# The H3ABioNet Reproducible WorkflowsProject

**www.h3abionet.org**

## International Data Week

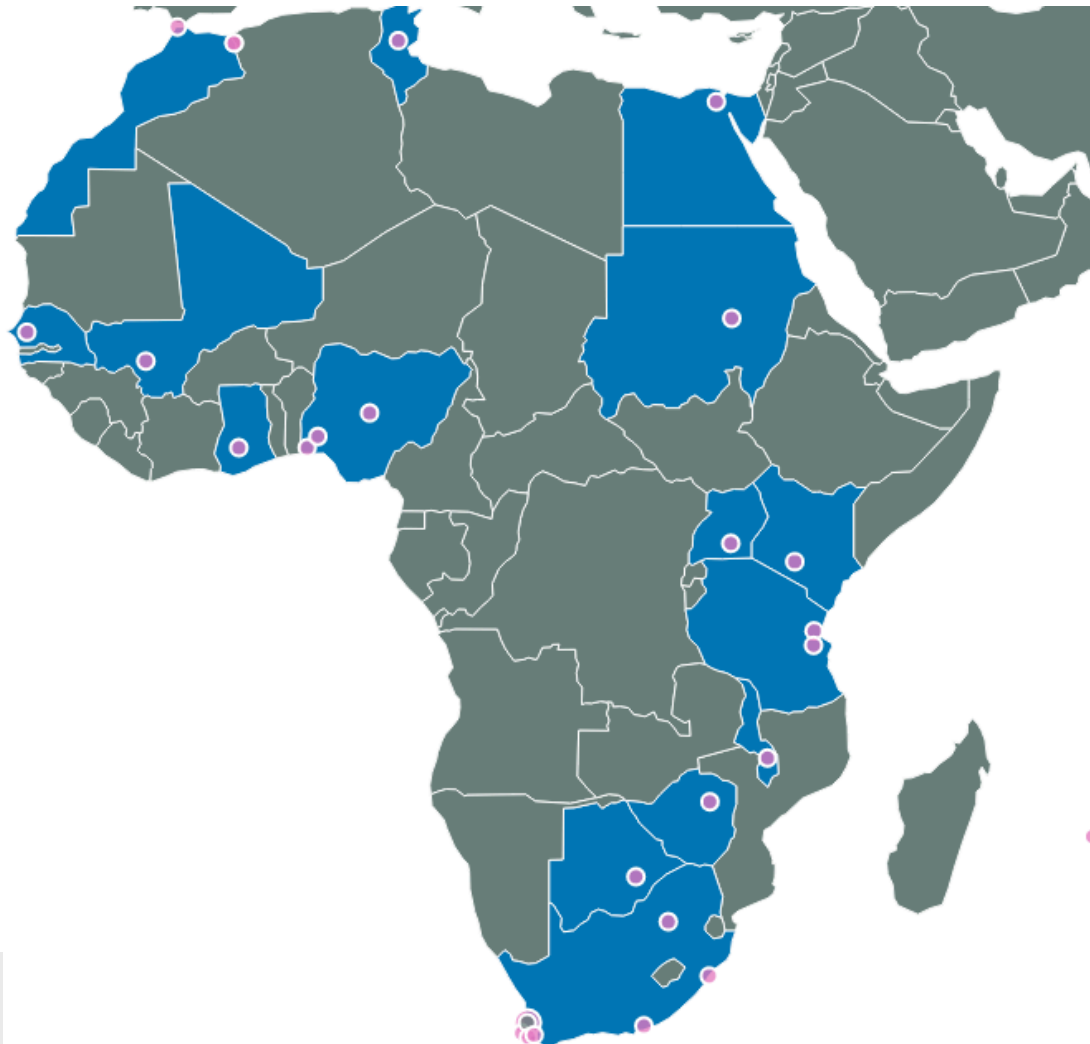## SciDataCon 2018

## Botswana

# H3ABioNet 2.0 Specific Aims

① To implement a Pan African informatics infrastructure

② To develop an H3Africa data coordinating center

③ To provide high quality informatics support to H3Africa

④ To enable and enhance innovative translational research

⑤ To address outreach, development and sustainability

**H3ABioNet**
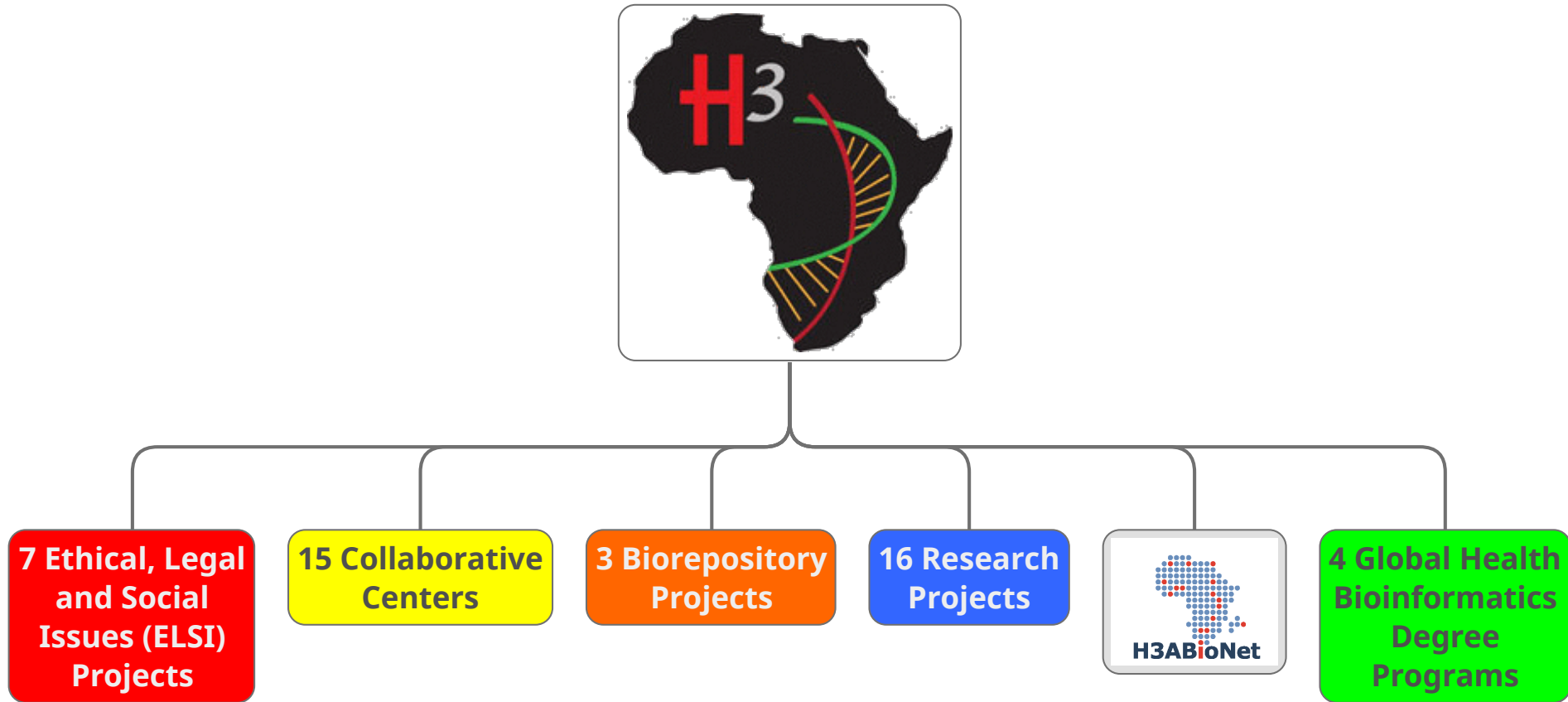Pan African Bioinformatics Network for H3Africa

# H3ABioNet 2.0 Specific Aims

①**To implement a Pan African informatics infrastructure**

② To develop an H3Africa data coordinating center

③**To provide high quality informatics support to H3Africa**

④ To enable and enhance innovative translational research

⑤**To address outreach, development and sustainability**
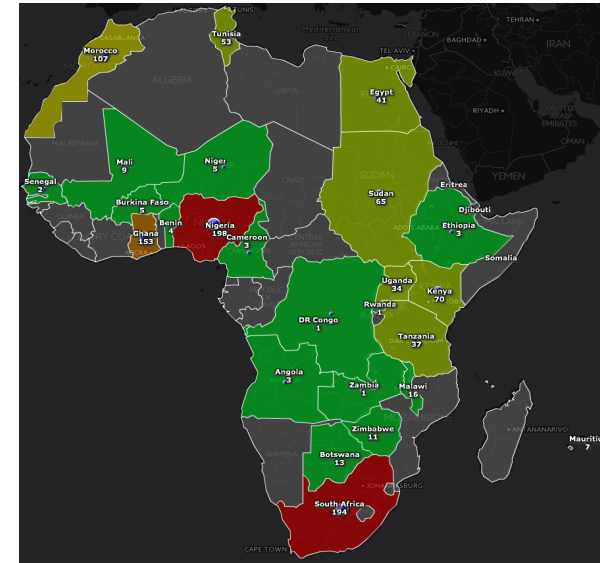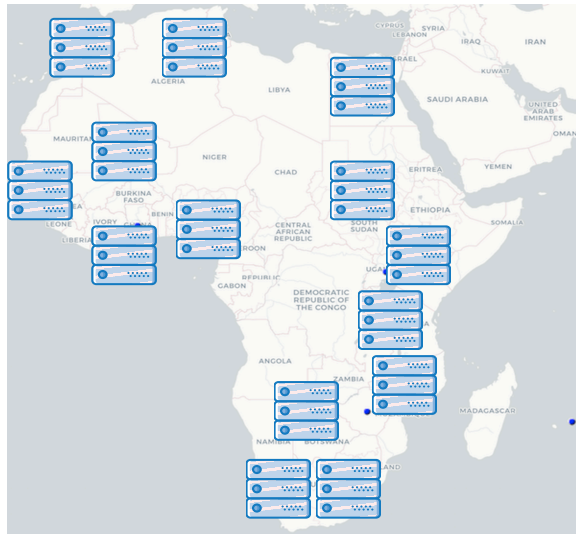
# H3Africa Bioinformatics Network (H3ABioNet)

- Pan African Bioinformatics Network to develop bioinformatics capacity in Africa and support the H3Africa research projects

# H3Africa Consortium



7 Ethical, Legal and Social Issues (ELSI) Projects

15 Collaborative Centers

3 Biorepository Projects

16 Research Projects

H3ABioNet

4 Global Health Bioinformatics Degree Programs

http://h3africa.org/consortium/projects

H3ABioNet
Pan African Bioinformatics Network for H3Africa

# H3ABioNet Workflows project



Compute infrastructure provided
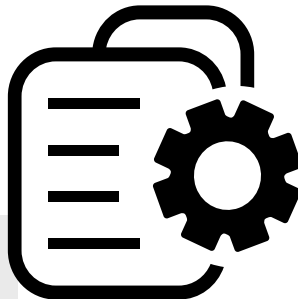


Training provided (sys admin, NGS)



Expertise by Nodes to run analysis to support H3Africa projects in region?
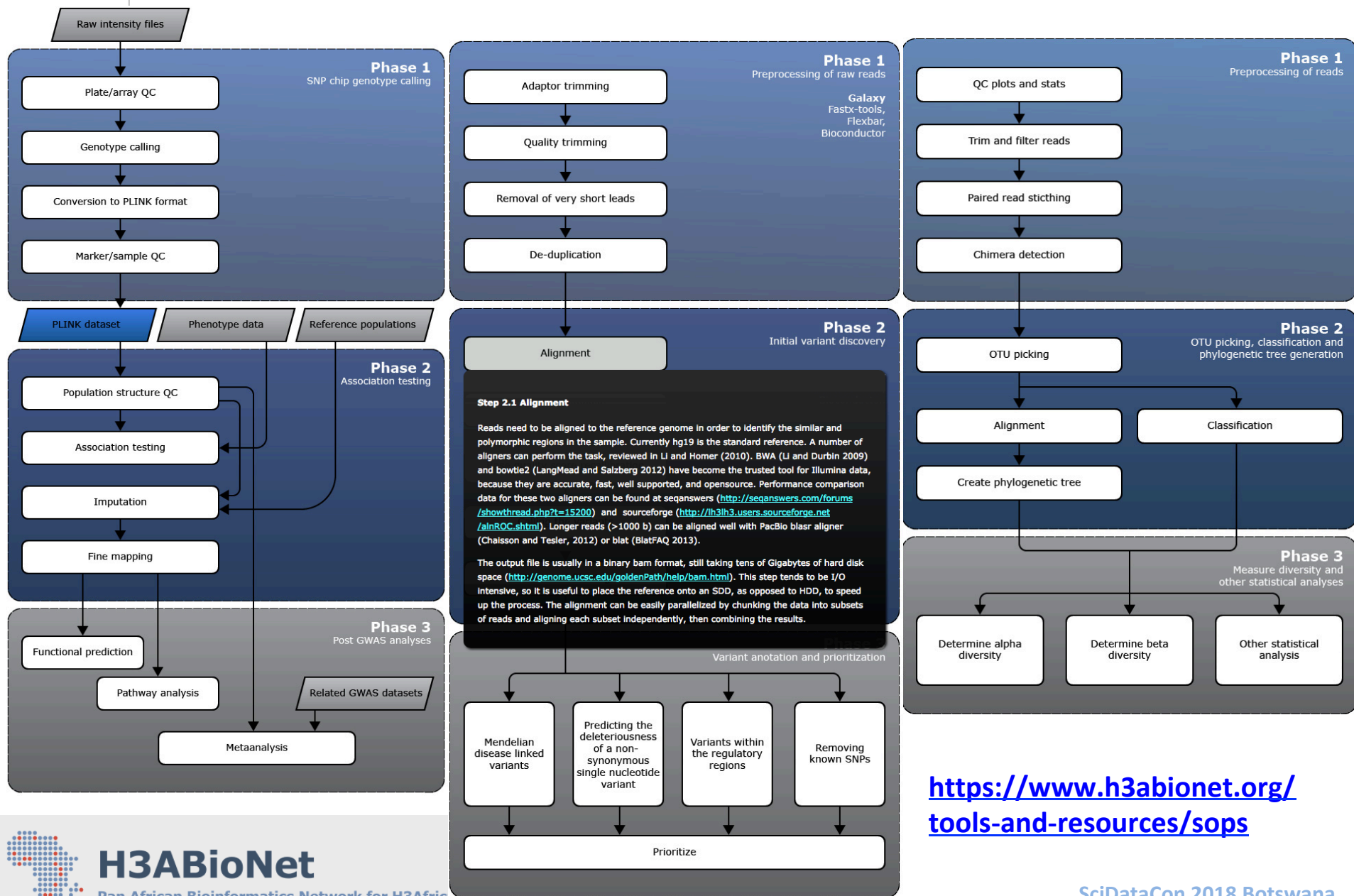
# H3ABioNet Node Accreditation Task force

Create a process for H3ABioNet Nodes to be assessed on the ability to undertake bioinformatics analyses of expected H3Africa genomics data

- Constituted by Victor Jongeneel to:
  - Create Standard Operating Procedures for various "omics" analysis
  - Create practice datasets for Nodes to work with
  - Create assessment datasets
  - Setup process for administrating a Node assessment
  - Constitute external board of reviewers to assess final report submitted

# H3ABioNet SOPs and Accreditation

Raw intensity files

**Phase 1**
SNP chip genotype calling

- Plate/array QC
- Genotype calling
- Conversion to PLINK format
- Marker/sample QC

PLINK dataset | Phenotype data | Reference populations

**Phase 2**
Association testing

- Population structure QC
- Association testing
- Imputation
- Fine mapping

**Phase 3**
Post GWAS analyses

- Functional prediction
- Pathway analysis
- Related GWAS datasets
- Metaanalysis

**Phase 1**
Preprocessing of raw reads

**Galaxy**
Fastx-tools,
Flexbar,
Bioconductor

- Adaptor trimming
- Quality trimming
- Removal of very short leads
- De-duplication

**Phase 2**
Initial variant discovery

- Alignment

**Step 2.1 Alignment**

Reads need to be aligned to the reference genome in order to identify the similar and polymorphic regions in the sample. Currently hg19 is the standard reference. A number of aligners can perform the task, reviewed in Li and Homer (2010). BWA (Li and Durbin 2009) and bowtie2 (LangMead and Salzberg 2012) have become the trusted tool for Illumina data, because they are accurate, fast, well supported, and opensource. Performance comparison data for these two aligners can be found at seqanswers (http://seqanswers.com/forums /showthread.php?t=15200) and sourceforge (http://lh3lh3.users.sourceforge.net /alnROC.shtml). Longer reads (>1000 b) can be aligned well with PacBio blasr aligner (Chaisson and Tesler, 2012) or blat (BlatFAQ 2013).

The output file is usually in a binary bam format, still taking tens of Gigabytes of hard disk space (http://genome.ucsc.edu/goldenPath/help/bam.html). This step tends to be I/O intensive, so it is useful to place the reference onto an SDD, as opposed to HDD, to speed up the process. The alignment can be easily parallelized by chunking the data into subsets of reads and aligning each subset independently, then combining the results.

Variant anotation and prioritization

- Mendelian disease linked variants
- Predicting the deleteriousness of a non-synonymous single nucleotide variant
- Variants within the regulatory regions
- Removing known SNPs

- Prioritize

**Phase 1**
Preprocessing of reads

- QC plots and stats
- Trim and filter reads
- Paired read stitching
- Chimera detection

**Phase 2**
OTU picking, classification and phylogenetic tree generation

- OTU picking
- Alignment
- Classification
- Create phylogenetic tree

**Phase 3**
Measure diversity and other statistical analyses

- Determine alpha diversity
- Determine beta diversity
- Other statistical analysis

https://www.h3abionet.org/
tools-and-resources/sops

## H3ABioNet
**Pan African Bioinformatics Network for H3Afric**

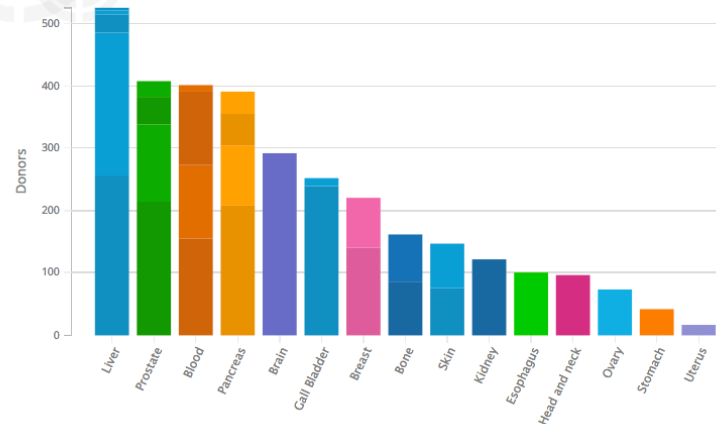# H3ABioNet Workflows project – Cloud



**Cloud Computing for BIG DATA Genomics**

Welcome to the **Cancer Genome Collaboratory**, an academic compute cloud resource that allows researchers to run complex analysis operations across large ICGC cancer genome data sets.

ABOUT OUR SERVICES →

**Collaboratory Data Repository: Donor Distribution by Primary Site**
34 projects and 15 primary sites

**The Collaboratory data consists of:**

| | | | |
|---|---|---|---|
| Collaboratory - Toronto | 3,248 donors | 108,973 files | 672.58 TB |
| PDC - Chicago | 885 donors | 21,989 files | 254.48 TB |
| Total | 4,133 donors | 130,962 files | 927.06 TB |

What We Offer

**Cloud Infrastructure**

The Collaboratory hosts an **OpenStack** cloud with more than **2592 CPU cores and over 7.7 PB of storage**, offering resources such as:

- Compute
- Storage
- Networking

ABOUT OUR RESOURCES AND FEES →

**Researchers Sharing Tools**

The Collaboratory offers multiple ways for researchers to share their tools including VM sharing through the OpenStack Console, and Docker container sharing through the Dockstore based on a GA4GH-compliant tool descriptor.

ABOUT DOCKSTORE →

https://www.cancercollaboratory.org

# H3ABioNet Workflows project - Containers



**https://dockstore.org**

# H3ABioNet Workflows project

- Cloud computing task force constituted to look at potential applications in March 2016

- Mainly learning with some groups having experience

- Brainstormed different ideas – convert pipelines to workflows and containerize for easy deployment

# H3ABioNet Workflows project - Containers

- Containers – means to package software, tools, dependencies



- Docker – containerization technology widely adopted by groups

- Portable from various Linux OS to Linux OS

**Image credits:**
https://www.zdnet.com/article/what-is-docker-and-why-is-it-so-darn-popular
https://blog.docker.com/2017/08/docker-101-introduction-docker-webinar-recap

# H3ABioNet Workflows project - Software



The Elements of Bioinformatics

# H3ABioNet Workflows project - Pipelines



- Bioinformatics analyses → directing files through a series of transformations and programs to a final output i.e. a computational pipeline

- Transformations typically done by third-party executable command line software written for Unix-compatible operating systems

- Manually started when previous transformation step completed e.g. qsub command run

# H3ABioNet Workflows project - Workflows

- A workflow is a description of a process (pipeline) that consists of a series of tasks connected in the form a directed graph

- Tasks can be defined as single units of work e.g. split files

- A workflow comprises of an initial unique task and ends with a unique terminal task

- Completion of a task can initiate one or more tasks

- Enables automation e.g. run a pipeline from start to finish without manual input (pipeline)

Image credit: Brief Bioinform. 2016;18(3):530-536. doi:10.1093/bib/bbw020

# Workflow languages and Workflow systems

- CWL is a workflow language that is explicit

- Provides a specification for describing workflows making them portable and scalable

- Used by Galaxy, International Cancer Genome Consortium, GATK

- Nextflow is a workflow language and system that integrates to with common resource management systems (work on CWL Toil and Cromwell integrations has progressed)

- Supports parallelism and has clear execution blocks

- Used by the Center for Genomic Regulation, Wellcome Trust Sanger Institute, Berkley, Wits' Bioinformatics, UCT CBIO

# H3ABioNet Workflows project - planning

- Workflow languages to use?  **next flow** COMMON WORKFLOW LANGUAGE

- Containerization technology?  docker

- Heterogeneous African compute environments (portability) – HPC 

- What workflows to develop? 

**Image credits: National Human Genome Research Institute (https://www.genome.gov/imagegallery/)**

# H3ABioNet Workflows project - execution

- Hackathon held at University of Pretoria (Prof. Fourie Joubert's lab)



- Four streams devised (Variant calling, 16S, GWAS and Imputation)

- Each stream had a mixture of skills from bioinformaticists, developers, sys admin, knowledge of tools for the pipeline and expertise in CWL (Michael R. Crusoe) and Nextflow (Prof. Scott Hazelhurst)

- Expert in Docker containerization shared between all streams (Dr. Brian O'Connor)

**H3ABioNet**
Pan African Bioinformatics Network for H3Africa

# H3ABioNet Workflows project - execution



METHOD ARTICLE

## Organizing and running bioinformatics hackathons within Africa: The H3ABioNet cloud computing experience [version 1; referees: 3 approved with reservations]

Azza E. Ahmed [1,2*], Phelelani T. Mpangase[3*], Sumir Panji[4], Shakuntala Baichoo [5], Gerrit Botha[4], Faisal M. Fadlelmola [1], Scott Hazelhurst[3,6], Peter Van Heusden[7], C. Victor Jongeneel [8], Fourie Joubert[9], Liudmila Sergeevna Mainzer[8,10], Ayton Meintjes [4], Don Armstrong[8], Michael R. Crusoe [11], Brian D. O'connor[12], Yassine Souilmi [13], Mustafa Alghali[1], Shaun Aron[3], Hocine Bendou [7], Eugene De Beste[7], Mamana Mbiyavanga[4], Oussema Souiai [14], Long Yi[7], Jennie Zermeno[10], ✉ Nicola Mulder [4]

* Equal contributors

+ Author details

## Abstract

The need for portable and reproducible genomics analysis pipelines is growing globally as well as in Africa, especially with the growth of collaborative projects like the Human Health and Heredity in Africa Consortium (H3Africa). The Pan-African H3Africa Bioinformatics Network (H3ABioNet) recognized the need for portable, reproducible pipelines adapted to heterogeneous compute environments, and for the nurturing of technical expertise in workflow languages and containerization technologies. To address this need, in 2016 H3ABioNet arranged its first Cloud Computing and Reproducible Workflows Hackathon, with the purpose of building key genomics analysis pipelines able to run on heterogeneous computing environments and meeting the needs of H3Africa research projects. This paper describes the preparations for this hackathon and reflects upon the lessons learned about its impact on building the technical and scientific expertise of African researchers. The workflows developed were made publicly available in GitHub repositories and deposited as container images on quay.io.

https://doi.org/10.12688/AASOPENRES.12847.1
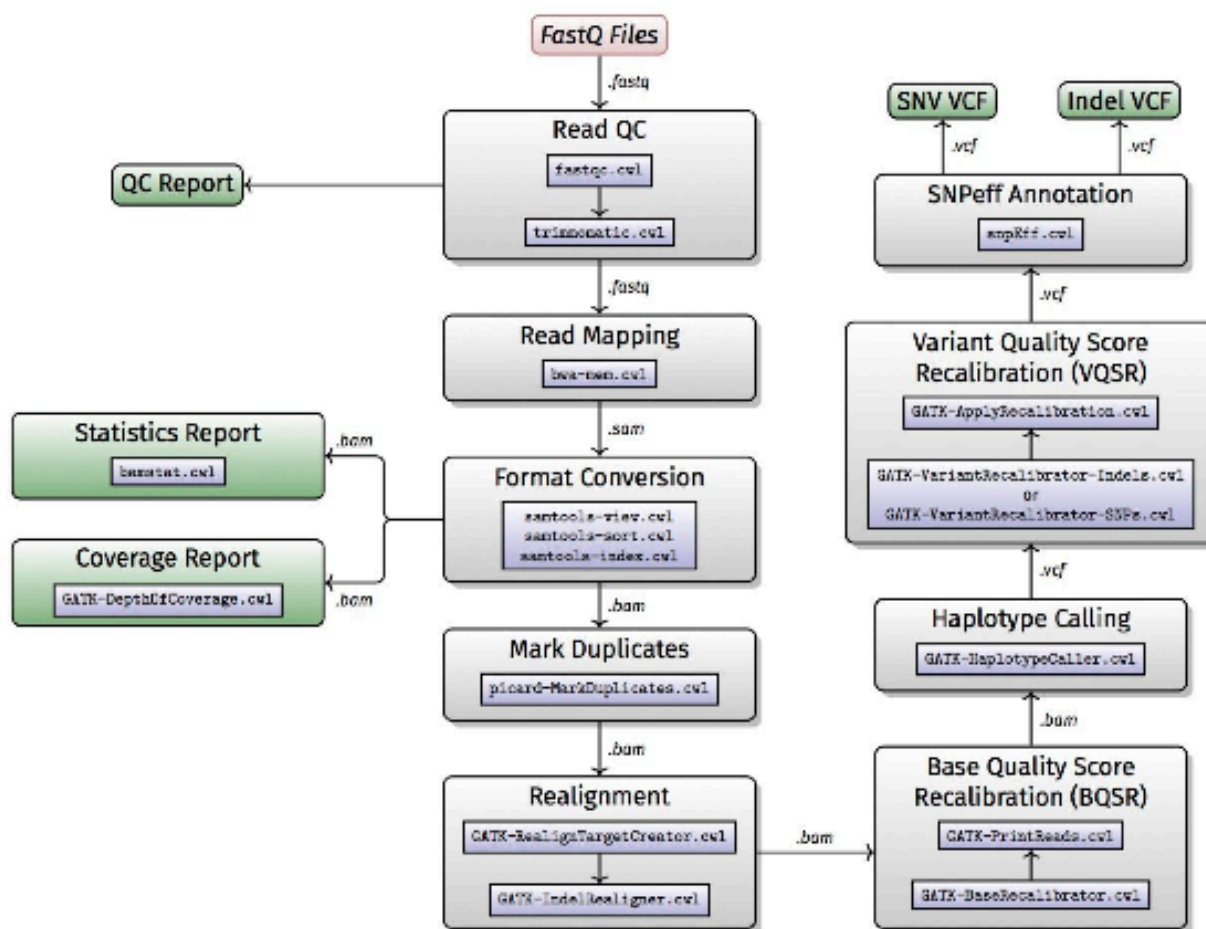
# H3ABioNet Workflows project - outputs



Figure 1 - Workflow A: Whole Genome/Exome NGS Data Analysis

① Sequencing adaptor, barcode and base QC trimming with *Trimmomatic*

② Quality control (QC) of the input fastq files with FastQC

③ Short reads mapping: BWA-MEM is used to perform paired-end Illumina reads

④ Quality control of the aligned reads using BAMstats

⑤ Quality control of the aligned reads using GATK's DepthOfCoverage to check the observed depth of coverage meets expected yield values.

⑥ Indels and single nucleotide variant (SNV) annotation: *SnpEff* extends the VCF file containing the variants with information relevant for downstream analysis. The information includes ranges from the SNP rsID, to clinically relevant variants from ClinVar.

**Availability**:
**https://github.com/h3abionet/h3agatk**

**https://dockstore.org/workflows/h3abionet/h3agatk**

**Image credits: Phelelani Mpangase**

H3ABioNet
Pan African Bioinformatics Network for H3Africa

# H3ABioNet Workflows project - outputs



Figure 2 - Workflow B: 16S rDNA Diversity Analysis

**Image credits: Phelelani Mpangase**

① FastQC for QC reporting of the data quality
② In house scripts: "UPARSE fastq renamer" to rename the FastQ files for compatibility with UPARSE scripts and "UPARSE derep" workaround and mapping back to OTUs
③ Helper scripts: "Combine fastqc reports", "fastq renamer" for UPARSE compatibility and "Fasta splitter" for splitting files
④ usearch for QC and OTU clustering
⑤ QIIME modules for demultiplexing, quality filtering, OTU picking, taxonomic assignment, phylogenetic reconstruction, diversity analyses of the data and the generation of OTU summary
⑥ R and relevant packages for statistical analysis using QIIME results, namely PhyloSeq

**Availability**:
https://github.com/h3abionet/h3abionet16S

https://quay.io/repository/h3abionet_org/h3a16s-qiime

# H3ABioNet Workflows project - outputs



Figure 3 - Workflow C: Genome Wide-association studies

The workflow (Figure 3) consists of 3 modules, which can be swapped in and out depending on the analysis needs:

① Conversion from Illumina TOP/BOTTOM call format to PLINK format.

② The core workflow carries out a set of QC steps, starting with standard PLINK files and resulting in quality controlled PLINK files.

③ Basic association testing and structure analysis.

In addition, we expect many researchers will use the imputation workflow after QC and before association testing.

**Availability**:
http://github.com/h3abionet/h3agwas

https://quay.io/organization/h3abionet_org

**Image credits: Phelelani Mpangase**

# H3ABioNet Workflows project - outputs



Figure 4 - Workflow D: SNPs Imputation: Boxed subgraphs indicate pathways which are executed in parallel (per chromosome and per region within each chromosome) as computational resources permit.

① Identify regions for imputation based on PLINK format an input file, output produced in IMPUTE haplotype format
② Ped and map input files split by chromosome using PLINK, chromosome extents are identified using a combination of awk and grep
③ Genotyped positions on individual chromosomes checked for strand flipping errors, improperly stranded positions are excluded using SHAPEIT
④ Genotyped positions prephased using SHAPEIT in parallel on each chromosome
⑤ IMPUTE2 run in parallel across all 500kB blocks in the entire genome
⑥ Imputed blocks combined into a single compressed haplotype file using custom perl script provided with the workflow.
⑦ File converted back to a PLINK dataset for integration back into the GWAS workflow

**Availability:**
https://github.com/h3abionet/chipimputation/

https://quay.io/repository/h3abionet_org/impute2

**Image credits: Phelelani Mpangase**

# H3ABioNet Workflows project - outputs

# H3ABioNet Workflows project - training

# H3ABioNet Workflows project - FAIR

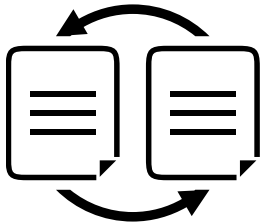Findable: https://github.com/h3abionet

Accessible: https://github.com/h3abionet

Interoperable:
https://quay.io/organization/h3abionet_org

Reusable (Reproducible):

https://quay.io/organization/h3abionet_org

# Acknowledgements

- Prof Nicky Mulder and H3ABioNet members



**H3ABioNet Consortium Members 2018**