

# Using molecular approaches to understand the drivers of population divergence and speciation.



**Prof. Ian R. Sanders, Dr. Luca Fumagalli, Prof. Nicolas Salamin**

**Assistants:** Dr. Soon-Jae Lee, Dr. Tristan Cumer, Eduard MasCarrio, Julie Guenat, Romain Feron, Tanuj Kafle, Dr. Sara Fonseca Costa, Dr. Daniel Jeffries

# Table of contents

<b>Introduction</b>	<b>3</b>
Objective of this course	3
The “theme” of the course	3
Advice on preparing your report	4
How to use this manual	5
<b>Project 1. Cryptic speciation in tree frogs.</b>	<b>7</b>
Overview	7
Wet-lab protocol	8
Extraction of DNA from froglets.	8
PCR amplification of genetic markers	9
Agarose gel electrophoresis & visualization	12
PCR purification before cycle sequencing reaction	13
Preparation of samples for Sanger sequencing.	15
Data analyses	17
Summary of raw RADseq data processing.	18
Reconstructing the phylogenetic history of Hylid frogs	19
Population genetics in R	22
Plotting in R	23
Population genetics of Hylid frogs in R	24
Population structure analyses of frogs in STRUCTURE	25
<b>Project 2. Genomic analyses of divergence between Lake Malawi cichlids</b>	<b>28</b>
Population genetics of lake Malawi cichlids	29
Detecting candidate loci under selection of lake Malawi cichlids	29
<b>Project 3. The curious case of the European Eel</b>	<b>30</b>
<b>Bibliography</b>	<b>31</b>

# Introduction

## Objective of this course

The main objective of this course is to give you an insight into how researchers use molecular techniques to address ecological and evolutionary questions. Although this course will contain a lot of technical information, **you must always keep the biological question in mind!** If you remember only one thing from this course, remember that (although we hope you remember more!).

As well as a clear question, **a scientist requires a powerful dataset on which to base solid conclusions.** In this practical, you will generate and analyse some data of your own, as well as some datasets we give you. You will learn how to assess a dataset for its suitability for the task, and recognise weaknesses inherent in some approaches. In the process, you will be introduced to several methods, both practical and analytical which are used on a daily basis by ecologists and evolutionary biologists.

## The “theme” of the course

Over this 8 day course, you will work on 3 projects. Each is an example of a natural biological system containing multiple populations / species. Although they come from different corners of the tree of life, **there is a common theme between these systems - they are all subject to ecological and evolutionary forces which promote or hinder population divergence and speciation.** The first (tree frogs) is an example of allopatric divergence (but with a twist), the second (cichlid fish) is a classic example of ecological divergence in sympatry, and the third (European eels) is an interesting case which we leave to you to categorise. You should already have heard about these modes of divergence in your fundamental evolutionary biology courses, but they are also briefly summarised below. Note however, that several more mechanisms of divergence and speciation exist. If only we had time to cover them all

**Allopatric divergence** results from reproductive isolation caused by the separation of previously connected populations of the same species in physical space. For example, it is thought that approximately 2.7-3.5 million years ago (MYA) the Isthmus of Panama closed. In doing so, it broke the connection between the Pacific and the Caribbean seas, meaning many previously connected species were isolated to an extent where gene flow was impossible. The resulting divergence has been studied in many species in this region, including sea urchins, fish and isopods. Note that there are many geographic mechanisms which can lead to isolation, not just the formation of land-bridges.



Source: [https://en.wikipedia.org/wiki/Allopatric\\_speciation](https://en.wikipedia.org/wiki/Allopatric_speciation)



**Ecological divergence** refers to cases of population divergence and speciation where reproductive isolation is not caused by an isolation of species in physical space - i.e. the species remain sympatric throughout the process of divergence speciation. Instead, reproductive isolation is caused by differences which arise in their ecology, allowing them to move into different ecological niches. A good example can be found in the skinks of North America, Richmond & Reeder (2002) found evidence for repeated evolution of increased body size which has allowed some species to occupy higher altitudes and withstand harsh

weather conditions.

### Advice on preparing your report

The marking for this practical course will be based predominantly on your written report, but, for those of you who take part in the wet lab section, a part of your final mark will also be based on your performance in the lab. During this wet lab portion of the course, make sure to keep strict notes of what you did and any auxiliary observations you made throughout so you can include them in your final report when you write it up. Below are some specific guidelines for writing the report but always remember to keep the work you report in the context of the biological question of the project.

The report should be organised according to the project, and the sub-sections should correspond to those found in a traditional scientific paper. I.e. Introduction, Methods, Results and Discussion.

**Introduction:** Here you should set the scene of the project. Describe the reasons for and significance of this study, introduce the biological system, outline the main questions or objectives.

**Methods:** Here you should document *in detail* the approaches you used. For Project 1. This will include both wet lab and dry lab. For those not doing the wet lab, you do not need to include this section in your report, but it is advisable to read through that section of the protocol below to give you an idea of the process preceding the data analyses. When writing the methods, keep the following points in mind: Are the lab protocols complete? It is ok to copy the protocol from the manual (as long as you cite it somewhere), just make your observations/discussions distinct from this (e.g. different colour). Clearly label samples, reagents, volumes, concentrations etc and properly track them through the report. Note the date, time. Be exhaustive in the information you record. Note important observations, including failures during lab work if there were any. Explain the logic of the analyses, i.e. why are you comparing observed and expected heterozygosity,  $F_{st}$ . . . Note all software names and versions. In the first class, there will be considerable emphasis put on why it is so important to clearly document what you did.



**Results:** Here you should thoroughly describe the results of your analyses. Include tables and figures to support the text description of your results. E.g. A table of your samples with relevant information when available like location, calculated population genetics statistics like  $H_e$ ,  $H_o$ . Display all results of analyses and make sure tables and figures are well labeled (e.g. title, legend, x-y axes labeled) and easy to understand? For Project 1. show analyses of data from the whole class, not just the samples you sequenced.

**Discussion:** Here you discuss the implications and significance of your results in the light of the biological questions you outlined in the introduction. Also discuss any relevant technical details, E.g. the differences between phylogenetic trees obtained with nuclear, mitochondrial and RADseq data, bootstrap values, suitability of these markers for our purpose (i.e. be critical of your data / study design). **Important:** Always look in the literature (especially literature mentioned in this manual) and include citations – how does your work fit with what is already known in phylogenetics, population structure in these species etc.

If you have any questions relating to the report, you can write them in this google doc: [https://docs.google.com/document/d/1IodRr2\\_IdI48csH-Em-eitBykXOGl26d9da8uvUHSdU/edit?usp=sharing](https://docs.google.com/document/d/1IodRr2_IdI48csH-Em-eitBykXOGl26d9da8uvUHSdU/edit?usp=sharing). One of the assistants will answer it in the course of the practical so that all students can learn from the answer.

## How to use this manual

In the sections below we will guide you through the wet-lab procedures as with any usual lab protocol of each project. For the analysis sections, although we will provide detailed instructions for some parts, we will not provide you the R code for the population genetic analyses. The reason for this is that, in the real world of molecular ecology and evolution, there is rarely someone to tell you exactly what to do. Thus, an integral part of a successful analysis involves looking around for the best approaches, software, parameter values etc. for your particular needs. So in the sections that follow, we will pose questions, suggest potential functions or approaches that you might use, but it is up to you to find and interpret tutorials and R documentation. Of course you will also have assistants to ask during the TP.

We have included questions throughout the manual - *look for text like this*. These questions are designed partially to help guide you through your analyses and also the report writing, so pay attention to these questions and try to incorporate as many answers as you can into your report. **Note** that it is not so important that you get these questions correct in your report, rather you will be marked on your ability to apply logic, based on the findings of your analyses.

# Schedule

Time	Monday Sept 28th	Tuesday Sept 29th	Wednesday Sept 30th	Thursday Oct 1st	Friday Oct 2nd		Monday Oct 5th	Tuesday Oct 6th	Wednesday Oct 7th	Thursday Oct 8th
8:00-8:45	DNA extraction	PCR Prepare Gel	Sequencing prep	Lecture IS			Lecture LF	Analysis	Analysis	
9:00-9:45										
10:15-11:00	Lecture IS	Lecture IS	Intro to R (JG)	Analysis			Lecture LF	Intro to R		
11:15-12:00										
12:15-13:00			Seminar of DEE	Microeconomics and Game Theory BEE				Seminar of DEE		
13:15-14:00	Introduction TP	Electrophoresis		Lecture IS	Lecture NS	Lecture NS				Analysis
14:15-15:00	DNA extraction Quantification Dilution		Lecture IS		R tutorial		Analysis	Analysis		
15:15-16:00		Lecture IS		Analysis	Analysis					
16:15-17:00										
17:15-18:00		PCR purification								
	= overlap									

= overlap

Lab experiment	BIO 4312
Lectures	POL 334
Analyses	Online

# Project 1. Cryptic speciation in tree frogs.

## Overview

In conservation biology, a fundamental piece of information required for effective management strategies is the distribution of a species and the population structure within that range. Molecular methods have, in the past few decades, proved invaluable for this task, and as technologies have advanced, we have been able to describe subtle distinctions between different lineages and species which were previously not possible.



In this project, we will work on a system of tree frogs from southern Switzerland and Italy. It was traditionally thought that these frogs were a single species, called *Hyla intermedia*. However, some phylogenetic evidence then arose (Stöck et al., 2008) that identified two distinct lineages within the range, which were proposed to be two different species. In a recent paper, researchers from the DEE took a closer look at this system using both traditional and modern genetic approaches (Dufresnes et al., 2018). We will recreate some of this work in order to try to decide if these really are two separate species and if conservation efforts should be made separately for both of them.

In the first few days of the course, we will work in the lab to extract some frog DNA and amplify some genetic markers (one nuclear and one mitochondrial). We will send these for sequencing and later in the course we will analyse this data to produce a phylogenetic tree to see if we can resolve the issue of whether these species are one species or two.

*Q. Why might we use both nuclear and mitochondrial markers in general and specifically in this case?*

In the second part of this project, we will give you some genomic data from the same samples, collected using Restriction Site Associated DNA sequencing (RADseq). We will run some general population genetics analyses and compare this to the sequence phylogeny we made to again see if it helps us determine the status of the frogs in this region.

*Q. What is RADseq, what are its advantages and why might we use RADseq here in addition to the gene sequences.*

## Wet-lab protocol

Ideally for the wet-lab stage of this course you should be working as a **TEAM (not individual)**, so please find a position at a workbench with partners.

### Extraction of DNA from froglets.

Before any genetic work can be done on an organism, we must first extract DNA from a tissue sample. However, when working with wild animals, especially those of conservation concern, an important consideration is how to take a tissue sample in a way that will cause as little damage / stress to the organism as possible. For frogs, the most non-invasive way of sampling is by using a buccal cell swab on the inside of their mouth. Unfortunately, these swabs can only be extracted from once, and so there are no swab samples left for this study. To give you experience of the extraction procedure, we have, therefore, provided you with froglets from a lab-reared clutch. The parents of this clutch were reintroduced to the site where they were collected and so the effect on the natural population was minimal. We will quantify your extractions, and then we will provide you with aliquots of DNA from the original study which we will use for the PCR and sequencing.

**Note: work at the molecular level is highly sensitive to contamination. It is thus necessary to wear gloves at all steps during lab work. Always sanitize your gloves with ethanol if you touch non-sterile material (for example, your glasses, pocket, cell phone, etc.).**

1. You will be given 4 froglets in purified water (deionized distilled H<sub>2</sub>O). These froglets were in long term storage in 70% ethanol at -20°C. However, ethanol can interfere with DNA extractions, so we soak them first in purified water.

For each froglet, remove it from the water and place it on the tissue paper provided. Roll it around to remove excess water and then transfer it to the plastic petri dish. With tweezers and a scalpel, cut off one of the back feet of the froglet, and place it in a 1.5 ml Eppendorf tube. Place the rest of the frog in a separate 1.5 ml tube. **Be sure to assign a unique label to each of your samples, note it in your lab book and write it clearly on the top and side of both tubes. Permanent markers can be removed by ethanol, please be cautious.**

Add 180 µL of Buffer ATL and 20 µL of proteinase K to the extraction tube (the one containing only the foot). Remember to keep the proteinase K on ice as enzymes can be very sensitive to temperature.

2. Mix thoroughly by vortexing, and incubate at 56°C in the thermomixer for 4 hours until the tissue is fully lysed.



**Note: during centrifugation steps, be careful to equilibrate the weight of the tubes in the centrifuge. Ask an assistant if you are unsure how to do this**

3. Collect your sample from the thermomixer and vortex for 15 sec. Centrifuge for a few seconds to remove the drops below the caps. Add 200 µL Buffer AL to the sample, and mix



thoroughly by vortexing and incubate at 56°C for 10 min. Add 200 µL ethanol 100% and mix again thoroughly by vortexing.

4. Pipet the mixture from step 3 into the DNeasy Mini spin column (students should label with the sample ID) placed in a 2 ml collection tube (provided). Centrifuge at 8000 rpm for 1 min. Discard flow-through and collection tube.
5. Place the DNeasy Mini spin column in a new 2 ml collection tube (provided), add 500 µL Buffer AW1, and centrifuge for 1 min at 8000 rpm. Discard flow-through and collection tube.
6. Place the DNeasy Mini spin column in a new 2 ml collection tube and add 500 µL Buffer AW2 and centrifuge for 3 min at 14,000 rpm. Discard-flow through and collection tube. (It is important to dry the membrane of the DNeasy Mini spin column, since residual ethanol may interfere with subsequent reactions. Following the centrifugation step, remove the DNeasy Mini spin column carefully so that the column does not come into contact with the flow-through in the collection tube).
7. Place the DNeasy Mini spin column in a clean 1.5 ml microcentrifuge tube (previously labeled with the sample ID) and carefully pipet 100 µL Buffer AE directly onto the DNeasy membrane (do not touch the membrane with the pipet tip, as this may break the membrane). Incubate at room temperature for 1 min, and centrifuge for 1 min at 8000 rpm to elute.
8. Take your sample to the nanodrop machine, where an assistant will help you check both quality and quantity of your extraction. Note down this concentration, 260/280 ratio and 260/230 ratio. You can then place your sample in a freezer box.

*Q. What does the 260/280 ratio or 260/230 ratio show, respectively? How can you conclude if your DNA extraction was successful or not?*

### PCR amplification of genetic markers

Let us say we want to compare the certain genetic markers from different population, first we need to collect the specific sequences from each of sample. We have extracted entire DNA molecules existing in organism, but we need to purify and amplify the sequences of our interest. How can we do this? The polymerase chain reaction (PCR) is a well-established technique to amplify specific DNA sequences in vitro. The polymerase chain reaction is based on the principle of DNA replication. In a PCR cycle, heat (94-96°C) is used to produce single stranded DNA (denaturation step). After this initial step, the temperature is reduced to ca. 50-65°C such that primers can anneal (annealing step). The annealing temperature is primer-specific. In a third step (elongation step; 68-72°C), the enzyme polymerase replicates the template in 5'-3' direction starting from each of the two primers. Thus, the stretch of DNA that is flanked by the primers is the target of amplification. PCR could be wisely used only after the finding of polymerases derived from those of heat-tolerant microorganisms such as *Thermus aquaticus*, so that the enzymes do not degrade during PCR cycles. Since during PCR also newly synthesized strands act as templates in following cycles, the amplification is exponential.

Therefore, only relatively few cycles (25-40) are necessary to yield enough product for further steps such as DNA sequencing.

Today we will prepare the PCR reactions to amplify regions of the mitochondrial *cytochrome b* gene and the nuclear gene *rag-1* of frogs sampled from northern and central Italy.

**Note: work at the molecular level is highly sensitive to contamination. It is thus necessary to wear gloves at all steps during lab work. Always sanitize your gloves with ethanol if you touched non-sterile material (for example, your glasses, pocket, cell phone, etc.).**

### Cytochrome b (~bp)

1. Prepare a 1.5ml tube for the master-mix, label it (e.g. MM*cyt-b*).
2. Prepare 0.2 ml PCR tubes labeling them with “*cyt-b*” + the DNA sample name.
3. Remove all the necessary reagents from the freezer and **place them on ice**. However, remember that enzymes are temperature-sensitive, so be sure to keep the Taq polymerase on ice at all times.
4. Prepare one master-mix for your PCR reactions using the following table (following the order from top to bottom). Do not forget to account for the volumes needed for both a positive control (which will be provided) and a negative control (replace DNA with water). The final volume of the individual reactions will be 25  $\mu$ L.

*Q. What are the functions of the positive and negative controls here?*

### PCR Master mix for *cyt b* gene

Reagent	Volume ( $\mu$ L) for 1 sample	Volume for (N samples +2 controls) + 1
H <sub>2</sub> O	16.05	
Buffer (5x)	5	
dNTP 25mM	0.2	
<i>cytb-F</i> 10 $\mu$ M	1.25	
<i>cytb-R</i> 10 $\mu$ M	1.25	
Taq Pol.	0.25	
	24	

5. Distribute **24  $\mu$ L master-mix** per well in the prepared PCR tubes.

6. Vortex the DNA extractions shortly and centrifuge them briefly to remove drops from the cap.
7. **Add 1  $\mu$ LDNA** per individual to the corresponding tubes. Close the tubes well. Vortex your PCR reactions, and spin them down. Place your tubes on ice right beside the PCR machine.
8. The machine has been prepared for the following cycling conditions:

Temp (°C )	Time	Cycles
95	2'	
95	30''	38 x
56	1'	
72	1'	
72	10'	
12	8'	

*rag-1* (Size ~ 780bp)

1. Follow the same procedure as used for *cyt b* except use the following volumes for the preparation of the master mix.

Reagent	Volume (µL) for 1 sample	Volume for (N samples +2 controls) + 1
H <sub>2</sub> O	15.05	
Buffer (5x)	5	
dNTP 25mM	0.2	
<i>rag-1-F</i> 10 µM	1.25	
<i>rag-1-R</i> 10 µM	1.25	
Taq Pol.	0.25	
	23	

2. This time, distribute **23 µL master-mix** per well in the prepared PCR tubes.
3. Vortex the DNA extractions shortly and centrifuge them briefly to remove drops from the cap.
4. **Add 2 µL DNA** per individual to the corresponding tubes. Close the tubes well. Vortex your PCR reactions, and spin them down. Place your tubes on ice right beside the PCR machine.
5. The *rag-1* PCR machine will be programmed as follows

Temp (°C )	Time	Cycles
------------	------	--------

95	2'	
95	30''	35 x
56	1'	
72	1'	
72	10'	
12	8'	

### Agarose gel electrophoresis & visualization

We will now check that amplification has been successful using gel electrophoresis, another classical method in molecular biology used on a daily basis. DNA molecules are negatively charged and therefore migrate towards the positive pole in an electric field. In gel electrophoresis, we suspend DNA in agarose gel and apply an electrical current. The agarose provides resistance to the movement of the DNA which is proportional to the size of the molecule, i.e. smaller molecules move faster through the agarose matrix than larger ones. As such, gel electrophoresis can be used to separate DNA strands by size, and to estimate their size through comparison to fragments of known lengths (size standard, or DNA ladder). We visualise DNA fragments in a gel using dyes which bind to the DNA molecules. When exposed to ultraviolet light these dyes will fluoresce. DNA fragments will then become visible as distinct bands in the gel.

**Note: There is a chemical in this experiment that can bind the DNA fragments to visualise it under UV. This means your DNA in your skin is not an exception. It is thus necessary to wear gloves at all steps during lab work.**

### Preparation of the agarose gel

1. Prepare 30 ml of a 1.5% agarose gel solution by dissolving 0.45g agarose in 30 ml 1x TBE buffer. Heat the agarose gel solution in the microwave until the agarose melts and the solution gets liquid and clear. Do not overheat the agarose gel solution as it may evaporate.
2. Let the agarose cool slightly and add 3 µL of GelRed. Mix well.
3. Prepare a gel tray with the respective comb and pour the agarose into the gel tray. Cover with an aluminium foil to protect the GelRed from the light.
4. Before each sample can be loaded into the well in the gel, it has to be **mixed with loading dye** so that it falls to the bottom of the well. Take a piece of parafilm and deposit on it the appropriate number separate drops of 1 µL loading dye each. Take 4 µL of PCR product per reaction and mix with an individual loading dye drop on the piece of parafilm. Make sure you remember which sample is where. **It is important to not make any bubble during the mixing. The bubble will pop during the loading and make the contamination during the electrophoresis.**



5. Carefully remove the comb and check if there is any damage on the gel. Place the tray with the solid agarose gel into a gel chamber filled with running buffer (1x TBE).
6. Carefully load the samples in the well, including 4  $\mu\text{L}$  of DNA ladder in **one of the wells** in the gel.
7. Close the chamber and run the gel for 35 min at 100 V.
8. When the gel has finished running, switch off the power and carefully remove the agarose gel and visualise under UV light. Take a photo and label it. What is the size of your PCR fragments? (See Fig. 2 for the reference sizes of the ladder).

*Q. What can you conclude from the gel, was your PCR successful and, if so, how do you know that you have amplified the correct region?*

### PCR purification before cycle sequencing reaction

To know the sequence differences of target genes among populations we need to sequence the amplified molecules from PCR. Before we can use the products from the PCR in the cycle sequencing reaction, we have to purify them from the remaining dNTPs, primers and salts which can inhibit the sequencing reaction. One way to do this is to run the PCR products on a gel, cut the band to be sequenced, and purify the gel. This method is usually applied when the PCR is not specific enough to exclusively amplify the target gene and instead amplifies other regions as well. In cases where the amplification is specific enough, such as the PCR we performed here, standard commercial kits are applied directly to the PCR product. Here we will use the QIAquick PCR purification kit.

**Note: work at the molecular level is highly sensitive to contamination. It is thus necessary to wear gloves at all steps during lab work. Always sanitize your gloves with ethanol if you touched non-sterile material (for example, your glasses, pocket, cell phone, etc.).**

1. For each successful PCR product of *cyt-b* and *rag1* (positive controls and negative controls DO NOT need to be purified), label a QIAquick column with the sample ID and gene name.
2. For each successful PCR product, prepare one new 1.5 ml tube, labeling it with “Purif” and the sample number and gene name.
3. Add 5 volumes of Buffer PB to 1 volume of the PCR sample and mix well (usually 80  $\mu\text{L}$  of Buffer PB if 16  $\mu\text{L}$  of PCR remain after the gel). Check that the color of the mixture is yellow.
4. Apply the mixture to the QIAquick column and centrifuge for 1 min at 13,000 rpm.
5. **Discard flow-through.** Place the column back into the same tube.

6. To wash, add 750  $\mu$ L of Buffer PE to the column and centrifuge for 1 min at 13,000 rpm.
7. **Discard flow-through.** Place the column back into the same tube. Centrifuge the column for an additional 1 min at 13,000 rpm.
8. Place the columns in the 1.5 ml tubes ("Purif") prepared in step 2.
9. To elute DNA, add 50  $\mu$ L Buffer EB to the center of the column. Avoid touching the membrane with the pipet tip! Let the column stand at room temperature for 1 min and centrifuge for 1 min at 13,000 rpm.
10. Quantify the purified DNA on the Nanodrop.

*Q. What can you conclude from the nanodrop? How did / could you increase the confidence of your concentration estimates?*

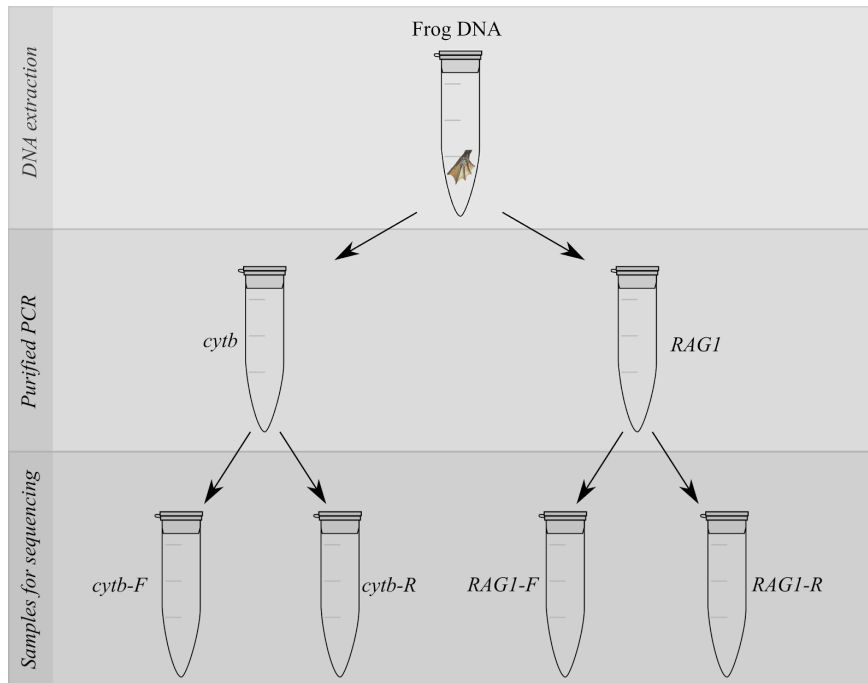
### Preparation of samples for Sanger sequencing.

The cycle sequencing reaction is based on the principle of PCR. Yet, in addition to the dNTPs, fluorescently labeled ddNTPs are added. These nucleotides lack the 3'-OH group that is essential to prolong the DNA strand. As soon as a ddNTP has been incorporated, the synthesis of the fragment is terminated. With this method, strands with different sizes (from primer length plus 1 bp to the entire length of the sequence) are produced, each ending with a fluorescently labeled ddNTP. To allow the determination of the sequence, **only one strand can be sequenced at a time** (i.e., only one primer can be used at a time). After the cycle reaction, all fragments of different sizes are separated by capillary electrophoresis. The identity of the ending ddNTP of each fragment (i.e. either A, T, C or G) is determined thanks to the fluorescence of each, and the sequence of the whole PCR product is reconstructed from the successive bases that are detected by capillary electrophoresis. Cycle sequencing and capillary electrophoresis are performed by commercial labs where we will send our samples.

**Note: work at the molecular level is highly sensitive to contamination. It is thus necessary to wear gloves at all steps during lab work. Always sanitize your gloves with ethanol if you touched non-sterile material (for example, your glasses, pocket, cell phone, etc.).**

### Preparation of samples

1. Prepare 0.2 ml tubes by labeling them with sample IDs. Each PCR product will be sequenced from both sides in separate reactions, i.e. from the forward primer side, and the reverse primer side. Therefore, each PCR product has to be prepared twice so that it can be mixed with both the F and R primers. We will send your samples to the commercial lab for sequencing specifying the individual sample names that you label the tubes with. It is therefore **very important** that you label the tubes with the number of the sample, the name of the gene, i.e. *cyt-b* or *RAG1*, and indicate the primer added, either forward (F) or reverse (R)). So for example, if one of your samples was number 14, label: "14*cyt-b*F".



2. Vortex the PCR purifications and centrifuge them briefly to remove drops from the cap.
3. The amount of DNA added to the sequencing reaction depends on the length of the PCR product. The volume of PCR product needed is 12  $\mu\text{L}$ . The required concentration of DNA in this volume is

$$\text{ng}/\mu\text{L of DNA in reaction} = 1.5 \times (\text{size PCR product in bp}/100)$$

4. Prepare each purified PCR by diluting it with ddH<sub>2</sub>O to achieve the correct concentration and volume. If the volume of PCR product is too small to pipet accurately (e.g. less than 0.5  $\mu\text{L}$ ), dilute your PCR product first.
5. Add the corresponding F or R primer to each tube: add 3  $\mu\text{L}$  Primer at 10  $\mu\text{M}$  to reach a final concentration of 2  $\mu\text{M}$ . The total volume of the sample is 15  $\mu\text{L}$ .
6. Vortex your samples briefly, and spin them down. Leave them on the rack prepared by your assistant.

## Data analyses

The first thing you must do is to **download all the data and code for this course** onto your computer. To do that go here: [https://github.com/cumtr/MolGen\\_2020](https://github.com/cumtr/MolGen_2020) and click on **Clone or Download >> Download ZIP**. Save the folder to your desktop and uncompress if needed. You should then have all of the course materials in the directory “C:/Users/yourusername/Desktop/MolGen\_2020-master/”. You can use this path with the `setwd()` function when working in R so you can easily find all of the files you need.

**Sanger sequence analysis:** Before we can use the sequences we produced for analysis, we had to quality check them and correct/remove errors. We aligned the forward and reverse sequences for each individual and checked them for congruence (note that this step controls for sequence reading errors, but not errors occurring during PCR). Once this was done, we built a consensus sequence for each individual. We will have to align all sequences such that homologous positions among individuals are identified. We will perform this step online, with the freely available program Clustal Omega (<https://www.ebi.ac.uk/Tools/msa/clustalo/>). Finally, using R, we will reconstruct the phylogenetic relationships among the samples.

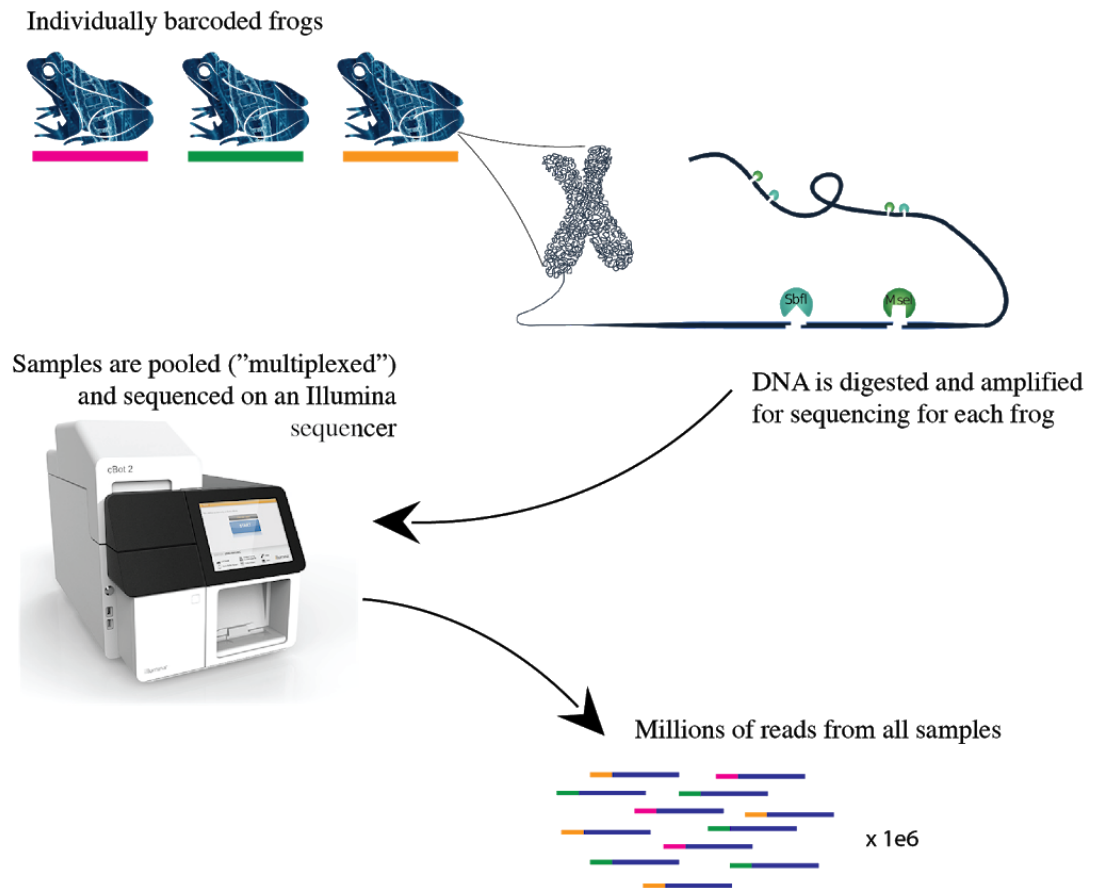
**RADseq analysis:** The initial steps in the RAD marker analysis involve a bioinformatics pipeline which can take several days to complete. As such you will not be asked to do these steps, however you should be aware of how they work. Below, and in the graphical summary on the next page, we outline the major steps in producing RADseq data.

The RADseq library is made in the lab by first digesting (cutting) the genome of each sample with a specific combination of enzymes. As enzymes cut in specific places in the genome, they will generally **cut in the same place across all samples**. We then **incorporate a sample-specific barcode** to the resulting fragments. All fragments from all samples are then pooled together and sequenced using a high throughput sequencing technology like Illumina HiSeq.

The data which is produced by the sequencer contains a mix of all of the pooled (multiplexed) samples, thus the first step of the data processing is to “demultiplex” the data, which means to split the sequences into subsets corresponding to the individual they come from. This step was performed using the program Stacks (<http://creskolab.uoregon.edu/stacks/>). Stacks contain several steps required for raw data processing, we will not describe them all here but encourage you to look at the pipeline yourselves and familiarise yourself with them.

One of the major outputs of the Stacks pipeline is a “VCF” file, which stands for “Variant Call Format”. This file contains all confirmed variants (e.g. Single Nucleotide Polymorphisms, SNPs) for all samples in the dataset. Once in VCF format, the polymorphism data can be reformatted into one of several commonly used data formats. Here we will give you the data in “FSTAT” or GENETOP formats, which are simply named after the software whose developers originally invented the formats. These formats are most easy to load into the R packages we will be using for most of our analyses.



Summary of raw RADseq data processing.

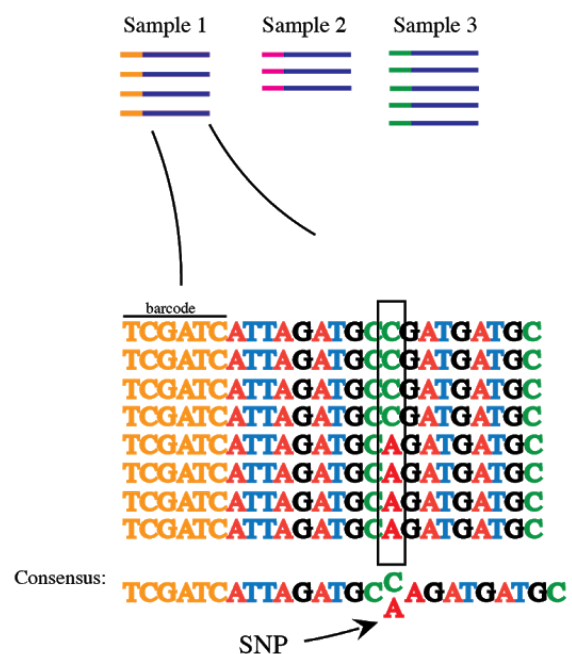
## Stacks

---

Raw sequencing reads are quality checked and "demultiplexed".

There are several other steps here which are required for identifying homology between loci of different individuals

SNP Calling: During Stacks, we also identify variant sites, they are positions where there is a heterozygous position in an individual. This process is done for all samples and the variants are then summarised for all loci and all samples in a .VCF file which stands for "Variant Call Format".



We will perform the majority of our genetics analyses in the computing language R, using the R studio interface. R is one of the most common scripting languages used for statistical analyses of data, especially in population genetics. Hopefully you are already at least a little familiar with R, but don't worry if you find it challenging, everyone does at first - that's what the assistants are for. Just ask us and we'll help you through.

While this course isn't designed specifically to improve your competence in R, a crucial skill in population genetic analyses is the ability to find the most suitable tool to analyse your data, to understand how that tool works and to properly use it to achieve your goals. With this in mind, we will not give you step by step instructions for this part of the course. Rather we will guide you towards the R packages and functions that you will need for each task, and it will be mostly up to you to write the code required. But you are not alone, you can look at the help documentation for any function in R by running it with a `?` in front of it. E.g. `?abline`

Note also that there are **TWO general R classes** scheduled in this course. One (Wednesday 10:15 - 12:00) is organised as part of the wider BEC masters program, and the other (Thursday 15:15 - 18:00) is supplementary for those that feel they need a bit more practice. The second course will specifically be designed for R beginners, although we will cover some more advanced things towards the end.

There are many packages in R which allow for various population genetics tests. For the analyses of the data we will give you, we suggest you use these packages, but you are, of course, free to find more if you wish.

### **APE and PHANGORN:**

Multiple packages may be used to make phylogenetic analyses in R. two of them are APE and PHANGORN, used to manipulate DNA sequences in R, calcul distances between these sequences and infer phylogenetic trees.

Some commands and functions from APE and PHANGORN that you might find useful are listed below but see here <https://www.rdocumentation.org/packages/apex> and here <https://www.rdocumentation.org/packages/phangorn>. Use also help pages of each function for more detailed information:

`read.dna()`: reads DNA sequences in a file, and returns a matrix or a list of DNA sequences with the names of the taxa read in the file as row names or names.

`phyDat()`: transform several DNA formats into the phyDat format.

`dist.ml()`: compute pairwise distances for an object of class *phyDat*.

`NJ()`: performs the neighbor-joining tree estimation.

`root()`: reroots a phylogenetic tree with respect to the specified outgroup or at the node specified.

**HIERFSTAT:**

HIERFSTAT is a package developed for R to study population structure with the help of genetic markers. Wright's F-statistics are commonly used to determine population structure and differentiation. HIERFSTAT allows you to estimate F-statistics and many other basic statistics important for the analyses of population-level genetic data.

Some commands and functions from HIERFSTAT that you might find useful are listed below but see here <https://www.rdocumentation.org/packages/hierfstat> and help pages of each function for more detailed information:

`read.fstat.data()`: import data formatted as a .dat file (the .dat file is the format used by the software FSTAT). The data will be read in as a dataframe.  
`basic.stats()`: calculate basic statistics such as expected heterozygosity, observed heterozygosity, Fis, etc for each population and each locus.  
`varcomp.glob()`: to calculate Fst and Fis  
`boot.vc()`: to calculate a 95% confidence interval for the overall Fst  
`pairwise.neifst()`: to calculate population pairwise Fst  
`indpca()`: to carry out a PCA on the centered matrix of individual allele frequencies

**ADEGENET:**

ADEGENET and HIERFSTAT are built using many of the same libraries and dependencies. There is, therefore, a lot of overlap between them and they complement each other very well. For example ADEGENET also uses the Genind object type and also contains functions for calculating heterozygosities, allele frequencies, Fst etc. However ADEGENET also implements methods which are very useful for spatial genetics, for example examining Isolation By Distance (IBD), and PCA.

Again, some functions that you will find useful are below. There is also a great tutorial for ADEGENET which you can find here <http://adegenet.r-forge.r-project.org/files/tutorial-basics.pdf>. Take a look specifically at the PCA and IBD sections which should prove very helpful.

`read.fstat()`: "Genind" object - a specific object class for this package. Take a look at the structure of the object, you can access slots within it using the `$` accessor, E.g. `mygenindobj$Pop` will show you the information in the Pop slot, which you can edit if you wish.  
`genind2genpop()`: transform a genind object to a genpop object, necessary for some functions  
`dist.genpop()`: Create a Euclidean distance matrix from genetic data  
`dist.geo()`: Create a Euclidean distance matrix from spatial coordinates  
`dist.genpop()`: Create a Euclidean distance matrix from a genpop  
`mantel.randtest()`: Perform a mantel test for isolation by distance

**LEA:**

LEA is an R package dedicated to landscape genomics and ecological association tests. LEA can run analyses of population structure and genome scans for local adaptation. It includes statistical methods

for estimating ancestry coefficients from large genotypic matrices and evaluating the number of ancestral populations (snmf, pca).

Some commands and functions that you might find useful are listed below but see here <https://www.rdocumentation.org/packages/LEA>, and help pages of each function for more detailed information:

**snmf()**: estimates admixture coefficients using sparse Non-Negative Matrix Factorization algorithms, and provides STRUCTURE-like outputs.

**cross.entropy()**: Return the cross-entropy criterion for runs of snmf with K ancestral populations.

### Plotting in R

One of the most powerful aspects in R is its ability to produce scientific plots quickly from the data and results that you are analysing. This is particularly useful in the data exploration phase of a project, but also for making the final plots for reports, theses and papers. Here are some of the most commonly used plotting functions:

**barplot()**: Used to plot 1 dimensional data, e.g. heights of all students in the class.

**plot()**: Used to plot 2 dimensional data. E.g. height vs age. This function may also be used to plot phylogenetic trees.

**heatmap()**: Used to produce a graphical representation of a 2D matrix (e.g. pairwise Fsts).

**abline()**: Adds a line to a plot (e.g. an average **abline(mean(x))**)

**text()**: Adds text to a plot (e.g. useful for adding population names to scatter plots etc.)

**barchart()**: This function displays a bar plot/bar chart representation of the Q-matrix computed from an snmf run.

There are many more plot types and functions, so feel free to explore the possibilities!

## Reconstructing the phylogenetic history of *Hylid* frogs

To reconstruct the phylogenetic relationships between our frog samples, we are going to compute a phylogenetic tree (one for each marker) using the simple, distance-based Neighbour-Joining (NJ) method. In order to root the tree, we will also need one or more outgroup samples which we know are more distantly related to the other frogs in our tree. We will supply you with some sequences from *Hyla arborea* for this purpose. Before estimating the NJ-tree, we have to integrate the outgroup sequence in the sequence alignment.

*Q. Why might we include an outgroup when constructing a phylogenetic tree?*

### Sequences Alignment:

1. Open the .fasta file containing the consensus sequences for *cytb* from the whole class with any text editor software.

The file may be found here : **./MolGen\_2020/1\_Frogs/Data/Phylo/AllSeq.fasta**

2. We have provided you with some representative sequences from other Hylid species. These sequences will firstly act as outgroups to your data, and secondly allow you to compare the divergence between our samples and the divergence between other, already-established species pairs. Import the outgroup sequences which can be found here:

**./MolGen\_2020/1\_Frogs/Data/Phylo/Outgroups/Hyla\_mtDNA\_Dufresnes\_2018.fasta**

Add these outgroup sequences to the fasta file with our sequences.

3. Run a Clustal Omega alignment online and export that alignment as *fasta*. To do so :
  - a. Connect to <https://www.ebi.ac.uk/Tools/msa/clustalo/>
  - b. Specify that your input sequences are DNA sequences.
  - c. load your file with all your sequences and the outgroup sequences.
  - d. select Pearson/FASTA format as output.
  - e. submit your alignment.
  - f. Copy and paste the results in a new *fasta* file.

### Simple phylogeny in R:

With the above R tools, and others in these packages, load the aligned sequences, compute the distance between these sequences and draw a Neighbour-Joining tree.

*Q. What can you say about your samples ? Do all populations form monophyletic groups ?*

*Q. Is this data sufficient or even useful for answering the question at hand?*



### Population genetics of Hylid frogs

With the above R tools, and others in these packages, explore and analyse the frog SNP dataset that can be found here: **`./MolGen_2020/1_Frogs/Data/RADpopgen/hyla_FSTAT.dat`**

We encourage you to explore the data thoroughly:

- Look at within population statistics like observed and expected heterozygosities, inbreeding (Fis) and expectations of Hardy Weinberg.
- Look at between population statistics - pairwise Fsts between populations, patterns of isolation by distance (i.e. a correlation between physical distance and genetic distance between populations (see also Mantel tests)).
- Look at population clustering (i.e. are some populations more closely related to each other than to any others?) using a PCA, or a neighbour-joining tree.

Remember to keep the biological question in mind all the way through - are these frogs one species, or two?

*Q. What do these calculations tell you about the divergence between populations?*

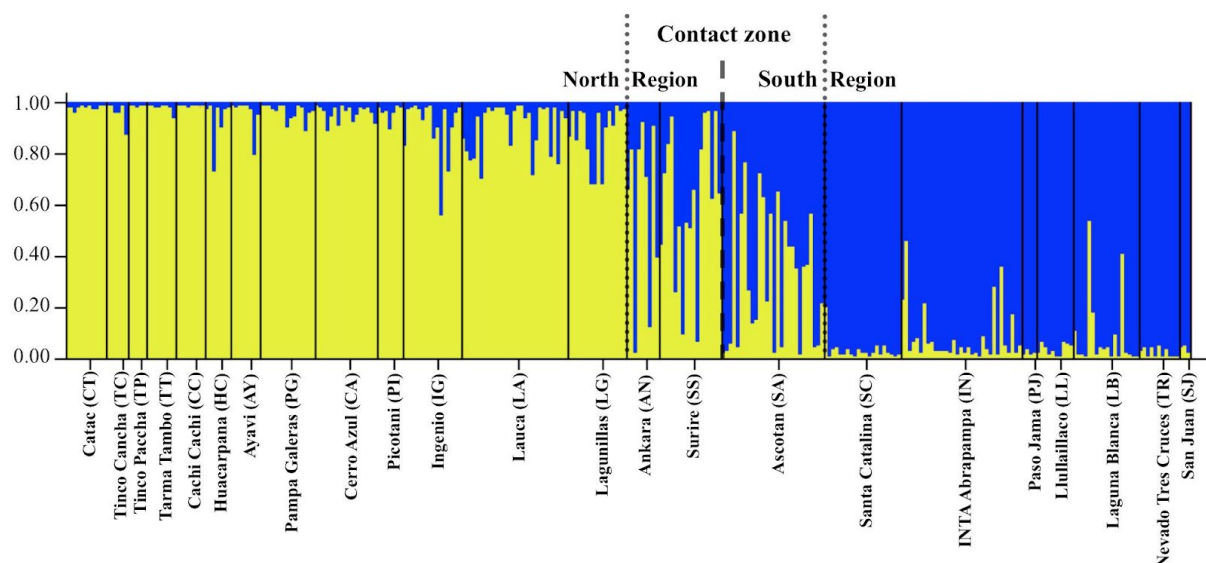
*Q. Are there some populations which are more diverged than others?*

*Q. Do the populations fall into discrete "groups" which might be considered as different species?*

*Q. What was the most likely evolutionary history of these species? Has the hybrid zone always been there? How did reproductive isolation occur?*

## Population structure analyses of frogs

We often talk about organisms sampled from nature as coming from “populations”, however, this is a reasonably arbitrary classification. What is a population - is it a single pond full of frogs? Is it a network of ponds that are close and always connected to each other? Is it a network of ponds which sometimes connect to each other during high rainfall? You can see that drawing a classification like this could be quite meaningless. However, one way to accurately determine a meaningful population genetics unit of organisms is to use an approach called genetic clustering. This is a method which groups individuals based on how closely related they are to each other. The advantage of genetic clustering is that it is not subject to inconsistent human impressions of what constitutes a cohesive population, but simply identifies groups (“clusters”) of individuals which are more closely related to each other than they are to any other sample in the dataset. STRUCTURE is the most widely known program for this task (<http://pritch.bsd.uchicago.edu/structure.html>, Pritchard et al, 2000) but many other tools have been developed since. The idea is to assign individuals to K clusters based on their genotypes, and estimates population allele frequencies. The program incorporates population genetics models, most importantly Hardy-Weinberg equilibrium, into its clustering and so has many advantages over more general multivariate clustering approaches like PCA.



As an example of the use of STRUCTURE (or similar software), this figure above shows individuals of two subspecies of vicuña (*Vicugna vicugna mensalis* and *V. v. vicugna*) which were genotyped at a set of 15 microsatellites (González et al., 2019), and assigned to K=2 clusters. Each vertical bar depicts an individual and is colored to denote the probability that each individual belongs to cluster 1 (yellow) or cluster 2 (blue). Individuals in this study are classified as either *V. v. mensalis* or *V. v. vicugna*, but you can see that the genetic data tells a more complicated story, with no convenient separation between these two subspecies and lots of hybridisation occurring in the contact zone between them.

With the above R tools, and others in these packages, run sNMF for K=2 for your samples. Use the results to plot a barchart of individuals. The dataset can be found here:

**`./MolGen_2020/1_Frogs/Data/RADstructure/Hyla.struct.geno`**

*Q. What do the sNMF results tell you about the distribution of diversity between these populations?*

*Q. What is "K" and how can you use this to help infer whether these frogs fall into two species or not?*

*Q. Do you think these two lineages should be considered separate species and why? Whatever your conclusion, remember to logically back up your decision!*

## Project 2. Genomic analyses of divergence between Lake Malawi cichlids

Cichlids are an extremely diverse group which make up over 10% of all Teleost (ray-finned) fish. The most famous are those found in the lakes of Africa, which have become textbook models for the study of speciation due to the extremely rapid radiation of thousands of species in only 1-2 Million years (Kocher, 2004). They are also extremely important models for the study of **ecological speciation**, as most of the divergence in this group has arisen while they were still in sympatry.



In this project we will analyse some data from Hahn et al. (2017) who studied the divergence between four species of cichlid. Not only did they do this at a population scale, but also on the genomic scale, that is they identified the regions of the genome which were most diverged between the species as a way of generating hypotheses about the mechanism by which they speciated.

We have provided you with two files. The first contains RADseq data similar to that which you analysed for Project 1. The second dataset is perhaps the most interesting and will give you an insight to the bleeding edge of population genomics - i.e. how do **different parts of the genome differ in their evolutionary history**. This dataset contains  $F_{st}$  which has been calculated **per locus**. I.e. for each locus, an  $F_{st}$  has been calculated which represents the difference in allele frequency between populations at that locus alone. Loci which are near (and so tightly linked) to loci under selection will show specific patterns of  $F_{st}$  which give clues to the type and direction of selection. For example, regions of the genome containing many loci with elevated  $F_{st}$  (relative to the genome average) are likely under diversifying selection. Plotting these values along the genome is therefore a powerful way to identify candidate regions/genes driving evolutionary processes such as speciation.

### Population genetics of lake Malawi cichlids

For the first part of this project, you will analyse and explore a SNP dataset from several cichlid species (find the data here: `./MolGen_2020/2_Cichlids/Data/cichlid_hierfstat.dat.`) Use your new found expertise in R to:

- Look at within population statistics like observed and expected heterozygosities, inbreeding (Fis) and expectations of Hardy Weinberg. Plot the mean observed and expected heterozygosities for the 4 populations against each other.
- Look at population pairwise Fst. Try to use the basic R plotting function `heatmap()` with the matrix you obtain using HIERFSTAT.
- Look at population clustering using a PCA.

*Q. Do you observe any difference in observed and expected heterozygosities for the 4 populations?*

*Q. Considering that these 4 species live together in the same lake, how can you interpret the populations pairwise Fst values you obtained? Which can be the biological explanation and evolutionary mechanism that drives this pattern?*

*Q. Are there some populations which are more diverged than others? If so can you hypothesise as to why?*

*Q. Do the individuals of the same populations fall into discrete "clusters"?*

### Detecting candidate loci under selection of lake Malawi cichlids

In the second part of this project you will identify candidate loci under selection by using a dataset of pairwise locus Fst values obtained using 2 out of the 4 population that can be found here: `./MolGen_2020/2_Cichlids/Data/Di_1-Di_2_global_mod.tsv`

- Look at the content of the file and try to understand the different information it provides.
- Generate a "Manhattan" plot using this data (i.e. plot Fst against the location in the genome)
- Use a statistical test to detect possible outliers.

*Q. Can you visually identify the presence and number of candidate loci under selection? If yes, in which chromosomes are located?*

*Q. How many of these loci are statistically identified as actual outliers?*

*Q. How would you use these loci to better understand the evolutionary mechanism underlying the observed genetic differentiation?*



## Project 3. The curious case of the European Eel

The final dataset we will give you is one from the European Eel, *Anguilla anguilla*. By this point you should be well acquainted with the techniques needed to explore a genetic dataset. Thus, we will not spoil the fun by telling you the story behind this system. Examine population diversity and especially differentiation between populations (i.e.  $F_{st}$ ). Compare the levels of divergence and the geographic distances with those of the previous two datasets and **try to come up with an explanation for what you see**. For clues, you can read the paper from which this data was taken (Pujolar et al., 2014). This final task represents the final step of any scientific project - once you have found an interesting result, you must go looking through the literature to see how it compares to what others have found.



# Bibliography

- Dufresnes, C., Mazepa, G., Rodrigues, N., Brelsford, A., Litvinchuk, S. N., Sermier, R., et al. (2018). Genomic Evidence for Cryptic Speciation in Tree Frogs From the Apennine Peninsula, With Description of *Hyla perrini* sp. nov. *Frontiers in Ecology and Evolution* 6.
- González, B. A., Vásquez, J. P., Gómez-Uchida, D., Cortés, J., Rivera, R., Aravena, N., et al. (2019). Phylogeography and Population Genetics of *Vicugna vicugna*: Evolution in the Arid Andean High Plateau. *Front. Genet.* 10, 445.
- Hahn, C., Genner, M. J., Turner, G. F., and Joyce, D. A. (2017). The genomic basis of cichlid fish adaptation within the deepwater “twilight zone” of Lake Malawi. *Evol Lett* 1, 184–198.
- Kocher, T. D. (2004). Adaptive evolution and explosive speciation: the cichlid fish model. *Nat. Rev. Genet.* 5, 288–298.
- Pujolar, J. M., Jacobsen, M. W., Als, T. D., Frydenberg, J., Munch, K., Jónsson, B., et al. (2014). Genome-wide single-generation signatures of local selection in the panmictic European eel. *Mol. Ecol.* 23, 2514–2528.
- Richmond, J. Q., and Reeder, T. W. (2002). Evidence for parallel ecological speciation in scincid lizards of the *Eumeces skiltonianus* species group (Squamata: Scincidae). *Evolution* 56, 1498–1513.
- Stöck, M., Dubey, S., Klütsch, C., Litvinchuk, S. N., Scheidt, U., and Perrin, N. (2008). Mitochondrial and nuclear phylogeny of circum-Mediterranean tree frogs from the *Hyla arborea* group. *Mol. Phylogenet. Evol.* 49, 1019–1024.