# PHYLOGENETIC AND PHYLOGEOGRAPHIC META-ANALYSIS OF COI SEQUENCES OF EAST AFRICAN ARTHROPODS IN THE BOLD DATABASE

## Objective:

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

  Identify the diversity of species that may not yet be recognized as potential crop pests or vectors of human and animal diseases

- Allow for better assessments of emerging threat to human, animal and crop health

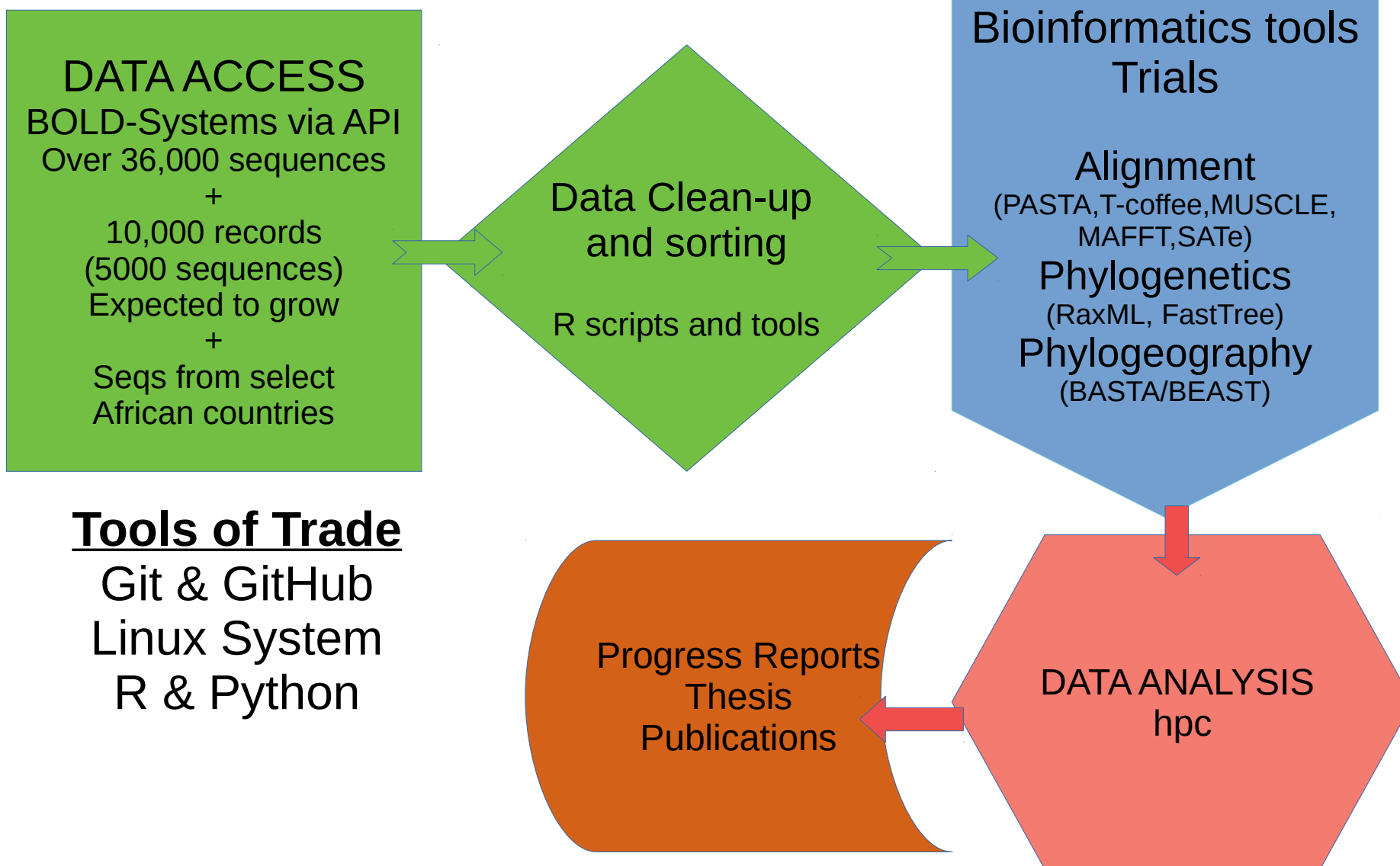## Kibet Gilbert (DRIP fellow)

- **Supervisors:**

  Dr Scott Miller, Dr Jandouwe Villinger, Dr Steven Ger

- Collaborators:

  Dr Caleb Kipkurui, Dr Jean-Baka, Dr Dan Masiga

# workflow

**DATA ACCESS**
BOLD-Systems via API
Over 36,000 sequences
+
10,000 records
(5000 sequences)
Expected to grow
+
Seqs from select
African countries

**Data Clean-up and sorting**

R scripts and tools

**Bioinformatics tools Trials**

**Alignment**
(PASTA,T-coffee,MUSCLE, MAFFT,SATe)
**Phylogenetics**
(RaxML, FastTree)
**Phylogeography**
(BASTA/BEAST)

**Tools of Trade**
Git & GitHub
Linux System
R & Python

Progress Reports
Thesis
Publications

DATA ANALYSIS
hpc

# Data Mining and Clean up

## Sequence retrival:

BOLD-Systems via API; using specific syntax:

"arthropoda Kenya Uganda Tanzania Rwanda Burundi"

## Data clean up:

R, Awk and Bash

| Metadata (80 columns) | | | | |
|---|---|---|---|---|
| 'COI-5P' = 35990  out of  37257 (1 sample) | | | | |
| #nucleotides / unaligned seqs / #ns | | | | |
| Over 700 (592) | Under 500 (705) | 500 - 700 (34693) | 650 - 660 (21886) | Over 500 (35285) |
| 1 smpl | 1 smpl | 3 smpls | 2 smpls | |
| Build.fasta: ProcessID \| order \| seq_len \| seq | | | | |

# MSA

## **Large dataset:**

- Accuracy

- Speed

### **Algorithmn:**

- Progressive (mafft/muscle)

- Progressive & transivity (pasta)

- Regressive (T_coffee)

| Muscle: default | Mafft: G_lins1 | t_coffee: -reg | Pasta: default | sate |
|---|---|---|---|---|
| Speed | Speed | speed | speed | *** |
| Low accuracy | Accurate | Accurate | Accurate | *** |
| (refine, align and merge) | (align, add, addfragments, and merge) | (align, evaluate: core_index, Tcs, | (align and addfragments) | *** |

# MSA evaluation

## Evaluation for Accuracy:

- T_coffee: consistency based scoring.

## Scores:

- Core index (html)

- TCS (html/ascii)

- Used to select the most suitable alignments.

- TCS ascii can be used to apply different weights to columns in phylogenetic analysis

# MSA eval...

# Current progress...

Build my Data set

Setting up a RAxML8 and FastTree pipelines

Improve on alignments:
translation-alignment-threading of DNA

HPC analysis

# ACKNOWLEDGEMENT:
Thank you for enabling a bioinformatics dream