other than tetrapods, please cite the appropriate manuscript for the organisms you are working on.

### References

Here are some additional references that have used UCE data for phylogenomic inference at both “deep” and “shallow” timescales. Most of these manuscripts used an older version of the code:

### Other UCE References

### License

### Documentation

The documentation for [phyluce](#) is available under a [CC-BY \(2.0\)](#) license. This license gives you permission to copy, distribute, and transmit the work as well as to adapt the work or use this work for commercial purposes, under the condition that you must attribute the work to the author(s).

If you use this documentation or the [phyluce](#) software for your own research, please cite both the software and (Faircloth et al. 2012). See the Citing section for more detail.

## Software

Copyright (c) 2010-2015, Brant C. Faircloth All rights reserved.

Redistribution and use in source and binary forms, with or without modification, are permitted provided that the following conditions are met:

- Redistributions of source code must retain the above copyright notice, this list of conditions and the following disclaimer.
- Redistributions in binary form must reproduce the above copyright notice, this list of conditions and the following disclaimer in the documentation and/or other materials provided with the distribution.
- Neither the name of the University of California, Los Angeles nor the names of its contributors may be used to endorse or promote products derived from this software without specific prior written permission.

THIS SOFTWARE IS PROVIDED BY THE COPYRIGHT HOLDERS AND CONTRIBUTORS “AS IS” AND ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE ARE DISCLAIMED. IN NO EVENT SHALL THE COPYRIGHT OWNER OR CONTRIBUTORS BE LIABLE FOR ANY DIRECT, INDIRECT, INCIDENTAL, SPECIAL, EXEMPLARY, OR CONSEQUENTIAL DAMAGES (INCLUDING, BUT NOT LIMITED TO, PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES; LOSS OF USE, DATA, OR PROFITS; OR BUSINESS INTERRUPTION) HOWEVER CAUSED AND ON ANY THEORY OF LIABILITY, WHETHER IN CONTRACT, STRICT LIABILITY, OR TORT (INCLUDING NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY OUT OF THE USE OF THIS SOFTWARE, EVEN IF ADVISED OF THE POSSIBILITY OF SUCH DAMAGE.

## Changelog

### v1.5 (April 2015)

- added configuration file system for binary dependencies
- moved many older scripts to attic
- prepended “phyluce” and purpose of program to all programs
- cleaned header and argparse entries
- added tutorial

## Attributions

A large number of people have worked on different aspects of the UCE approach, including creating the laboratory methods to collect the data and the computational methods to analyze the data. Below, we have identified a list of approximately which people/groups did what.

### Maintainer/Author

- [Brant Faircloth](#) (brant **at** faircloth-lab **dot** org)

## Contributed to the code

- Nick Crawford (ngcrawford at gmail dot com)
- Mike Harvey (mharve9 at lsu dot edu)
- Carl Oliveros (carl at ku dot edu)
- Tobias Hofmann (tobiashofmann at gmx dot net)

## Developed the UCE approach

- Brant Faircloth (LSU)
- Travis Glenn (UGA)

## Contributed to the UCE approach

- John McCormack (Occidental College)
- Robb Brumfield (LSU)
- Mike Alfaro (UCLA)
- Nick Crawford (Boston Univ.)
- Mike Harvey (LSU)
- Roger Nilsen (UGA)
- Brian Smith (LSU)
- Laurie Sorenson (LSU)
- Kevin Winker (U. Alaska - Fairbanks)

## Contributed samples, funding, time, comments, etc.

Additionally, the following individuals have contributed samples, funding, laboratory work, computer code, documentation, or all of the above:

- Mike Braun (Smithsonian)
- Noor White (Smithsonian)
- Ed Braun (U. Florida)
- Rebecca Kimball (U. Florida)
- David Ray (MsState)
- Seán Brady (Smithsonian)
- Jesus Maldonado (Smithsonian)
- Jonathan Chang (UCLA)

## Funding

### Primary Sources

The [National Science Foundation \(NSF\)](#) and the [Smithsonian Institution](#) have supported a large portion of our work. The specific programs and proposal identifiers are below:

- [NSF DEB-0841729](#)
- [NSF DEB-0956069](#)
- [NSF DEB-1242241](#)
- [NSF DEB-1242260](#)
- [NSF DEB-1242267](#)
- [Smithsonian Institution Consortium for Understanding and Sustaining a Biodiverse Planet](#)

### Secondary Sources

We have also received funding for computational support from and/or materials from the following organizations:

- [Amazon Web Services](#) (Education grants to BCF, NGC, JEM, and TCG)
- [IDTDNA](#)

## Acknowledgements

We thank the following people, each of whom made contributions ensuring the success of our work. These include:

- Claire Mancuso
- Brant Peterson
- Brent Pederson
- Chris Moran
- Ken Jones
- Joe DeYoung
- LSU Genomics Facility
- UCLA Neuroscience Genomics Core
- 13. Reasel

## List of Programs

Below is a list, by function (directory), of the various programs included in the \* `phyluce` package. This list includes ALL programs in `phyluce`, and some may be lightly tested. This list of programs **does not** include programs currently in the `ATTIC` directory of the repository, because these programs are deprecated.

### **assembly**

- `phyluce_assembly_assemblo_abys`
- `phyluce_assembly_assemblo_trinity`
- `phyluce_assembly_assemblo_velvet`
- `phyluce_assembly_copy_trinity_symlinks`
- `phyluce_assembly_explode_get_fastas_file`
- `phyluce_assembly_extract_contigs_to_barcodes`
- `phyluce_assembly_get_fasta_lengths`
- `phyluce_assembly_get_fastas_from_match_counts`
- `phyluce_assembly_get_fastq_lengths`
- `phyluce_assembly_get_match_counts`
- `phyluce_assembly_get_trinity_coverage`
- `phyluce_assembly_get_trinity_coverage_for_uce_loci`
- `phyluce_assembly_match_contigs_to_barcodes`
- `phyluce_assembly_match_contigs_to_probes`
- `phyluce_assembly_parse_trinity_coverage_for_uce_loci_log`

- `phyluce_assembly_parse_trinity_coverage_log`
- `phyluce_assembly_screen_probes_for_dupes`

## align

- `phyluce_align_add_missing_data_designators`
- `phyluce_align_convert_one_align_to_another`
- `phyluce_align_explode_alignments`
- `phyluce_align_extract_taxa_from_alignments`
- `phyluce_align_extract_taxon_fasta_from_alignments`
- `phyluce_align_filter_alignments`
- `phyluce_align_filter_characters_from_alignments`
- `phyluce_align_format_concatenated_phylip_for_paml`
- `phyluce_align_format_nexus_files_for_mrbayes`
- `phyluce_align_format_nexus_files_for_raxml`
- `phyluce_align_get_align_summary_data`
- `phyluce_align_get_bed_from_lastz`
- `phyluce_align_get_gblocks_trimmed_alignments_from_untrimmed`
- `phyluce_align_get_incomplete_matrix_estimates`
- `phyluce_align_get_indels_from_alignments`
- `phyluce_align_get_informative_sites`
- `phyluce_align_get_only_loci_with_min_taxa`
- `phyluce_align_get_smilogram_from_alignments`
- `phyluce_align_get_taxon_locus_counts_in_alignments`
- `phyluce_align_get_trimmed_alignments_from_untrimmed`
- `phyluce_align_move_align_by_conf_file`
- `phyluce_align_output_list_of_taxon_counts`
- `phyluce_align_parallel_sate`
- `phyluce_align_randomly_sample_and_concatenate`
- `phyluce_align_remove_empty_taxa`
- `phyluce_align_remove_locus_name_from_nexus_lines`
- `phyluce_align_screen_alignments_for_problems`
- `phyluce_align_seqcap_align`
- `phyluce_align_split_concat_nexus_to_loci`

## genetrees

- phyluce\_genetrees\_phybase
- phyluce\_genetrees\_reformat\_raxml\_output
- phyluce\_genetrees\_rename\_tree\_leaves
- phyluce\_genetrees\_run\_raxml\_bootstraps
- phyluce\_genetrees\_run\_raxml\_genetrees
- phyluce\_genetrees\_run\_raxml\_multilocus\_bootstraps
- phyluce\_genetrees\_split\_models\_from\_genetrees

## ncbi

- phyluce\_ncbi\_chunk\_fasta\_for\_ncbi
- phyluce\_ncbi\_example-prep.conf
- phyluce\_ncbi\_prep\_uce\_align\_files\_for\_ncbi
- phyluce\_ncbi\_prep\_uce\_fasta\_files\_for\_ncbi

## snps

- phyluce\_snp\_bwa\_align
- phyluce\_snp\_convert\_vcf\_to\_snapp
- phyluce\_snp\_convert\_vcf\_to\_structure
- phyluce\_snp\_find\_snps\_in\_bed\_interval
- phyluce\_snp\_get\_dbsnp\_freq\_stats
- phyluce\_snp\_get\_dbsnp\_variability\_for\_all\_ucscs
- phyluce\_snp\_prep\_interval\_list\_file\_for\_picard
- phyluce\_snp\_screen\_vcf\_files
- phyluce\_snp\_summarize\_vcf\_file

## utilities

- phyluce\_utilities\_combine\_reads
- phyluce\_utilities\_filter\_bed\_by\_fasta
- phyluce\_utilities\_merge\_multiple\_gzip\_files
- phyluce\_utilities\_merge\_next\_seq\_gzip\_files
- phyluce\_utilities\_replace\_many\_links
- phyluce\_utilities\_sample\_reads\_from\_files
- phyluce\_utilities\_unmix\_fasta\_reads

## probes

- `phyluce_probe_easy_lastz`
- `phyluce_probe_get_clusters_from_bed`
- `phyluce_probe_get_clusters_from_taxa`
- `phyluce_probe_get_genome_sequences_from_bed`
- `phyluce_probe_get_locus_bed_from_lastz_files`
- `phyluce_probe_get_multi_fasta_table`
- `phyluce_probe_get_multi_merge_table`
- `phyluce_probe_get_probe_bed_from_lastz_files`
- `phyluce_probe_get_screened_loci_by_proximity`
- `phyluce_probe_get_subsets_of_tiled_probes`
- `phyluce_probe_get_tiled_probe_from_multiple_inputs`
- `phyluce_probe_get_tiled_probes`
- `phyluce_probe_query_multi_fasta_table`
- `phyluce_probe_query_multi_merge_table`
- `phyluce_probe_remove_duplicate_hits_from_probes_using_lastz`
- `phyluce_probe_remove_overlapping_probes_given_config`
- `phyluce_probe_run_multiple_lastzs_sqlite`
- `phyluce_probe_slice_sequence_from_genomes`
- `phyluce_probe_strip_masked_loci_from_set`



---

## Bibliography

---

- [BCF2015] Faircloth BC. 2016. PHYLUCES is a software package for the analysis of conserved genomic loci. *Bioinformatics* 32:786-788. doi:10.1093/bioinformatics/btv646.
- [BCF2012] BC Faircloth, McCormack JE, Crawford NG, Harvey MG, Brumfield RT, Glenn TC. 2012. Ultraconserved elements anchor thousands of genetic markers spanning multiple evolutionary timescales. *Systematic Biology* 61: 717–726. doi:10.1093/sysbio/SYS004.
- [JEM2012] McCormack JE, Faircloth BC, Crawford NG, Gowaty PA, Brumfield RT, Glenn TC. 2012. Ultraconserved elements are novel phylogenomic markers that resolve placental mammal phylogeny when combined with species tree analysis. *Genome Res* 26:746-754. doi:10.1101/gr.125864.111.
- [NGC2012] Crawford NG, Faircloth BC, McCormack JE, Brumfield RT, Winker K, Glenn TC. 2012. More than 1000 ultraconserved elements provide evidence that turtles are the sister group of archosaurs. *Biol Lett* 8:783-786. doi:10.1098/rsbl.2012.0331.
- [BCF2013] Faircloth BC, Sorenson L, Santini F, Alfaro ME. 2013. A phylogenomic perspective on the radiation of ray-finned fishes based upon targeted sequencing of ultraconserved elements. *PlosONE* 8:e65923. doi:10.1371/journal.pone.0065923.
- [JEM2013] McCormack JE, Harvey MG, Faircloth BC, Crawford NG, Glenn TC, Brumfield RT. 2013. A phylogeny of birds based on over 1,500 loci collected by target enrichment and high-throughput sequencing. *PlosOne* 8:e54848. doi:10.1371/journal.pone.0054848.
- [BTS2013] BT Smith, MG Harvey, BC Faircloth, TC Glenn, RT Brumfield. 2013. Target capture and massively parallel sequencing of ultraconserved elements (UCEs) for comparative studies at shallow evolutionary time scales. *Syst Biol* 63:83-95. doi:10.1093/sysbio/syt061.
- [MGH2014] Sequence capture versus restriction site associated dna sequencing for phylogeography. MG Harvey, BT Smith, TC Glenn, BC Faircloth, RT Brumfield. arXiv:1312.6439.
- [GB2004] Bejerano G, Pheasant M, Makunin I, Stephen S, Kent WJ, et al. (2004) Ultraconserved elements in the human genome. *Science* 304: 1321–1325. doi:10.1126/science.1098119.
- [AS2004] Sandelin A, Bailey P, Bruce S, Engström PG, Klos JM, et al. (2004) Arrays of ultraconserved non-coding regions span the loci of key developmental genes in vertebrate genomes. *BMC Genomics* 5: 99. doi:10.1186/1471-2164-5-99.
- [ED2005] Dermitzakis ET, Reymond A, Antonarakis SE (2005) Opinion: Conserved non-genic sequences — an unexpected feature of mammalian genomes. *Nat Rev Genet* 6:151–157. doi:10.1038/nrg1527.

- [AS2005] Siepel A, Bejerano G, Pedersen JS, Hinrichs AS, Hou M, et al. (2005) Evolutionarily conserved elements in vertebrate, insect, worm, and yeast genomes. *Genome Res* 15: 1034–1050. doi:10.1101/gr.3715005.
- [AW2005] Woolfe A, Goodson M, Goode D, Snell P, McEwen G, et al. (2005) Highly conserved non-coding sequences are associated with vertebrate development. *PLoS Biol* 3:116–130. doi:10.1371/journal.pbio.0030007.
- [LP2006] Pennacchio LA, Ahituv N, Moses AM, Prabhakar S, Nobrega MA, et al. (2006) In vivo enhancer analysis of human conserved non-coding sequences. *Nature* 444: 499–502. doi:10.1038/nature05295.
- [NA2007] Ahituv N, Zhu Y, Visel A, Holt A, Afzal V, et al. (2007) Deletion of Ultraconserved Elements Yields Viable Mice. *PLoS Biol* 5: e234. doi:10.1371/journal.pbio.0050234.
- [WM2007] Miller W, Rosenbloom K, Hardison RC, Hou M, Taylor J, et al. (2007) 28-way vertebrate alignment and conservation track in the UCSC Genome Browser. *Genome Res* 17: 1797–1808. doi:10.1101/gr.6761107.
- [AG2009] Gnirke A, Melnikov A, Maguire J, Rogov P, LeProust EM, et al. (2009) Solution hybrid selection with ultra-long oligonucleotides for massively parallel targeted sequencing. *Nature Biotechnology* 27: 182–189. doi:10.1038/nbt.1523.
- [BB2010] Blumenstiel B, Cibulskis K, Fisher S, DeFelice M, Barry A, et al. (2010) Targeted exon sequencing by in-solution hybrid selection. *Curr Protoc Hum Genet Chapter 18: Unit18.4*. doi:10.1002/0471142905.hg1804s66.