

# HybSeq course — from raw data to species trees

Practical processing of HybSeq target enrichment sequencing data on computing grids like MetaCentrum — pre-processing, HybPiper, alignments, gene trees, species trees

Tomáš Fér, Roswitha Elisabeth Schmickl and **Vojtěch Zeisek**

Department of Botany, Faculty of Science, Charles University, Prague  
Institute of Botany, Czech Academy of Sciences, Průhonice  
<https://trapa.cz/>, [zeisek@natur.cuni.cz](mailto:zeisek@natur.cuni.cz)

June 10–18, 2024



# Outline I

Target enrichment for plant/animal systematics — methodological workshop (MB120C117)

## 1 Introduction

- Test data

- Data processing overview

- Software needed

- MetaCentrum computing environment

## 2 Data and scripts

- General data structure

- Data download and start

- HybSeq scripts

## 3 Data preprocessing

- Trimming and deduplication

## 4 HybPiper

- Preparing data for HybPiper

# Outline II

Target enrichment for plant/animal systematics — methodological workshop (MB120C117)

Processing input files

Retrieving sequences

## 5 Alignments

Alignments of all contigs (and their trimming)

Sorting alignments

## 6 Gene trees

Gene trees from all alignments

Post-processing gene trees

## 7 Comparing gene trees

Visualizing differences among trees

Filtering trees

Species trees

# Outline III

Target enrichment for plant/animal systematics — methodological workshop (MB120C117)

Phylogenetic networks

Comparing trees

8 The end

The very end

## Resources before we start

- Course [Git \(slides with all links\)](#) and information in SIS (česky)
- **Scripts** <https://github.com/V-Z/hybseq-scripts> **shown in the course** — clone the Git repository and adapt to your needs
- Download presentation from <https://trapa.cz/en/hybseq-course-2024>
- Most of the work is done in Linux/UNIX (macOS, ...) command line, so that good knowledge of work in command line is essential, good starting point can be my Linux and MetaCentrum course <https://soubory.trapa.cz/linuxcourse/>
- Many tasks are done in R, so that at least basic knowledge of R is needed, good starting point can be my R course <https://soubory.trapa.cz/rcourse/>
- Processing HybSeq data is computationally demanding (it requires plenty of resources), during the course we use [MetaCentrum, Czech National Grid Infrastructure \(česky\)](#) (slide 14), but any computing cluster or powerful desktop (for patient users;-) can be used

# HybSeq and its data

- **HybSeq** combines target enrichment and genome skimming (see lesson by RS) and especially in larger plant genomes it allows to efficiently select only  $\sim 1000$  single/low copy genes
- It requires sequencing probes, general or group specific (can be design using pipelines like **Sondovač**)
- From sequencing laboratory we get demultiplexed raw FASTQ files
- Steps leading to lists of gene trees require plenty of computing resources and disk space
  - Even simple operation can take significant time — think twice before every step
  - User can select how much resources provide for each step — depends on data size and available resources (more resources like CPU and memory will speed up processing)

## Zingiberaceae test data set

- Family Zingiberaceae has altogether ca. 1600 species throughout tropical Africa, Asia, and the Americas
- We selected 35 members from TF's dataset as test data
- The probes used for sequencing were introduced in [Fér et Schmickl 2018](#) (reduced to 250 probes)
- For data structure see slide [21](#)
- If you did not yet do so, download data (slide [24](#)) – it can take some time...

# Steps from sequencing files to species trees I

- 1 Trimming of raw sequencing FASTQ files (removal of adaptors, ...) e.g. by **Trimmomatic**
- 2 Deduplication of FASTQ reads e.g. by **BBDup**
  - Not strictly required, duplicates mainly provide wrong insight into real coverages of loci
- 3 Checking of FASTQ files in **FastQC** or similar tool and removal of low-quality files
- 4 Preparing probe reference FASTA file and list of samples for processing by HybPiper
- 5 Processing every sample with with **HybPiper** (or alternatively **HybPhyloMaker** — see lessons by TF — or similar tool)
  - 1 Mapping of FASTQ reads with **BWA** to FASTA reference
  - 2 Distributing (sorting) of reads according to successful hits (using **Samtools**) into FASTA files for assembly
  - 3 Assembly of sorted reads with **SPAdes**
  - 4 Alignment of SPAdes contigs against the target sequence
    - Contigs are not expected to overlap much
    - Initial exonerate search is filtered for hits that are above a certain threshold



## Steps from sequencing files to species trees II

- Contigs that pass this filter are arranged in order along the alignment
  - All contigs that pass the previous steps are concatenated into a “supercontig” and the exonerate search is repeated
- 5 Search for paralogs — if SPAdes assembler generates multiple contigs that contain coding sequences representing 75% of the length of the reference protein, HybPiper will print a warning for that gene
  - 6 Recovering of the individual sequences
  - 7 Statistics of the recovery
  - 8 Cleanup of temporary files (especially SPAdes produces huge amount of data unneeded for further processing)
- 6 Statistics of sequence lengths in all samples and more information about recovered contigs
  - 7 Creation of heatmaps of sequence recoveries
  - 8 Summary of warnings about paralogs
  - 9 Retrieve of sequences of exons, introns and supercontigs for all samples

## Steps from sequencing files to species trees III

- 10 Alignment of all contigs (e.g. by **MAFFT**; or **MUSCLE**, **Clustal**, ... e.g. using **R** and packages **ape** and/or **ips**)
  - All alignments must be trimmed – columns/rows with too much missing data (e.g. beginning and end of the alignment) must be removed (e.g. using **R** and package **ape**)
  - It is also useful to create simple NJ tree graphical check of alignment (e.g. using **R** and package **ape**)
- 11 Sorting of alignments, statistics of their length and quality, discarding of poor (too short, too few individuals, too few variable positions, ...) alignments
- 12 Reconstruction of gene trees from all aligned contigs (e.g. using **IQ-TREE**, or **ExaML**, **MrBayes**, **PhyML**, **RAxML-NG**, ...)
- 13 Post-processing of gene trees
  - Identification, inspection and possible removal of gene trees with significantly different topology (e.g. by **R** and packages **ape** and **kdetrees**; **TreeShrink**, etc.)

## Steps from sequencing files to species trees IV

- Comparison of gene trees (e.g. heatmaps and PCoA by R and packages `ade4`, `ape`, `distory`, `phytools`, etc.)
- Comparison of (several) (species) trees (e.g. by R and packages `ape` or `phytools`)
- 14 Construction of species trees (e.g. by `ASTRAL-III`, `ASTER*` or `ASTRID-2`)
  - Comparison of species tree and gene trees (e.g. by `phyparts` and `MJPythonNotebooks`)
- 15 Phylogenetic networks (e.g. by `PhyloNet`)
- 16 And more...

### Note...

- This general scheme can be significantly altered...
- There are plenty of technical as well as biological problems (HGT, ILS, ...) and new software keep being developed...
- Much more analysis possible...

## List of software used during the course I

- **ASTRAL** (see lesson by TF) — species trees from gene trees
- BASH 4 or later and GNU core utils (“Linux command line”)
- **BMap** — deduplication of FASTQ
- **BLAST+** (used by **HybPiper**)
- **BWA** (used by **HybPiper**)
- **Dendroscope** — visualize outputs of **PhyloNet**
- **Exonerate** (used by **HybPiper**)
- **GNU Parallel** (used by **HybPiper** and in BASH scripts)
- **HybPiper** — recovering genes from targeted sequence capture data
- **IQ-TREE** — gene trees
- **MAFFT** — alignment

## List of software used during the course II

- MJPythonNotebooks (used by phyparts)
- PhyloNet — phylogenetic networks
- phyparts — comparison of species tree vs. gene trees
- Python 3.7 or later and Biopython 1.80 or later (used by HybPiper)
- R 4.0 or later and packages ade4, adegenet, ape, corrplot, distory, ggplot2, gplots, ips, kdetrees, pegas, phangorn, phytools and scales
  - Used for alignment of contigs, post-processing of alignments, post-processing and comparison of gene trees, etc.
- Samtools (used by HybPiper)
- SPAdes (used by HybPiper)
- TreeShrink — detection of outlier long branches in collections of phylogenetic trees
- Trimmomatic — trimming of FASTQ

# CESNET and MetaCentrum I

- **CESNET (česky)** is organization of Czech universities, Academy of Science and other organizations taking care about Czech backbone Internet, one of world leading institutions of this type
- CESNET provides various **services (česky)**
  - Massive computations — **MetaCentrum (česky)** — **this we need to process our HybSeq data**
  - Large **data storage (česky)** — **this we can use to store our HybSeq data**
  - And **much more (česky)** — note especially **FileSender (česky)** and **ownCloud (česky)**...
  - See also my course [https://soubory.trapa.cz/linuxcourse/linux\\_bash\\_metacentrum\\_course.pdf](https://soubory.trapa.cz/linuxcourse/linux_bash_metacentrum_course.pdf) (chapter “MetaCentrum”)
- Information about data storage <https://du.cesnet.cz/en/start> (česky) contains detailed usage instructions
- Information about MetaCentrum <https://www.metacentrum.cz/en/> (česky)
- Most of practical information for users are at <https://docs.metacentrum.cz/>
  - Old **wiki** is deprecated

## CESNET and MetaCentrum II

- To start work see at least **access**, **computing**, **work with data** and some **tutorial**
- Of course, good knowledge of work in Linux command line (BASH) is needed...

### MetaCentrum vs. other grids and clusters...

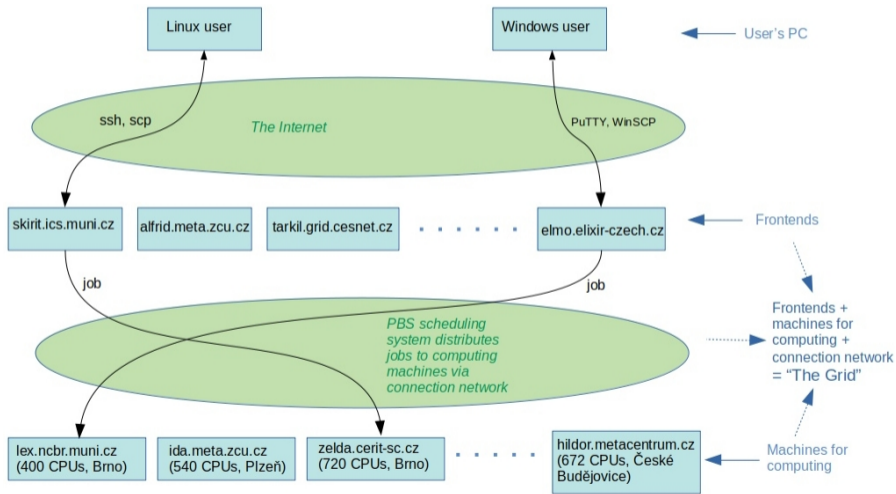
- I show processing on MetaCentrum Czech National Grid Infrastructure, as it is readily available, well maintained and contains all needed applications, but it's possible to use any computing cluster in similar way (the scripts are easily modifiable)
- For other clusters than Czech MetaCentrum, user must in the **scripts** edit at least commands to load modules and handling with scratch directory
- MetaCentrum is distributed — it has multiple front end and storages, which might be confusing (be careful with paths in scripts)

# MetaCentrum

- Find all needed information at <https://docs.metacentrum.cz/>
- Current state and usage as available at <https://metavo.metacentrum.cz/>
- Manage your user account at <http://metavo.metacentrum.cz/en/myaccount/> (česky)
- Personal view on actual resources and running tasks is at <https://metavo.metacentrum.cz/pbsmon2/person>
- Listing of available applications <https://docs.metacentrum.cz/software/search-soft/>
- It has several **front ends** where users log, various **storages**, and thousands of computers (nodes) doing the calculations — they are not accessed directly to run task
  - Distributed nature and number of front ends and storages may be confusing for beginners
- Most of computers are running **Debian GNU/Linux**



# Basic workflow



## Launching of tasks

- <https://docs.metacentrum.cz/computing/run-basic-job/>
- Personal view <https://metavo.metacentrum.cz/pbsmon2/person> has nice overview of available resources and tasks and allows comfortable construction of submission command

```
1 # Task in metacentrum.sh script will run up to 4 days (96 h),
2 # require one physical computer with 8 CPU threads, 24 GB of
3 # RAM, 10 GB of disk space and we get all information mails
4 qsub -l walltime=96:0:0 -l \
5     select=1:ncpus=8:mem=24gb:scratch_local=10gb -m abe metacentrum.sh
6 # Check how the task is running (above web) and
7 qstat -u $USER # Information about $USER's jobs
8 qstat 123456789 # The task ID is available from qstat
9 qstat -f 123456789 # Print a lot of details
10 qdel 123456789 # Terminate scheduled or running task
```

## Key MetaCentrum commands

- MetaCentrum is “just” normal Linux server — work as usually
- Command `module` loads/unloads selected `application` (`module add r`)
- Tasks (BASH scripts) are submitted for computing by `qsub` — the script must copy the data into `$SCRATCHDIR` and do all calculations there
  - It has plenty of options how to specify requirements (see `help`)
- Queued and running jobs can be seen by `qstat -u $USER` (`qstat` has much more options) and any job can be terminated by `qdel 123456789` (number from `qstat`)

```

1 module avail XXX/ # List available modules for XXX
2 module add <TAB><TAB> # Load some module (or use 'module load XXX')
3 module list # List of currently loaded modules
4 module rm XXX # Unload selected module
5 # All jobs for $USER ('-x' also shows tasks ended within last 24 h)
6 qstat -w -n -1 -u $USER -t -p -a @pbs-m1.metacentrum.cz -x

```

## Running R tasks on MetaCentrum

- There are only some R packages, to get more create own package library and use it in scripts (see e.g. `.libPaths()` within R)
- **Be careful about paths!**
- In the `metacentrum.sh` script load R, e.g.  
`module add r/4.1.3-gcc-10.2.1-6xt26dl` and start there R script as usually `R CMD BATCH script.r`
- ① Login to selected front node via SSH and load R module
- ② Create somewhere new directory for R packages `mkdir rpkg` (or use default `~/R/`)
- ③ Start R `R` and install **all** R packages needed for the task – install them e.g. into the `rpkg` directory `install.packages(pkgs=..., ..., lib="rpkg")`
- ④ In the R script `*.r` load the packages from the `rpkg` directory  
`library(package=..., lib.loc="/storage/.../rpkg")` or so
- ⑤ Ensure all needed outputs are saved from the R script

# Test data directory structure I

```
1 |— hybseq # BASH scripts to process the data
2 |I |— bin # Scripts themselves
3 |I |— ref # References for HybPiper
4 |I |— rpackages # R library with installed R packages
5 |└─ hybseq_course_zingibers # Data
6 |   |— 1_data # Sequencing libraries
7 |   |I |— lib_01 # Sequencing library 1
8 |   |I |I |— 0_data # Raw FASTQ sequences
9 |   |I |I |— 1_trimmed # Trimmed FASTQ sequences
10 |  |I |I |— 2_dedup # Deduplicated FASTQ sequences - for HybPiper
11 |  |I |I |└─ 3_qual_rep # FastQC quality reports
12 |  |I |└─ lib_02 # Possible sequencing library 2
13 |  |I |   |— 0_data # Raw FASTQ sequences
14 |  |I |   |└... # Same structure as lib_01...
15 |  |... |   ... # Next slide...
```

## Test data directory structure II

```
1 ... # Previous slide...
2 |— 2_seqs # Outputs of HybPiper
3 |   |— Aframomum-alboviolaceum_S118_L001.dedup # Sample output dir
4 |   |   |— Assembly_1 # Gene output directory
5 |   |   |— ... # More gene output directories...
6 |   |   |— Assembly_16358 # Gene output directory
7 |   |— ... # More samples...
8 |   |— Zingiber-officinale_S242_L001.dedup # Sample output dir
9 |       |— Assembly_1 # Gene output directory
10 |       |— ... # More gene output directories...
11 |       |— Assembly_16358 # Gene output directory...
12 |   # Plenty of reconstructed contigs, statistics, etc.
13 |— 3_aligned # Aligned contigs
14 |   |— exons # Aligned exons
15 |   |— introns # Aligned introns
16 |   |— supercontigs # Aligned supercontigs
17 ... # Next slide...
```

## Test data directory structure III

```
1 ... # Previous slide...
2 |— 4_gene_trees # Reconstructed gene trees
3 |   |— exons # Gene trees of exons
4 |   |— introns # Gene trees of introns
5 |   |— supercontigs # Gene trees of supercontigs
6 |   # Possibly create some substructure for outputs of kdetrees,
7 |   # TreeShrink and another gene trees filtration... e.g.
8 |   # mkdir {exons,introns,supercontigs}_{0_trees,1_unfiltered,
9 |   #   2_pcoa_kdetrees,3_treeshrink}
10 |   # To keep everything sorted
11 |— 5_species_trees # Processing of the trees
12 |   |— # Outputs of various analysis like ASTRAL...
13 |   ... # More downstream analysis...
```

- Of course, every user can figure different directory structure, but HybSeq produces a lot of data and plenty of software packages are used, so keep some logical structure...

## Data download

- MetaCentrum storages have sometimes too limited quota for such a large data — see your [quotas \(česky\)](#), see [details](#)
- The [pipeline](#) produce a lot of data (especially [HybPiper](#); gene trees can be also large) — ensure to have enough space to store everything
  - Especially **HybPiper produces a lot of files** — **user may reach quota for number of files**, not necessarily (only) amount of data

```
1 # Login to any MetaCentrum front end
2 ssh USER@skirit.metacentrum.cz # Or any other front end
3 # Download the data
4 wget https://botany.natur.cuni.cz/zeisek/hybseq_course_zingibers_0_data.zip
5 # Unpack it
6 unzip hybseq_course_zingibers_0_data.zip
```

- The archive contains input raw data and reference
- Further steps are in other archives (see later)



## Scripts and other resources to process the data

- The archive `hybseq_course_zingibers_0_data.zip` contains directory `hybseq_course_zingibers` with...
  - Raw data `*.fastq.gz` in `1_data/lib_01/0_data/`
  - Reference for HybPiper `curcuma_HybSeqProbes_first958_concat.fasta` in `ref/`

Do not blindly copy-paste commands...

Commands show typical way how to proceed, but should not be using without understanding, and other ways how to work are possible... Command and submission scripts also **do require** edits prior launching them.

# Getting HybSeq scripts

- See <https://github.com/V-Z/hybseq-scripts>
- Install the scripts as follows:

```
1 # Ensure to be in home folder
2 cd
3 # Clone repository with scripts into "hybseq" directory
4 git clone https://github.com/V-Z/hybseq-scripts.git hybseq
5 # Install all needed R packages
6 cd hybseq/ # Go to hybseq directory
7 ls -lha # See content...
8 module add r/4.1.3-gcc-10.2.1-6xt26dl # Load R module
9 R # Start R and install needed packages (next slide)
```

# Install needed R packages

```
1 getwd() # Ensure to be in ~/hybseq/ And if not go there by
2 setwd("~/hybseq/") # Set path into HybSeq directory
3 dir() # Ensure to see directory "rpackages"
4 # Package Matrix (required by ips) requires R 4.4 and newer while we have
5 # 4.1.3 - we have to install its dependency lattice and older version of
6 # Matrix manually
7 install.packages(pkgs="lattice", lib="rpackages",
8   repos="https://mirrors.nic.cz/R/", dependencies="Imports")
9 install.packages(pkgs="https://cran.r-project.org/src/contrib/Archive/
10   Matrix/Matrix_1.6-5.tar.gz", lib="rpackages", repos=NULL,
11   dependencies="Imports")
12 # Install needed packages in R
13 install.packages(pkgs=c("ape", "codetools", "cpp11", "farver",
14   "ips", "RcppArmadillo", "scales"), lib="rpackages",
15   repos="https://mirrors.nic.cz/R/", dependencies="Imports")
16 dir("rpackages/") # Ensure all packages are correctly installed
```

# Trimming and deduplication

- Raw demultiplexed FASTQ sequences must be trimmed (sequencing adaptors removed, ...) and should be deduplicated (removal of artificial duplicates to get correct statistics of coverage)
- There are plenty of software packages available, **Trimmomatic** use to be used for trimming and e.g. **BBMap** for deduplication
- Usually, libraries are processed as they are delivered from sequencing company
- Quality of all FASTQ files should be checked by e.g. **FastQC**
- It is practical to obtain simple statistics — number of sequences in original files, after trimming and after deduplication
- Low quality files should be discarded...
- Everything can be easily coded into simple BASH script processing all files (see following slides)

## Scripts to trim and deduplicate FASTQ sequences

- Script `~/hybseq/bin/hybseq_1_prep_1_qsub.sh`
  - Submits via `qsub` script `~/bin/hybseq_1_prep_2_run.sh` FASTQ data for trimming, deduplication and quality checks
  - Variables `WORKDIR` (location of HybSeq scripts) and `DATADIR` (data to process) must point to correct existing location – script contains settings for running following script
  - If changing parameters for `hybseq_1_prep_2_run.sh`, parameters for `qsub` must be changed here accordingly
  - The script is started multiple times to process all sequencing libraries
- Script `~/hybseq/bin/hybseq_1_prep_2_run.sh`
  - See `./hybseq_1_prep_2_run.sh -h` for usage help
  - Script checks if all needed parameters and tools are available and in simple `for` loop trims all sequences (one by one), deduplicates them, does quality checking, prints simple statistics, and prepares list of samples for HybPiper
- Output is prepared as input for HybPiper itself

## Submission of `hybseq_1_prep_1_qsub.sh`

- Pre-processing data for HybPiper
- `~/hybseq/bin/hybseq_1_prep_1_qsub.sh` must be edited before submission via `qsub` – change `WORKDIR` and `DATADIR`
- After it runs for a while (everything had been copied to the computing node), it is possible to change `DATADIR` and submit processing of the second library

```
1 # After edition of WORKDIR and DATADIR, go to data directory...
2 cd ~/hybseq_course_zingibers/1_data/lib_01/ # ...and submit the job...
3 qsub -l walltime=24:0:0 -l \
4     select=1:ncpus=4:mem=64gb:scratch_local=100gb -m abe \
5     ~/hybseq/bin/hybseq_1_prep_1_qsub.sh
6 # Or similar command
7 # Monitor running $USER's tasks with details
8 qstat -w -n -l -u $USER # Last column contains machine name
9 # See your processes running the machine (from above list)
10 ssh exec_host "ps ux" # Replace exec_host by hostname!
```

# Trimm and deduplicate all FASTQ files

## Tasks to pre-process FASTQ data for HybPiper

- 1 Inspect `~/hybseq/bin/hybseq_1_prep_1_qsub.sh` and `~/hybseq/bin/hybseq_1_prep_2_run.sh` and be sure to understand what the scripts do, including syntax used.
- 2 Edit declarations of variables `WORKDIR` and `DATADIR` in `hybseq_1_prep_1_qsub.sh` so that they point to correct locations.
- 3 Submit via `qsub` `hybseq_1_prep_1_qsub.sh` to process the data and monitor the jobs during processing.
- 4 Inspect outputs of `hybseq_1_prep_2_run.sh`, including statistics and FASTQ checks. What do they show?
- 5 Can you run the task on your computer (desktop or notebook) without `qsub`? If so, how?

## When the pre-processing job is done

- Wait for task to finish, or **download** and inspect test data
- Output directory `job_123456789.pbs-m1.metacentrum.cz` will be in `DATADIR (0_data)` – go there and move directories `1_trimmed`, `2_dedup` and `3_qual_rep` and log file `hybseq_prepare.log` to same level as `0_data`, e.g.  

```
mv 1_trimmed 2_dedup 3_qual_rep hybseq_prepare.log ../../
```
- Check log files `hybseq_prepare.log` and `hybseq_1_prep_1_qsub.sh.[eo]123456789`
- Directory `1_trimmed` contains statistics in `report_trimming.tsv`
- Directory `2_dedup` contains statistics in `report_filtering.tsv` and everything needed for HybSeq (see following chapter)
- See HTML files with quality reports in directory `3_qual_rep`

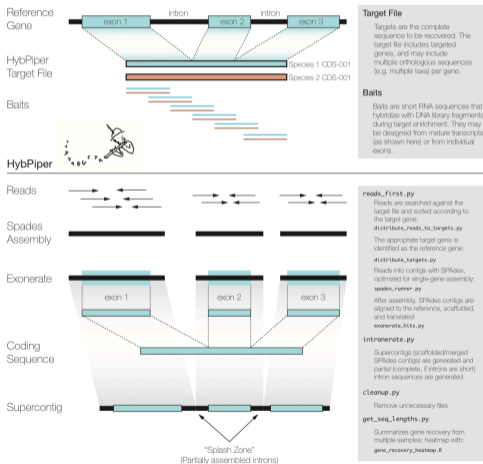


# HybPiper

- See [README](#) and [wiki](#) documentation
- HybPiper was designed for targeted sequence capture, in which DNA sequencing libraries are enriched for gene regions of interest, especially for phylogenetics
- The HybPiper pipeline starts with high-throughput sequencing reads, and assigns them to target genes using BLASTx/DIAMOND or BWA. The reads are distributed to separate directories, where they are assembled separately using SPAdes. The main output is a FASTA file of the (in frame) CDS portion of the sample for each target region, and a separate file with the translated protein sequence.
- Includes commands to extract the intronic regions flanking each exon, and investigate putative paralogs
- End with contigs (one for each probe) ready to align

# Overview of HybPiper

## Target file and bait design (pre-HybPiper)



## Requirements to run HybPiper

- See <https://github.com/mossmatters/HybPiper#readme> — HybPiper has plenty of requirements, the easiest is probably installation using `conda`
- It is possible to run HybPiper on your computer (with Linux or macOS), but it requires plenty of CPU and memory and creates huge output directories...
- `hybpiper assemble` is processing single sample — to process multiple (by `while` loop or by `~/hybseq/bin/hybseq_2_hybpiper_1_submitter.sh`) there must be list of sample base names without suffixes like `*[.]R{1,2}[.]*f*q*` (here `samples_list.txt` created by `hybseq_1_prep_2_run.sh`)
- Reference bait FASTA file **must** have sequences named as `>Species_name-gene_id` (see `help`) (**note** order and dash in between)
- **HybPiper** is processing individual files with given baits FASTA sequences — batch processing must be scripted

## Before running HybPiper

- Directory `~/hybseq/ref/` contains prepared reduced bait file `curcuma_HybSeqProbes_first958_concat.fasta` (it will be used to process test data)
- To run **HybPiper** we need the target reference file above and all required software, and BASH script to process all input files in batch — see `hybseq_2_hybpiper_1_submitter.sh` (submitting via `qsub` all input files), `hybseq_2_hybpiper_2_qsub.sh` (preparing individual job to run) and `hybseq_2_hybpiper_3_run.sh` (processing every input file — doing the job)
- **Settings** `WORKDIR`, `DATADIR`, `SAMPLES`, `BAITFILE` and `NCPU` **must be changed** in `~/hybseq/bin/hybseq_2_hybpiper_1_submitter.sh`

# Prepare to run HybPiper

## Tasks before running HybPiper

- 1 Check `samples_list.txt` in `2_dedup` output directory of the test data after running `~/hybseq/bin/hybseq_1_prep_1_qsub.sh`, and check how it was created.
- 2 Compare [https://github.com/tomas-fer/HybPhyloMaker/blob/master/HybSeqSource/curcuma\\_HybSeqProbes\\_test.fa](https://github.com/tomas-fer/HybPhyloMaker/blob/master/HybSeqSource/curcuma_HybSeqProbes_test.fa) (output of Geneious assembler) and `curcuma_HybSeqProbes_first958_concat.fasta` — see form required by HybPiper
- 3 Do we have everything needed to run HybPiper?

# Understanding how presented scripts run HybPiper I

- Script `~/hybseq/bin/hybseq_2_hybpiper_1_submitter.sh` goes to directory set by `DATADIR` variable and for every sample (pair of forward and reverse FASTQ files) listed in `samples_list.txt` (variable `SAMPLES`, created by `hybseq_1_prep_2_run.sh`) and submits via `qsub` individual task for each input sample (they are processed independently in parallel)
- Script `~/hybseq/bin/hybseq_run_2_hybpiper_2_qsub.sh` drives submission and manipulation of processing of each individual file — it checks if all required data are provided, copy input files to `SCRATCHDIR`, loads required software modules, runs processing itself by `hybseq_2_hybpiper_3_run.sh` and copy results back
- Script `~/hybseq/bin/hybseq_2_hybpiper_3_run.sh` runs all HybPiper steps for individual samples (as submitted by previous scripts)

## Understanding how presented scripts run HybPiper II

- Outputs of `~/hybseq/bin/hybseq_2_hybpiper_3_run.sh` should be checked and moved into special directory (`~/hybseq_course_2023_zingibers/2_seqs`) where new list of samples must be created
- Script `~/hybseq/bin/hybseq_run_3_hybpiper_postprocess_1_qsub.sh` can be then used to run `~/hybseq/bin/hybseq_3_hybpiper_postprocess_2_run.sh` which uses HybPiper to retrieve sequences of exons, introns and supercontigs; prints statistics, creates heatmaps, and summaries paralogs warnings — retrieved sequences can then be aligned (see further)

## Submitting jobs to run HybPiper

- To retrieve probe sequences from every input FASTQ file
- `~/hybseq/bin/hybseq_2_hybpiper_1_submitter.sh` must be edited before running it — change `WORKDIR`, `DATADIR`, `SAMPLES`, `BAITFILE` and `NCPU` according to your needs
- After submission of all input files is done, it is possible to change `DATADIR` and submit processing of the second library
- It will start job for every sample, so that output of `qstat` will notably prolong...

```
1 # After edition of WORKDIR, DATADIR, SAMPLES, BAITFILE and NCPU, go to
2 # DATADIR AND run simply (There are no further parameters on command line)
3 cd ~/hybseq_course_zingibers/1_data/lib_01/2_dedup/ # Go to DATADIR
4 ~/hybseq/bin/hybseq_2_hybpiper_1_submitter.sh # Submit all samples
5 # Monitor running $USER's tasks with details
6 qstat -w -n -1 -u $USER # Last column contains machine name
7 # See your processes running the machine (from above list)
8 ssh exec_host "ps ux" # Replace exec_host by hostname!
```



# Running HybPiper

## Run HybPiper

- 1 Inspect scripts `~/hybseq/bin/hybseq_2_hybpipe_1_submitter.sh`, `~/hybseq/bin/hybseq_2_hybpipe_2_qsub.sh` and `~/hybseq/bin/hybseq_2_hybpipe_3_run.sh` and be sure to understand what they are doing, including syntax used.
- 2 Check <https://github.com/mossmatters/HybPiper/wiki/Full-pipeline-parameters#10-hybpipe-assemble>
- 3 Edit declarations of variables `WORKDIR`, `DATADIR`, `SAMPLES` and `BAITFILE` in `hybseq_2_hybpipe_1_submitter.sh` so that they point to correct locations.
- 4 Run in directory `2_dedup` script processing all samples, simply `~/hybseq/bin/hybseq_2_hybpipe_1_submitter.sh`

# Other tasks with HybPiper

For interested students

## HybPiper tasks for advanced users

- 1 Think how to process all samples on single computer. Use `while` BASH loop and feed it by `samples_data.txt`. How would such script look like? What had to be changed?
- 2 Check requirements of software used by HybPiper (BWA, SPAdes, ...) and think if changing of number of CPU threads and memory would significantly speed up processing of individual files.
- 3 See `man qsub` and think if there is another option how to pass options from `hybseq_2_hybpiper_1_submitter.sh` for individual `qsub` commands (apart of usage of exported variables as is used now).
- 4 Think about HybPiper parameters in `hybseq_2_hybpiper_3_run.sh`.

## After running HybPiper for every sample...

- Previous step processed by HybPiper every single sample individually and independently
- HybPiper postprocessing will retrieve contigs for every exon, intron and supercontig containing all individuals where particular genetic region was found; and do some statistics
- Script `~/hybseq/bin/hybseq_3_hybpiper_postprocess_1_qsub.sh` is submitted via `qsub` and it runs  
`~/hybseq/bin/hybseq_3_hybpiper_postprocess_2_run.sh` which uses HybPiper to obtain table of contig lengths, some statistics, heatmaps, information about paralogs, and retrieves from individual directories contigs to be aligned
- New list of samples for all samples (all libraries) must be prepared (next slide)
- All outputs (a lot of files) are in same directory where HybPiper outputs are (typically `XXX/2_seqs/`)

## Sorting data after running HybPiper for all samples

```
1 # Ensure You are in 2_dedup
2 cd ~/hybseq_course_zingibers/1_data/lib_01/2_dedup/
3 # Move MetaCentrum logs into output directories
4 while read L; do mv HybPiper."$L".[eo]* "$L"/; done < samples_list.txt
5 # Move from ~/hybseq_course_zingibers/1_data/lib_01/2_dedup all outputs of
6 # HybPiper to hybseq_course_zingibers/2_seqs/ for next steps
7 mkdir ../../../../2_seqs # Create directory for HybPiper outputs
8 mv *.dedup ../../../../2_seqs/ # Move all HybPiper output directories there
9 cd ../../../../2_seqs/ # Go to that directory
10 # NOTE: The mkdir command is needed only once and directory 2_seqs can
11 # contain results from multiple sequencing libraries
```

- In `~/hybseq_course_zingibers/1_data/lib_#/2_dedup/` can be multiple libraries — all HybPiper outputs should be moved into `~/hybseq_course_zingibers/2_seqs/` where HybPiper postprocessing will retrieve all sequences and statistics

## HybPiper postprocessing — retrieving of contigs and statistics

- Code below is exemplary and requires edits, as well as script

```
~/hybseq/bin/hybseq_run_3_hybpiper_postprocess_1_qsub.sh —  
change WORKDIR, BAITFILE, DATADIR and SAMPLES
```

- This process can be repeated every time new sequencing library is added
- All outputs will be directly in DATADIR ( 2\_seqs )

```
1 # Create in ~/hybseq_course_zingiberess/2_seqs/ new samples_list.txt  
2 find . -maxdepth 1 -type d | sed 's/^\.\.\/' | sort | tail -n+2 > \  
3 samples_list.txt  
4 # Edit WORKDIR, BAITFILE and DATADIR in  
5 # ~/hybseq/bin/hybseq_3_hybpiper_postprocess_1_qsub.sh  
6 # When ready, submit task to post-process HybPiper results  
7 qsub -l walltime=12:0:0 -l select=1:ncpus=1:mem=8gb:scratch_local=100gb \  
8 -m abe ~/hybseq/bin/hybseq_3_hybpiper_postprocess_1_qsub.sh  
9 # Note that for large data (~1000 probes and hundreds of samples) the  
10 # postprocessing can take long time and longer walltime can be needed
```

# Post-processing HybPiper outputs and retrieving contig sequences I

## Tasks to post-process HybPiper outputs I

- 1 Inspect `~/hybseq/bin/hybseq_3_hybpiper_postprocess_1_qsub.sh` and `~/hybseq/bin/hybseq_3_hybpiper_postprocess_2_run.sh` and be sure to understand what the scripts do, including syntax used.
- 2 Edit declaration of variables `WORKDIR`, `BAITFILE` and `DATADIR` in `~/hybseq/bin/hybseq_3_hybpiper_postprocess_1_qsub.sh` so that they point to correct locations.
- 3 Submit via `qsub`  
`~/hybseq/bin/hybseq_3_hybpiper_postprocess_1_qsub.sh` to post-process all HybPiper outputs.
- 4 Inspect outputs, including statistics, heatmaps and log. What do they show?

# Post-processing HybPiper outputs and retrieving contig sequences II

For interested students

- Wait for tasks to finish, or **download** and inspect test data

## Tasks to post-process HybPiper outputs II

- 1 Which part does take the longest time on this task? Does this task take plenty of resources (CPU, memory)?
- 2 Check parameters of `hybpiper stats`, `recovery_heatmap`, `retrieve_sequences` and `paralog_retriever` at <https://github.com/mossmatters/HybPiper/wiki/Full-pipeline-parameters>
- 3 Can you run the task on your computer (desktop in office or notebook) without `qsub`? If so, how? Can you run the task directly on MetaCentrum front end?

# Alignment of all contigs I

- All sequences retrieved in the previous step with HybPiper must be aligned (by any aligner)
- Alignments must be post-processed
  - Rows (individuals) and/or columns (positions within alignment) with more than  $\sim 10\text{--}50\%$  of missing data must be removed (e.g. beginning and the end of the alignment)
  - Too short alignments or alignments with too few variable sites or too few individuals could be removed (their usefulness is questionable)
  - Exact thresholds are to be discussed (according to purpose, phylogenetic scale etc.)...
- Script `~/hybseq/bin/hybseq_4_alignment_1_submitter.sh` goes to directory set by `DATADIR` variable and for every contig (all `*.fasta` or `*.FNA` files) in that directory (retrieved by `hybseq_3_hybpiper_postprocess_2_run.sh`) submits via `qsub` individual alignment task



## Alignment of all contigs II

- Script `~/hybseq/bin/hybseq_run_4_alignment_2_qsub.sh` drives submission and manipulation of processing of each individual file – it checks if all required data are provided, copy input files to `SCRATCHDIR`, loads required software modules, runs processing itself by `hybseq_4_alignment_3_run.r` R script and copy results back
- R script `~/hybseq/bin/hybseq_4_alignment_3_run.r` aligns every input contig with `MAFFT`, trims the alignment, creates NJ tree (NWK and PNG), creates image of alignment and reports alignment details
- Script `~/hybseq/bin/hybseq_4_alignment_4_postprocess.sh` is short, can be runned on the front end, it will sort outputs into subdirectories for exons, intron and supercontigs, create statistics of alignments and lists of NJ gene trees

## Submitting alignment jobs

- To align all retrieved contigs (sequences)
- `~/hybseq/bin/hybseq_4_alignment_1_submitter.sh` must be edited before running it — change `WORKDIR` and `DATADIR`
- It will start job for every sample, so that output of `qstat` will be very long (3 jobs — respective exon, intron and supercontig — for every probe)...
- For small contigs, jobs are very fast

```
1 # After edition of WORKDIR and DATADIR run simply
2 # NOTE: There are no further parameters on command line
3 ~/hybseq/bin/hybseq_4_alignment_1_submitter.sh
4 # Monitor running $USER's tasks with details
5 qstat -w -n -1 -u $USER # Last column contains machine name
6 # Something went wrong? Cancel all running or queued tasks by
7 qdel $(qstat -u $USER | grep -o "[0-9]\+" | tr "\n" " ")
```

# Running alignments

## Run alignments

- 1 Inspect scripts `~/hybseq/bin/hybseq_4_alignment_1_submitter.sh`, `~/hybseq/bin/hybseq_4_alignment_2_qsub.sh` and `~/hybseq/bin/hybseq_4_alignment_3_run.r` and be sure to understand what they are doing, including syntax used.
- 2 Edit declarations of variables `WORKDIR` and `DATADIR` in `hybseq_4_alignment_1_submitter.sh` so that they point to correct locations.
- 3 Process all `*.FNA` and `*.fasta` files.
- 4 Monitor progress of the jobs
- 5 Inspect outputs (including images) and log files.

# Other tasks with alignments

For interested students

## Alignment tasks for advanced users

- 1 Think how to process all samples on single computer. Use `find` to list all `*.FNA` and `*.fasta` files and pass them to `GNU Parallel`. How would such script look like? What had to be changed?
- 2 Check requirements of `MAFFT` and/or another aligners and think if changing of number of CPU threads and memory would significantly speed up processing of individual files.
- 3 In `hybseq_4_alignment_3_run.r` replace usage of `MAFFT` by another aligner like e.g. `MUSCLE` or `Clustal`.
- 4 Think about parameters for functions `deleteGaps()`, `del.rowgaponly()` and `del.colgaponly()`. How do they influence output (`*.aln.fasta` files)?

## After all alignment jobs are done

- Results (alignments named `*.aln.fasta` and other files) are in newly created `aligned` directory created by `hybseq_4_alignment_1_submitter.sh` in `DATADIR`
- Other outputs are images with alignment checks (`*.aln.check.png` and `*.aln.png`), NJ trees (`*.nwk` and `*.tree.png`), saturation plots (`*.saturation.png`) and logs (`*.log` from `R` and `HybSeq.alignment.*.[eo]*` from `qsub`)
- Outputs should be sorted by `hybseq_4_alignment_4_postprocess.sh` – it requires two arguments: directory to process (`aligned`) and path to list of samples (i.e. `samples_list.txt` in `2_seqs`)

## Sorting data after running alignments for all samples

- Outputs of alignments are in directory `aligned` which was created in the input directory, files with alignments itself are named `*.aln.fasta`
- It should be sorted by `hybseq_4_alignment_4_postprocess.sh`
- All outputs should be then moved to dedicated directory (cf. slide 21)
- **Code below is exemplary and requires edits**

```
1 # Move qsub logs to the 'aligned' directory - go to directory XXX/2_seqs
2 mv HybSeq.alignment.* aligned/
3 # Everything should be moved from XXX/2_seqs/aligned to XXX/3_aligned
4 mv aligned ../3_aligned && cd ../
5 # Post-process (sort into subdirectories and get simple statistics)
6 # all alignments - provide path to directory with aligned files
7 ~/hybseq/bin/hybseq_4_alignment_4_postprocess.sh -p 3_aligned \
8   -s 2_seqs/samples_list.txt | tee hybseq_align_postprocess.log
9 mv hybseq_align_postprocess.log 3_aligned/ && cd 3_aligned/
```

# Post-processing of alignments

## Tasks to post-process alignments

- 1 Inspect `~/hybseq/bin/hybseq_4_alignment_4_postprocess.sh` and be sure to understand what the script does, including syntax used.
  - 2 Run `~/hybseq/bin/hybseq_4_alignment_4_postprocess.sh` with correct paths to post-process all aligned outputs (previous slide).
  - 3 Inspect outputs, including statistics, images and logs. Open in spreadsheet (e.g. LibreOffice Calc) the `*.tsv` tables. What do they show?
- Wait for tasks to finish, or **download** and inspect test data

# Gene trees from all alignments I

- Gene trees must be computed from all aligned sequences
- Gene trees must be post-processed
  - Trees should be sorted into subdirectories for exons, introns and supercontigs and lists of gene trees created
  - Trees with significantly different topology must be identified and inspected (and possibly removed) – this will be later done in `R`
  - Long branches within trees must be identified and respective trees inspected – artificial long branches can be discarded from given trees (e.g. by `TreeShrink`)
- Script `~/hybseq/bin/hybseq_5_gene_trees_1_submitter.sh` goes to directory set by `DATADIR` and for every aligned contig (all `*.aln.fasta` files) in that directory (created in previous step) submits via `qsub` individual task to reconstruct gene tree



## Gene trees from all alignments II

- Script `~/hybseq/bin/hybseq_5_gene_trees_2_qsub.sh` drives submission and manipulation of processing of each individual file — it checks if all required data are provided, copy input files to `SCRATCHDIR`, loads required software modules, runs processing itself by `hybseq_5_gene_trees_3_run.sh` and copy results back
- Script `~/hybseq/bin/hybseq_5_gene_trees_3_run.sh` computes gene tree for every input file with `IQ-TREE 2` (see how easily it can be converted for another tree builder, e.g. `RAXML-NG`)
- Script `~/hybseq/bin/hybseq_5_gene_trees_4_postprocess.sh` is short, can be runned on the front end, it will sort outputs into subdirectories for exons, intron and supercontigs and create lists of maximum likelihood and consensus (after bootstrapping) gene trees

## Submitting gene trees jobs

- To reconstruct gene trees from every aligned and trimmed sequence
- `~/hybseq/bin/hybseq_5_gene_trees_1_submitter.sh` must be edited before running it – change `WORKDIR` and `DATADIR`
- It will start job for every sample, so that output of `qstat` will be very long (3 jobs – respective exon, intron and supercontig – for every of  $\sim 1000$  probes)...
- Larger alignments can take long time to compute – set `walltime` in `hybseq_5_gene_trees_1_submitter.sh` accordingly

```
1 # After edition of WORKDIR and DATADIR run simply
2 # NOTE: There are no further parameters on command line
3 ~/hybseq/bin/hybseq_5_gene_trees_1_submitter.sh
4 # Monitor running $USER's tasks with details
5 qstat -w -n -1 -u $USER # Last column contains machine name
6 # See your processes running the machine (from above list)
7 ssh exec_host "ps ux" # Replace exec_host by hostname!
```

# Running gene trees

## Run gene trees

- 1 Inspect scripts `~/hybseq/bin/hybseq_5_gene_trees_1_submitter.sh`, `~/hybseq/bin/hybseq_5_gene_trees_2_qsub.sh` and `~/hybseq/bin/hybseq_5_gene_trees_3_run.sh` and be sure to understand what they are doing, including syntax used.
- 2 Edit declarations of variables `WORKDIR` and `DATADIR` in `hybseq_run_5_gene_trees_1_submitter.sh` so that they point to correct locations.
- 3 Run it to process all `*.aln.fasta` files.
- 4 Monitor progress of the jobs
- 5 Inspect outputs and log files.

# Other tasks with gene trees

For interested students

## Gene trees tasks for advanced users

- 1 See **IQ-TREE help**, check its parameters in `hybseq_5_gene_trees_3_run.sh` and think about possible changes to fit better your needs.
- 2 Think how to process all samples on single computer. Use `find` to list all `*.aln.fasta` files and pass them to **GNU Parallel**. How would such script look like? What had to be changed?
- 3 Check requirements of **IQ-TREE** and/or another tree builder and think if changing of number of CPU threads and memory would significantly speed up processing of individual files.
- 4 Replace usage of **IQ-TREE** in `hybseq_5_gene_trees_3_run.sh` by another tree builder like e.g. **ExaML**, **MrBayes**, **PhyML** or **RAxML-NG**.

## Quick visualization of gene trees I

- Individual gene trees can be quickly visualized by simple **R** script processing in **for** loop each gene tree input file
- Adjust code below and save it as **mltreepplot.r** in directory with gene trees

```
1 library(ape, lib.loc=~ /hybseq/rpackages")
2 # [1] file.aln.fasta.treefile / file.aln.fasta.raxml.support
3 fnames <- commandArgs(TRUE)
4 tr <- read.tree(file=fnames[1])
5 png(filename=fnames[2], width=1800, height=1800, units="px", bg="white")
6   plot.phylo(x=tr, type="unrooted", edge.color="blue", edge.width=3,
7             cex=1.25)
8   title("Maximum-likelihood phylogenetic tree")
9   nodelabels(text=tr$node.label, frame="none", col="red")
10  edgelabels(text=round(x=tr$edge.length, digits=3), frame="none",
11            col="brown", cex=0.75)
12  add.scale.bar()
13  dev.off()
```

## Quick visualization of gene trees II

- Run following `for` loop to process `R` script from previous slide to get quick visualization of individual gene trees
- Of course, this is very exemplary and especially graphical parameters must be highly adjusted according size of phylogeny etc.

```

1 module add r/4.1.3-gcc-10.2.1-6xt26dl # Load R. Ensure 'ape' is installed.
2 # IQ-TREE
3 for L in *.aln.fasta.treefile; do echo "${L}"; R CMD BATCH --no-save \
4   --no-restore "--args ${L} ${L%.aln.fasta.treefile}.mltree.png" \
5   mltreeplot.r; done
6 # RAXML-NG
7 for L in *.aln.fasta.raxml.support; do echo "${L}"; R CMD BATCH --no-save \
8   --no-restore "--args ${L} ${L%.aln.fasta.raxml.support}.mltree.png" \
9   mltreeplot.r; done
10 # Note > " < around arguments to pass to R

```

## After all gene trees jobs are done

- Results are in newly created `trees` directory created by `hybseq_5_gene_trees_1_submitter.sh` in `DATADIR`
  - Records of what IQ-TREE did (`*.log`) and `HybSeq*` from `qsub`
  - Results are in `*.iqtree` (IQ-TREE report), maximum likelihood tree is in `*.treefile` and likelihood distances in `*.mldist`
  - Ultrafast bootstrap approximation results contain split support values in `*.splits.nex`, consensus tree in `*.contree`, bootstrap trees `*.ufboot` and likelihood mapping plots in `*.ckp.gz`; and more...
- Outputs should be sorted by `hybseq_5_gene_trees_4_postprocess.sh` – it requires only path to the `trees` directory
- Outputs of gene trees are in directory `trees` which was created in the input directory

## Sorting data after running gene trees for all samples

- Results should be sorted by `hybseq_5_gene_trees_4_postprocess.sh`
- All outputs should be then moved to dedicated directory (cf. slide 21)
- Code below is exemplary and requires edits

```
1 # Move qsub logs to the 'trees' directory
2 cd ~/hybseq_course_zingibers/3_aligned/ && mv HybSeq.genetree.* trees/
3 # Everything should be moved from XXX/2_seqs/aligned to XXX/3_aligned
4 mv trees ../4_gene_trees && cd ../
5 # Post-process (sort into subdirectories and get lists of gene trees)
6 # all gene trees - provide path to directory with gene trees files
7 ~/hybseq/bin/hybseq_5_gene_trees_4_postprocess.sh 4_gene_trees | tee \
8   hybseq_gene_trees_postprocess.log
```

- Wait for tasks to finish, or **download** and inspect test data



# Post-processing of gene trees

## Tasks to post-process gene trees

- 1 Inspect `~/hybseq/bin/hybseq_5_gene_trees_4_postprocess.sh` and be sure to understand what the script does, including syntax used.
  - 2 Run `~/hybseq/bin/hybseq_5_gene_trees_4_postprocess.sh` with correct path to post-process all gene trees outputs.
  - 3 Inspect outputs of, including logs. What do they show?
- Now we have gene trees for all alignments — we can directly construct species tree (using e.g. *ASTRAL-III*, *ASTER\** or *ASTRID-2* from lists of trees `trees_ml_*.nwk`), or inspect and compare topologies of gene trees and possibly discard outliers

```
1 # Some SW like ASTRAL dislike tree names on the beginning of lines
2 # - discard them e.g. by (i.e. keeping only tree topology):
3 sed -i 's/^[[:graph:]]\+ //' trees_*.nwk
```

## Seeing trees in forest

- Comparison of gene trees start with identifying trees with significantly different topology
- There are several distance matrices allowing compare topological differences among trees (and subsequently plot heatmap, PCoA, etc.)

### How to recognize artifact and real biological feature?

- Without good reference genome it is hard to tell if long branches, weird topologies, etc. are some artifacts or biological reality...
- Problems commonly start with low-quality DNA in lab and subsequent high number of missing data
- Statistically, most of “weird” gene trees topologies are rather from technical issues, so that most of people filter them out...
- Results can vary according to strictness with trimming raw FASTQ, sensitivity of various HybPiper settings, settings of aligner and tree builder...

## Install needed R packages

- The tasks will be done in R (use e.g. RStudio or RKWard)
- Use your notebook or e.g. MetaCentrum OnDemand RStudio or Jupyter notebook
- Following R code is available at [https://github.com/V-Z/hybseq-course/blob/master/trees\\_filtration.r](https://github.com/V-Z/hybseq-course/blob/master/trees_filtration.r) – **download** it and **edit** according to your needs

```
1 # Install R packages
2 install.packages(pkgs=c("ape", "ade4", "distory", "gplots", "ggplot2",
3   "phangorn", "phytools"), repos="https://mirrors.nic.cz/R/",
4   dependencies="Imports")
5 # Install kdetrees package (removed from CRAN)
6 # Ensure package 'devtools' is installed
7 if( ! 'devtools' %in% installed.packages()) {install.packages('devtools')}
8 # Install 'kdetrees' from https://github.com/V-Z/kdetrees Git repository
9 devtools::install_github('V-Z/kdetrees')
```

# Distances comparing trees I

## Single number to compare each pair of complex topologies?

- To compare topology of trees, we need some appropriate distance matrix
  - There is no general agreement which is the best, all have issues...
  - If the distance matrix is not **Euclidean**, we run into another issues...
- 
- Download e.g. `trees_ml_exons.nwk` (or another final list of gene trees) and work in **R** in your notebook
    - Comparing plenty of individual gene trees, finding different topologies, construction of consensual species tree topology
    - Robinsons-Foulds distance in `phytools::multiRF`
      - The index adds 1 for each difference between pair of trees
      - Well defined only for fully bifurcating trees – if not fulfilled, some results might be misleading
      - Allow comparison of trees created by different methods

## Distances comparing trees II

- If the difference is very close to root, RF value can be large, even there are not much differences in the tree at all – `dist.multiPhylo` from package `distory` can be an alternative, although interpretation of that geodesic distance is sometimes not so straightforward as simple logic of RF
- Methods implemented in `ape::dist.topo` allow comparison of trees with polytomies (`method="PH85"`) or use of squared lengths of internal branches (`method="score"`)
- Final matrices are commonly not **Euclidean** – may be problematic for usage in methods like PCoA (it can show misleading results)
  - Test it with `ade4::is.euclid`, can be scaled (forced to become Euclidean) by functions like `quasieuclid` or `cailliez` in `ade4` – carefully, it can damage meaning of the data
  - We get matrix of pairwise differences among trees (from multiple genes), we need display and analyze it
- Set of tools for identifying discordant phylogenetic trees are e.g. in package `kdetrees`
- Filtered trees (with removed outlying topologies) are input for further species tree reconstruction method

# Loading trees into R

```
1 # Load libraries
2 library(ape)
3 library(ade4)
4 library(disty)
5 library(gplots)
6 library(ggplot2)
7 library(kdetrees)
8 library(phangorn)
9 library(phytools)
10 # Set working directory
11 setwd("~/dokumenty/vyuka/hybseq/") # For example
12 # Load the list of trees
13 trees <- read.tree(file="trees_ml_exons.nwk")
14 trees # See it
15 print(trees, details=TRUE)
```

# Heatmap of topological distances

- There are several heatmap functions, try also at least `heatmap`
- Edit settings to fit your needs and preferences

```
1 # Compute distance of topological similarities
2 # Set number of cores according to your computer
3 trees.d <- dist.topo(x=trees, method="score", mc.cores=4)
4 # Plot the heatmap (package gplots)
5 png(filename="trees_dist.png", width=10000, height=10000)
6 heatmap.2(x=as.matrix(trees.d), Rowv=FALSE, Colv="Rowv",
7           dendrogram="none", symm=TRUE, scale="none", na.rm=TRUE,
8           revC=FALSE, col=rainbow(15), cellnote=as.matrix(trees.d),
9           notecex=1, notecol="white", trace="none",
10          labRow=rownames(as.matrix(trees.d)), labCol=colnames
11          (as.matrix(trees.d)), key=FALSE, main="Correlation
12          matrix of topographical distances")
13 dev.off() # Saves the image
```

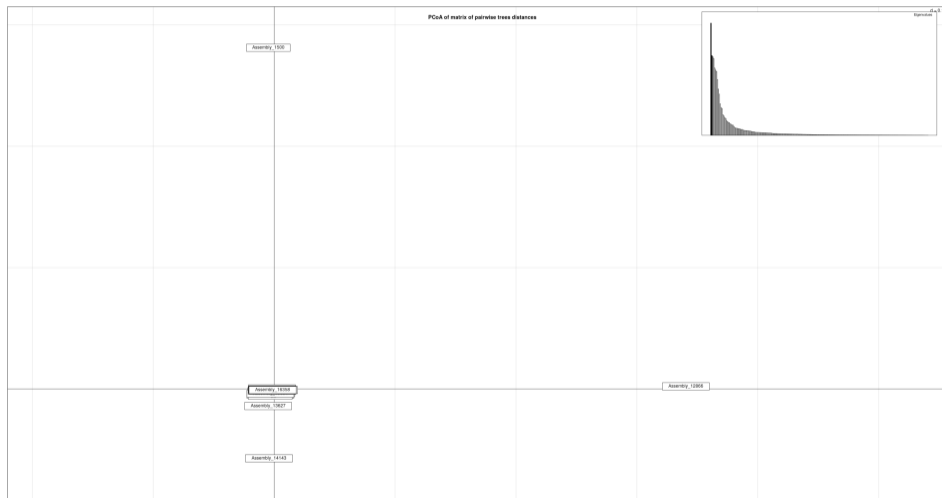
# PCoA of topological distances I

- Requires **Euclidean distance matrix** (`is.euclid()`)
- Non-Euclidean matrices can be forced to become Euclidean by e.g. `quasiEuclid()` or `cailliez()`
- There are plenty of options how to display it

```
1 # Test if the distance matrix is Euclidean
2 is.euclid(distmat=as.dist(trees.d), plot=TRUE, tol=1e-05)
3 # PCoA
4 trees.pcoa <- dudi.pco(d=trees.d, scannf=FALSE, nf=5)
5 trees.pcoa
6 # Plot PCoA
7 s.label(dfxy=trees.pcoa$li)
8 s.kde2d(dfxy=trees.pcoa$li, cpoint=0, add.plot=TRUE)
9 add.scatter.eig(trees.pcoa[["eig"]], 3, 1, 2, posi="topright")
10 title("PCoA of matrix of pairwise trees distances")
```



# PCoA of topological distances II



## Filtration of PCoA

```
1 trees # See original trees
2 # Remove trees identified in the PCoA plot
3 trees[c("Assembly_12866", "Assembly_14143", "Assembly_1500")] <- NULL
4 trees # See new object
5 # Possibly remove trees with too few tips
6 print(trees, details=TRUE)
7 trees[c(1, 2, 3, 4)] <- NULL
8 # Possibly remove rare tips
9 trees <- lapply(X=trees, FUN=drop.tip, tip=c("Amomum-sp7_S308_L001"))
10 class(trees) <- "multiPhylo" # Use after usage of lapply to multiPhylo
```

- Now you can repeat recalculation of distance matrix and PCoA and possibly remove more trees... — or use another method like `kdetrees` (next slide) etc.
- Calculation of distance matrix for large tree set can be very time demanding...
- See more details in <https://soubory.trapa.cz/rcourse/>, chapter “Trees” and subchapters “Seeing trees in forest” and “Comparisons”

# Kdetrees

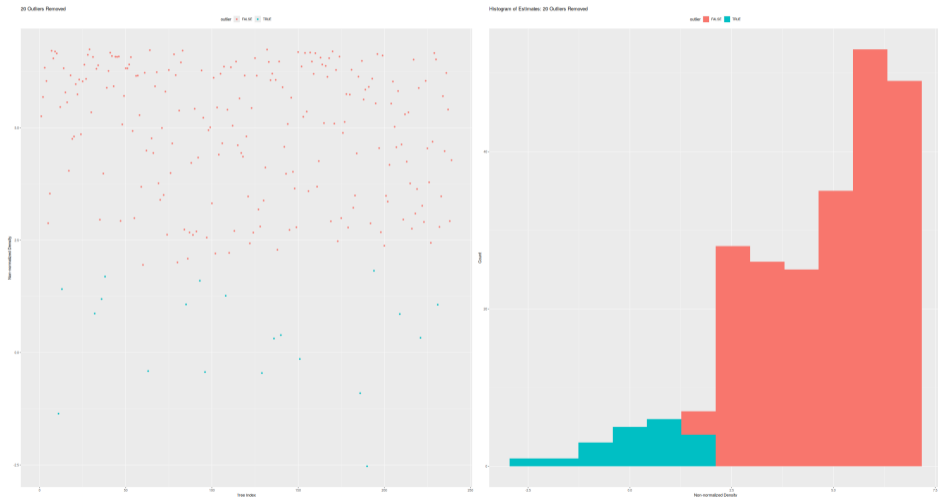
- Finds discordant phylogenetic trees
- Produces relative scores — high are relatively similar to each other, low dissimilar (discordant with the others)
- Produces scores, list of passing/discarded trees and graphical outputs
- In `kdetrees()`, value of `k` is responsible for threshold for removal of outliers — play with it

```
1 # Run kdetrees to detect outliers - play with k
2 ?kdetrees # See options for kdetrees
3 trees.kde <- kdetrees(trees=trees, k=0.9 distance="dissimilarity",
4   topo.only=FALSE, greedy=TRUE)
5 # See text results with list of outlying trees
6 trees.kde
```

## Outputs of kdetrees

```
1 # See graphical results
2 plot(x=trees.kde)
3 hist(x=trees.kde)
4 # See removed trees
5 plot.multiPhylo(trees.kde[["outliers"]])
6 # Save removed trees
7 write.tree(phy=trees.kde[["outliers"]], file="trees_outliers.nwk")
8 # Save kdetrees report
9 write.table(x=as.data.frame(x=trees.kde), file="trees_scores.tsv",
10           quote=FALSE, sep="\t")
11 # Extract passing trees
12 trees.good <- trees[names(trees) %in% names(trees.kde[["outliers"]])
13   == FALSE]
14 trees.good
15 # Save passing trees
16 write.tree(phy=trees.good, file="trees_good.nwk")
```

# Outputs of kdetrees — graphs



# TreeShrink

- Algorithm for detecting abnormally long branches in one or more phylogenetic trees [Mai et Mirarab 2018](#) (see [details and install instructions](#))
- Requires Python 3.X (available on every Linux) and R 4.X to be installed
- Work in command line of your notebook (Linux, macOS), or on MetaCentrum

```
1 # On MetaCentrum load Python and R modules
2 module add python/3.9.12-gcc-10.2.1-rg2lpmk r/4.1.3-gcc-10.2.1-6xt26dl
3 git clone https://github.com/uym2/TreeShrink.git # Download TreeShrink
4 cd TreeShrink/
5 python3 setup.py install --user # Install it "classically"
6 # Go to directory with trees_good.nwk and run TreeShrink
7 ~/.local/bin/run_treeshrink.py -h # See help
8 ~/.local/bin/run_treeshrink.py -t trees_good.nwk -o treeshrink_exons_good \
9   -O treeshrink_exons_good | tee treeshrink.log
10 ls -lha treeshrink_exons/ # Results
```

## Outputs of TreeShrink

- Output (3 files) is saved into directory (in our case) `treeshrink_exons_good`
- File `treeshrink_exons_good.nwk` contains new list of phylogenetic trees in `NEWICK` which can be then used as an input for any species tree reconstruction software (ASTRAL, ...)
- `treeshrink_exons_good.txt` can be bit hard to read – it has one line for every tree in the input list and every line contains list of removed tips – if there is an empty line, no tip was removed from that particular tree; trees are not named, only in same order as in the original input file
- `treeshrink_exons_good_summary.txt` contains statistics for affected taxa

```
1 # Find out how many times particular sample was removed from the list of
2 # the trees
3 grep -o "\<[[:graph:]]+\>" treeshrink_exons_good/treeshrink_exons_good.txt \
4 | sort | uniq -c | sort -gr
```

# Filtration tasks

## Practice filtration of gene trees

- 1 Create PCoA at least from `trees_ml_exons.nwk`. Consider repeated running of PCoA for same input set of trees.
  - Try several different distance matrices (see slide 68). What are differences? Are all appropriate for all cases?
- 2 Run `kdetrees` on list of trees after removal of outliers detected by PCoA and create `trees_good.nwk` (or more similar files).
- 3 Process by TreeShrink at least `trees_ml_exons.nwk` and `trees_good.nwk`. Compare differences.
  - See also slide 65 for differences between both input files, i.e. prior running TreeShrink run:  

```
sed -i 's/^[[:graph:]]\+ //' trees_ml_exons.nwk
```
- 4 What are main differences between PCoA/`kdetrees` and TreeShrink?



## Parsimony super tree

- Parsimony has plenty of implementation, example below is from **R** package **phangorn**
- Usage, principles and accuracy are more or less same...

```
1 # Compute parsimony super tree
2 ?superTree # See help first...
3 tree.sp <- superTree(tree=trees.good, method="NNI", rooted=TRUE,
4   trace=2, start=NULL, multicore=TRUE)
5 tree.sp # See details
6 tree.sp <- root(phy=tree.sp, outgroup=c("Riedelia-arfakensis_S49_L001",
7   "Zingiber-officinale_S242_L001"), resolve.root=TRUE) # Root it
8 # Save parsimony super tree
9 write.tree(phy=tree.sp, file="parsimony_sp_tree.nwk")
10 # Plot parsimony super tree
11 plot.phylo(x=tree.sp, type="phylogram", edge.width=2,
12   label.offset=0.01, cex=1.2)
13 add.scale.bar()
14 # Tune display of the tree...
```

## Other options for species tree estimation in R

- Similar approach as `superTree` is implemented in `phytools::mrp.supertree`
- Distance-based tree reconstruction is in `ape::speciesTree`
- Coalescence model handling multiple individuals per species is in `phangorn::coalSpeciesTree`

```
1 ?ape::speciesTree # See help...
2 ?phytools::mrp.supertree
3 ?phangorn::coalSpeciesTree
4 # All trees must be ultrametric - chronos scale them
5 trees.ultra <- lapply(X=trees.good, FUN=chronos, model="correlated")
6 class(trees.ultra) <- "multiPhylo"
7 # Calculate the species tree
8 tree.sp.mean <- speciesTree(x=trees.ultra, FUN=mean)
9 tree.sp2 <- mrp.supertree(tree=trees.good, method="optim.parsimony",
10   rooted=TRUE)
```

## ASTRAL and related tools

- **ASTRAL-III** and related tools like **ASTER\*** or **ASTRID-2** are popular easy to use methods to construct species tree out of set of gene trees
- From <https://github.com/smirarab/ASTRAL#installation> download **ZIP** archive, unpack **astral.5.7.8.jar** and edit **~/hybseq/bin/hybseq\_6\_sp\_tree\_2\_run.sh** so that **java -jar /storage/.../astral.5.7.8.jar** points to correct location
- **~/hybseq/bin/hybseq\_6\_sp\_tree\_1\_qsub.sh** must be edited before submission via **qsub** – change **WORKDIR** and **DATADIR**
- It goes to **DATADIR** and process all **\*.nwk** files there
- The scripts can be easily edited to use with another similar tool
- Run the task (next slides), or **download** and inspect test data

## Prepare lists of trees for ASTRAL and similar tools and submit job

```
1 # Create new directory and go there
2 cd ~/hybseq_course_zingibers/ && mkdir 5_sp_trees && cd 5_sp_trees/
3 # Copy to 5_sp_trees also results from R and TreeShrink, e.g.
4 cp ../4_gene_trees/trees_*.nwk ../4_gene_trees/treeshrink_*/*.nwk .
5 # Trees lists exported from R and in dir. 4_gene_trees contain gene name
6 # Assembly_10014 (Afr...eum_S118:0.0504,Afr...eta_S701:0.054,((((((Am...
7 # Remove gene (contig) names from lists of trees
8 sed -i 's/^[[:graph:]]\+ //' \
9     trees_{cons,ml}_{exons,introns,supercontigs}.nwk
10 # Edit WORKDIR and DATADIR in ~/hybseq/bin/hybseq_6_sp_tree_1_qsub.sh and
11 # path to ASTRAL in ~/hybseq/bin/hybseq_6_sp_tree_2_run.sh and submit job
12 qsub -l walltime=4:0:0 -l select=1:ncpus=1:mem=4gb:scratch_local=1gb -m \
13     abe ~/hybseq/bin/hybseq_6_sp_tree_1_qsub.sh
14 # Monitor running $USER's tasks with details
15 qstat -w -n -1 -u $USER # Last column contains machine name
16 # See your processes running the machine (from above list)
17 ssh exec_host "ps ux" # Replace exec_host by hostname!
```

# Running species trees

## Run species trees

- 1 Inspect scripts `~/hybseq/bin/hybseq_6_sp_tree_1_qsub.sh` and `~/hybseq/bin/hybseq_6_sp_tree_2_run.sh` and be sure to understand what they are doing, including syntax used.
- 2 Edit declarations of variables `WORKDIR` and `DATADIR` in `hybseq_6_sp_tree_1_qsub.sh` so that they point to correct locations.
- 3 Ensure to have correct path to ASTRAL `jar` file in `hybseq_6_sp_tree_2_run.sh`.
- 4 Run it to process all `*.nwk` files.
- 5 Inspect outputs and log files.

# Other tasks with gene trees

For interested students

## Tasks for advanced users with ASTRAL and relatives

- 1 Think how to run `ASTRAL`, `ASTER*`, `ASTRID-2` or similar tool in your computer.
- 2 Edit `hybseq_6_sp_tree_2_run.sh` (and possibly `hybseq_6_sp_tree_1_qsub.sh`) to use different species tree builder.
- 3 Check parameters of `ASTRAL`, `ASTER*` and `ASTRID-2` and other similar tools listed on their pages. What are differences? What are their advantages and disadvantages? Think about various possible settings.

## Quick visualization of species trees I

- Individual species trees can be quickly visualized by simple **R** script processing in **for** loop each species tree input file
- Adjust code below and save it as **treepLOT.r** in directory with species trees

```
1 library(ape)
2 library(scales)
3 fnames <- commandArgs(TRUE) # [1] file.nwk [2] file.png
4 tr <- read.tree(file=fnames[1])
5 tr <- root(phy=tr, outgroup="Zingiber-officinale_S242_L001",
6   resolve.root=TRUE)
7 tr <- di2multi.phylo(phy=tr, tol=1e-08)
8 tr$edge.length[is.nan(tr$edge.length)] <- 1
9 png(filename=fnames[2], width=1500, height=1500, units="px", bg="white")
10 plot.phylo(x=tr, type="phylogram", edge.color="blue", edge.width=4,
11   cex=1.25, root.edge=TRUE, label.offset=0.1, align.tip.label=TRUE)
12 title("Phylogenetic tree") # Ends on next slide ...
```

## Quick visualization of species trees II

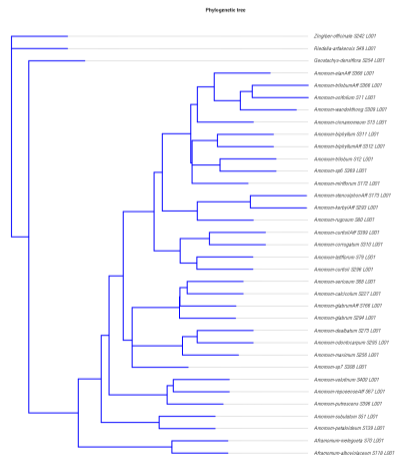
```
1 edgelabels(text=round(x=tr$edge.length, digits=3), frame="rect",  
2   col="brown", bg=alpha(colour="white", alpha=0.5), cex=0.75)  
3 nodelabels(text=round(x=tr$node.label, digits=3), frame="rect",  
4   col="red", bg=alpha(colour="white", alpha=0.5), cex=1.2)  
5 axisPhylo()  
6 dev.off() # ... starts on previous slide
```

- Run following `for` loop to process `R` script from previous slide to get quick visualization of individual species trees
- Of course, this is very exemplary and especially graphical parameters must be highly adjusted according size of phylogeny etc.

```
1 module add r/4.1.3-gcc-10.2.1-6xt26d1 # Load R. Packages: ape, scales  
2 for L in sp_*.nwk; do echo "${L}"; R CMD BATCH --no-save --no-restore \  
3   "--args ${L} ${L%.nwk}.png" treeplot.r; done  
4 # Note > " < around arguments to pass to R
```



# Final species tree



# Consensus network

- Available in R package `phangorn`
- Requires same set of tips in all trees

```
1 # See help
2 ?consensusNet
3 # Compute consensus network
4 tree.net <- consensusNet(obj=trees.good, prob=0.25)
5 # Plot 2D or 3D
6 plot(x=tree.net, planar=FALSE, type="2D", use.edge.length=TRUE,
7      show.tip.label=TRUE, show.edge.label=TRUE, show.node.label=TRUE,
8      show.nodes=TRUE, edge.color="black", tip.color="blue") # 2D
9 plot(x=tree.net, planar=FALSE, type="3D", use.edge.length=TRUE,
10     show.tip.label=TRUE, show.edge.label=TRUE, show.node.label=TRUE,
11     show.nodes=TRUE, edge.color="black", tip.color="blue") # 3D
```

# Phylonet

- Requires as input NEXUS file with **PhyloNet commands** (see below) — e.g. export from **R**:

```
ape::write.nexus(trees.good, file="trees_good.nex", translate=FALSE)
```

- Take **trees\_good.nex** and append following code (or another **command**) to its end:

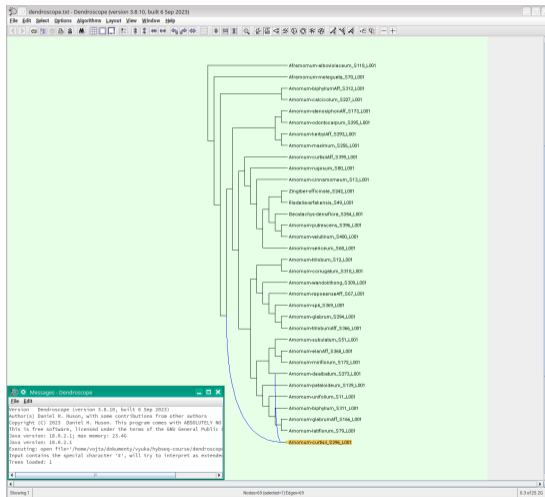
```
1 ... # End of content of trees_good.nex
2 END; # Last line of trees_good.nex - append the PHYLONET block below:
3 BEGIN PHYLONET;
4 InferNetwork_MP (all) 1 -b 50 -x 5 -pl 2 -di;
5 END;
```

```
1 # Download binary JAR file (ready to run)
2 wget https://phylogenomics.rice.edu/media/PhyloNet.jar
3 java -Xmx8g -jar PhyloNet.jar trees_good.nex | tee phylonet_exons.log
```

# Running phylonet

- **PhyloNet** has various options how to create phylogenetic network
- Can be computationally very demanding, running long time — preferably run on MetaCentrum in standard batch job
- The **PHYLONET** section of the input NEXUS contains settings according to **list of commands**
- TreeID can be completely random, or simple consecutive sequence like GT0001–GT####
- PhyloNet can be computationally very demanding, calculating more than 1–3 reticulations can be unrealistic in terms of time needed...
- It **does not save output** file, the network in special NWK format for **Dendroscope** is on the end (after **Visualize in Dendroscope :**) — copy it from terminal or log file (e.g. **tee** from previous slide), extract and save as tiny TXT, which can be opened in **Dendroscope**

# PhyloNet in Dendroscope



# Installing phyparts and other tools

- Requires `maven` and several Python packages, installation can be complicated...

```
1 cd ~/bin/ || { mkdir ~/bin && cd ~/bin/; }
2 # Install Phyparts
3 git clone https://bitbucket.org/blackrim/phyparts.git
4 cd phyparts/
5 # Install dependencies
6 ./mvn_cmdline.sh
7 # Install PhyParts_PieCharts
8 git clone https://github.com/mossmatters/MJPythonNotebooks.git
9 # Or on MetaCentrum (above steps are not needed then)...
10 module add phyparts/0.0.1
11 # Split list of trees into individual files
12 mkdir trees_good
13 split -a 4 -d -l 1 trees_good.nwk trees_good/trees_good_
14 ls trees_good/
15 # If applicable, open parsimony_sp_tree.nwk and remove 'Root' directive
```

## Producing phyparts and phypartspiecharts.py outputs

```
1 # Remove IQTREE ultrafast bootstrap values from gene trees
2 sed -i 's/\/[0-9]\{1,3\}\/g' trees_good/trees_*
3 # Analysis with phyparts
4 java -jar \
5     ~/bin/phyparts/target/phyparts-0.0.1-SNAPSHOT-jar-with-dependencies.jar \
6     --help
7 java -jar \
8     ~/bin/phyparts/target/phyparts-0.0.1-SNAPSHOT-jar-with-dependencies.jar \
9     -a 1 -d trees_good -m parsimony_sp_tree.nwk -o trees_good_res -s 0.5 -v
10 # Copy phypartspiecharts.py to directory with trees
11 cp ~/bin/phyparts/MJPythonNotebooks/phypartspiecharts.py .
12 # See help for phypartspiecharts.py
13 python phypartspiecharts.py --help
14 # Pie chart: concordance (blue) top conflict (green), other conflict
15 # (red), no signal (gray). Run phypartspiecharts.py to get the graphics:
16 python phypartspiecharts.py --svg_name trees_good_res.svg \
17     parsimony_sp_tree.nwk trees_good_res 219
```

## Comparing two trees — cophyloplots

- Slightly different implementation in R packages `ape` (`cophyloplot`) and `phytools` (`cophylo`)
- See help pages and play with graphical parameters

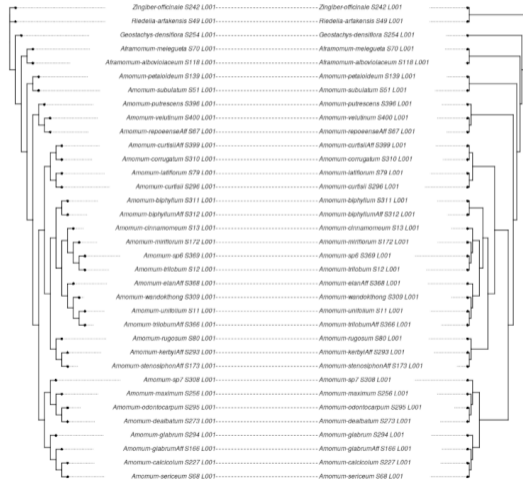
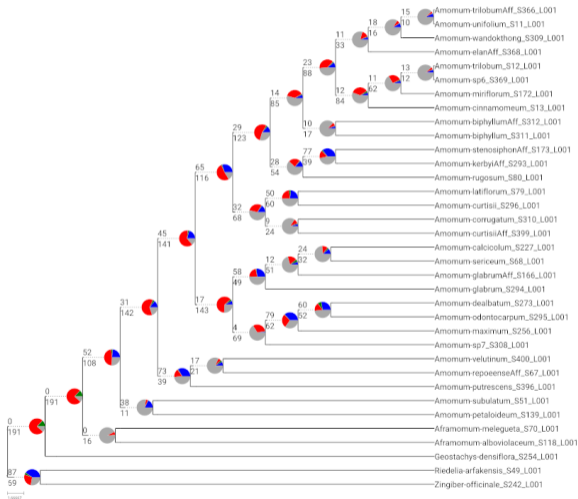
```

1 # We need 2 column matrix with tip labels
2 tips.labels <- matrix(data=c(sort(tree.sp[["tip.label"]]), sort
3   (tree.sp2[["tip.label"]])), nrow=length(tree.sp[["tip.label"]]), ncol=2)
4 # Draw the tree, play with graphical parameters
5 # Click to nodes to rotate them to get better display
6 cophyloplot(x=tree.sp, y=tree.sp2, assoc=tips.labels, use.edge.length=
7   FALSE, space=60, length.line=1, gap=2, type="phylogram", rotate=TRUE,
8   col="red", lwd=1.5, lty=2)
9 # Slightly better display in phytools::cophylo
10 trees.cophylo <- cophylo(tr1=tree.sp, tr2=tree.sp2, assoc=tips.labels,
11   rotate=TRUE)
12 plot.cophylo(x=trees.cophylo, lwd=2, link.type="curved")

```



# Phyparts and cophyloplot



# Density tree

- The trees should be (otherwise plotting works, but may be more ugly) rooted, ultrametric and binary bifurcating
- Implementations are in `phangorn` (`densiTree`) and `phytools` (`densityTree`)

```
1 is.rooted.multiPhylo(trees.ultra) # rooted
2 is.ultrametric.multiPhylo(trees.ultra) # ultrametric
3 is.binary.multiPhylo(trees.ultra) # binary bifurcating
4 # See help page
5 ?phangorn::densiTree
6 # Plotting density trees
7 densiTree(x=trees.ultra, scaleX=TRUE, col=rainbow(6), width=5, cex=1.5)
8 densiTree(x=trees.ultra, direction="upwards", scaleX=TRUE, width=5)
9 densiTree(x=trees.ultra, scaleX=TRUE, width=5, cex=1.5)
10 densiTree(x=trees.ultra[1:10], scaleX=TRUE, width=5, cex=1.25)
```

## Different display for multiple trees

- `phytools::densiTree` requires same number of tips in all trees
- Note various ways how to select trees to display
- Nodes of the trees are not rotated (the display might be suboptimal)
- All trees must have same number of tips

```
1 # See help page
2 ?phytools::densityTree
3 # Plotting density trees
4 densityTree(trees=c(tree.sp, tree.sp2), fix.depth=TRUE, lwd=4)
5 densityTree(trees=trees.ultra, fix.depth=TRUE, use.gradient=TRUE,
6   alpha=0.5, lwd=4)
7 densityTree(trees=trees.ultra[1:3], fix.depth=TRUE, use.gradient=TRUE,
8   alpha=0.5, lwd=4)
9 densityTree(trees=trees.ultra[c(2, 4, 6)], fix.depth=TRUE,
10  use.gradient=TRUE, alpha=0.5, lwd=4)
```

## This is just beginning of long journey...

- I presented common basic tools how to process HybSeq data and deal with various problems, but new tools keep emerging...
- Stay updated and keep trying new tools
- The scripts presented are not the only rigid way how to proceed, rather very general guideline, which should be subject of heavy modifications according to your needs...

# The end

Our course is over...

...I hope it was helpful for You...

...any feedback is welcomed...

...happy playing with the data...

...any final questions?

Typesetting using X<sub>Y</sub>L<sup>A</sup>T<sub>E</sub>X on openSUSE GNU/Linux June 10, 2024