# Phylogenetic and Phylogeographic meta-analysis of Cytochrome c Oxidase I barcode sequences of East African arthropods submitted into the Barcode of Life Database

**Presented by:**

Gilbert Kibet
HSB316-5106/2016

**Supervisors:**

Dr Scott Miller
Dr Jandouwe Villinger
Dr Steven Ger

**Collaborators:**

Dr Caleb Kipkurui
Dr Jean-Baka Domelevo
Dr Daniel Masiga

# Background

**Identification and classification of organisms:**

**Morphology-based identification systems** – extensive information (ecology, anatomy, physiology); expensive, slow and needs expertise

**Molecular-based system** – efficient (fast and effective); dependent on reference libraries of DNA barcode -*short and standardized genes or regions thereof used in identification and discovery of species*

The Consortium for the Barcode of Life (CBOL), May 2004: To aid rapid and inexpensive identification of millions of species using DNA barcodes

- **International Nucleotide Sequence Database Collaborative (INSDC):** GenBank, the European Molecular Biology Lab in Europe, and the DNA Data Bank of Japan

- **Barcode of Life Database (BOLD):** University of Guelph in Ontario

# Background

A **_658 base-pair 5' region of mitochondrial cytochrome c oxidase subunit I (COI/COXI) gene_** is the standard the barcode for <u>most</u> animal groups
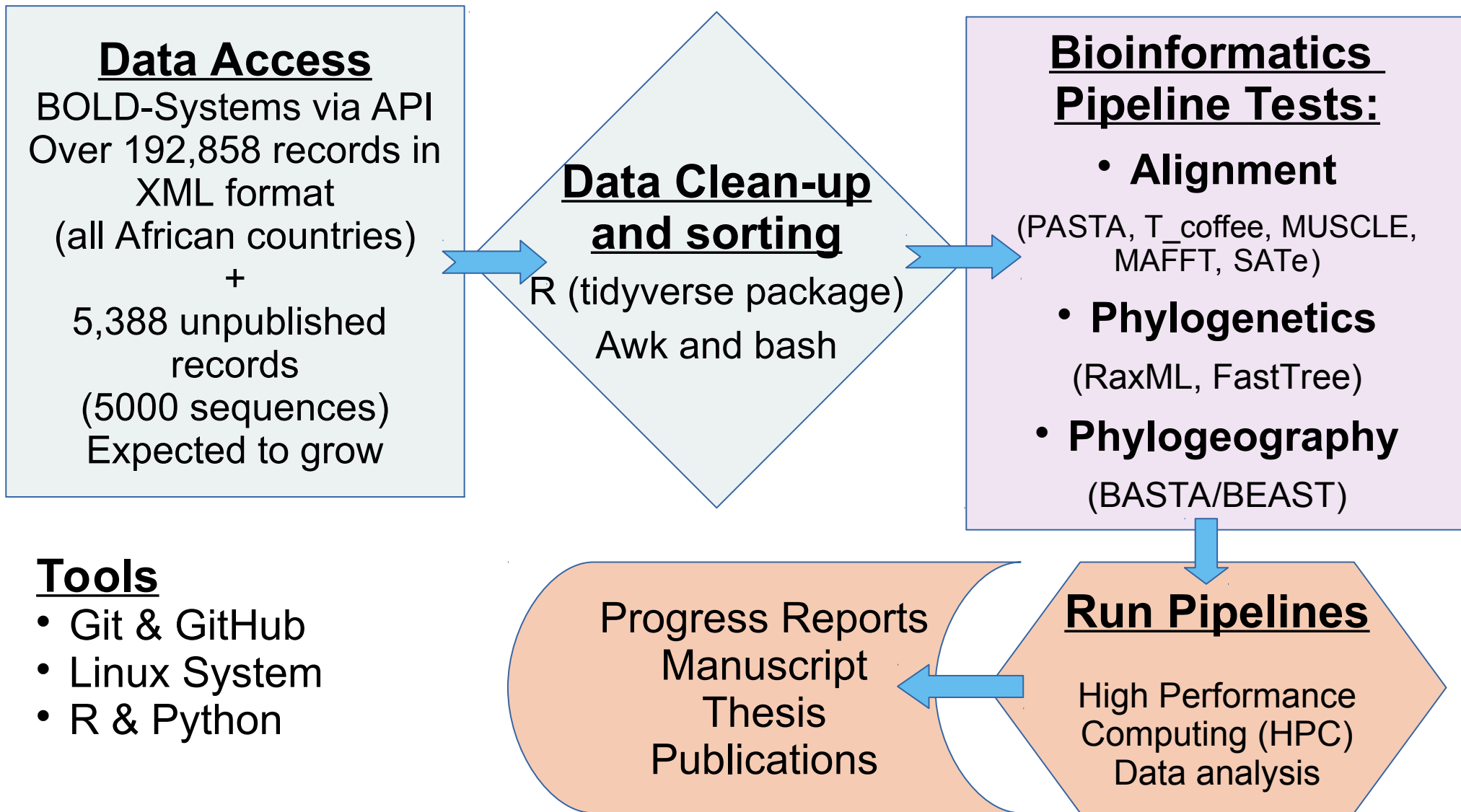
**Problem statement:**

Thousands of East African COI sequences from voucher arthropods submitted into the BOLD database and are yet to be analysed comprehensively: *__phylogenetic__* **diversity** and **<u>phylogeographic distribution</u>**

**Objectives:**

- Improve phylogeographic and *phylogenetic* diversity descriptions of arthropod in East Africa.

- Identify the cryptic species that may not yet be recognized and may be potential crop pests or vectors of human and animal diseases.

# Workflow

**Data Access**
BOLD-Systems via API
Over 192,858 records in XML format
(all African countries)
+
5,388 unpublished records
(5000 sequences)
Expected to grow

**Data Clean-up and sorting**
R (tidyverse package)
Awk and bash

**Bioinformatics Pipeline Tests:**
- **Alignment**

(PASTA, T_coffee, MUSCLE, MAFFT, SATe)

- **Phylogenetics**

(RaxML, FastTree)

- **Phylogeography**

(BASTA/BEAST)

**Tools**
- Git & GitHub
- Linux System
- R & Python

Progress Reports
Manuscript
Thesis
Publications

**Run Pipelines**

High Performance Computing (HPC)
Data analysis

# Data Mining and Cleansing

**<u>Automated Sequence retrival:</u>**

From BOLD-Systems via a URL based API using wget; XML files: Bash

Parsed using BeautifulSoup4 and lxml converted to 80 column text files (.tsv) with pandas: Python3

# <u>Data clean up:</u> R (tidyverse package), Awk and BASH

| Metadata (80 columns) | | | | |
|---|---|---|---|---|
| 'COI-5P' = 198148 (43484 East African (1 sample)) | | | | |
| #nucleotides / unaligned seqs / #ns | | | | |
| Over 700 (1607) | Under 500 (6715) | 500 -700 (184855) | 650 -660 (99698) | Over 499 (186458) |
| 1 sample | 1 sample | 3 samples | 2 samples | |
| Build.fasta: >ProcessID\|order\|genus\|species\|sub_species\|country\|exactsite\|lat\|lon\|elev\|seq_len AGGTTCATCCCAA----- | | | | |

# Multiple Sequence Alignment

## Large dataset:

- Accuracy

- Speed

## Algorithms:

- Progressive (mafft/muscle)

- Progressive & transivity (pasta)

- Regressive (T_coffee)

| MUSCLE: default | MAFFT: --large G-INS-1 | T_coffee: -reg | PASTA: default | SATe |
|---|---|---|---|---|
| Fast speed | Fast speed | Fast speed | Fast speed | NA |
| Low accuracy | High accuracy | High accuracy | High accuracy | NA |
| • Refine<br>• Align<br>• Merge | • Align<br>• Add sequence<br>• Add_fragments<br>• Merge | • Align<br>• Evaluate: (CORE index TCS) | • Align<br><br>• Add fragments | NA |

# Multiple Sequence Alignment evaluation

**Evaluation for Accuracy:**

T_coffee: **consistency based scoring**

- CORE index (html)
- Transitive Consistency Scores (TSC) (html/ascii)

Purpose:

Used to select the most suitable alignments.

TCS ascii to used in applying different weights to columns in phylogenetic analysis

## Visualization:

- Seaview

- Jalview

- SuiteMSA_Package1.3.22B

# T_coffee consistency based
# Multiple Sequence Alignment evaluation

## MAFFT



## MUSCLE



## T_coffee



## PASTA

# Phylogenetic analysis: Basic information



enafroCOI_all_data.tre

Nodes: 386351

External nodes: 193176

Internal nodes: 193175

Branches: 386351

Depth: 518

Maximum distance to root: 10.26635

Archaeopteryx visualization

# Current progress...

What next?

➜ Studying/analysing the phylogenetic tree using: RAxML8, FastTree and T_Coffee

➜ Phylogeographic analysis

Points to Note:

➜ Open science principles: Collaborative supervision through GitHub, Open source software, and aspire to make scripts/codes transparent, available, free and accessible

# Timeline

| Activity | Time in months (2018-2019) | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Sept | Oct | Nov | Dec | Jan | Feb | March | April | May | June | July | Aug |
| Proposal writing and Literature Review | █ | █ | █ | | | | | | | | | |
| Data Mining and Sorting | | | █ | █ | █ | █ | | | | | | |
| Pipeline Development and Testing | | | | █ | █ | █ | █ | █ | | | | |
| Data Analysis on HPC | | | | | | | | █ | █ | █ | █ | |
| Manuscript Writing and submission | | | | | | | | | | █ | █ | █ |
| Thesis writing and Defence | | | | | | | | | | █ | █ | █ |

# Acknowledgement

## Thank you for enabling a bioinformatics dream