

Building Bioinformatics Capacity in Africa

Nicky Mulder
CBIO Group, UCT



H3ABioNet

Pan African Bioinformatics Network for H3Africa



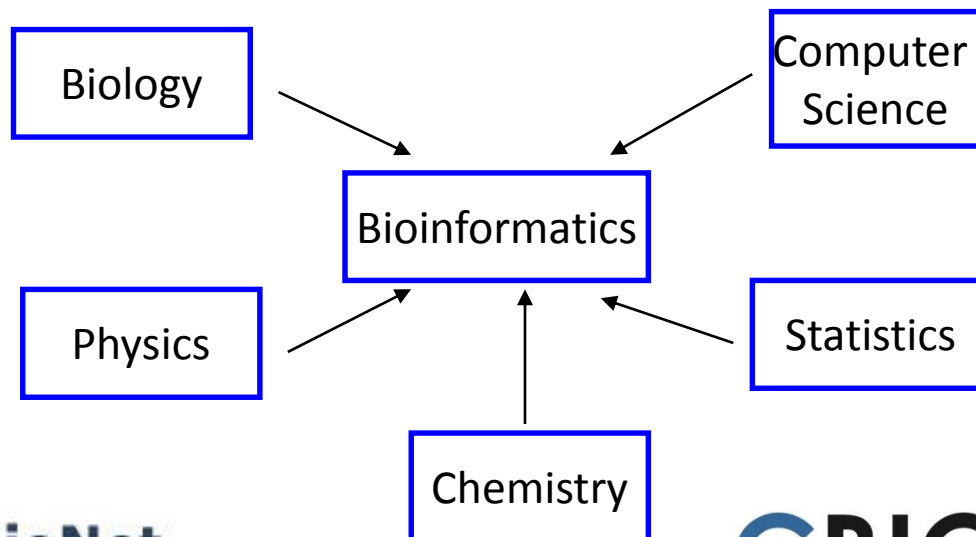
Outline

- What is bioinformatics?
- Why do we need IT infrastructure?
- What e-infrastructure does it require?
- How we are developing this

What is Bioinformatics?

- The analysis of biological information using computers and statistical techniques; the science of developing and utilizing computer databases and algorithms to accelerate and enhance biological research.

www.niehs.nih.gov/nct/glossary.htm



H3ABioNet

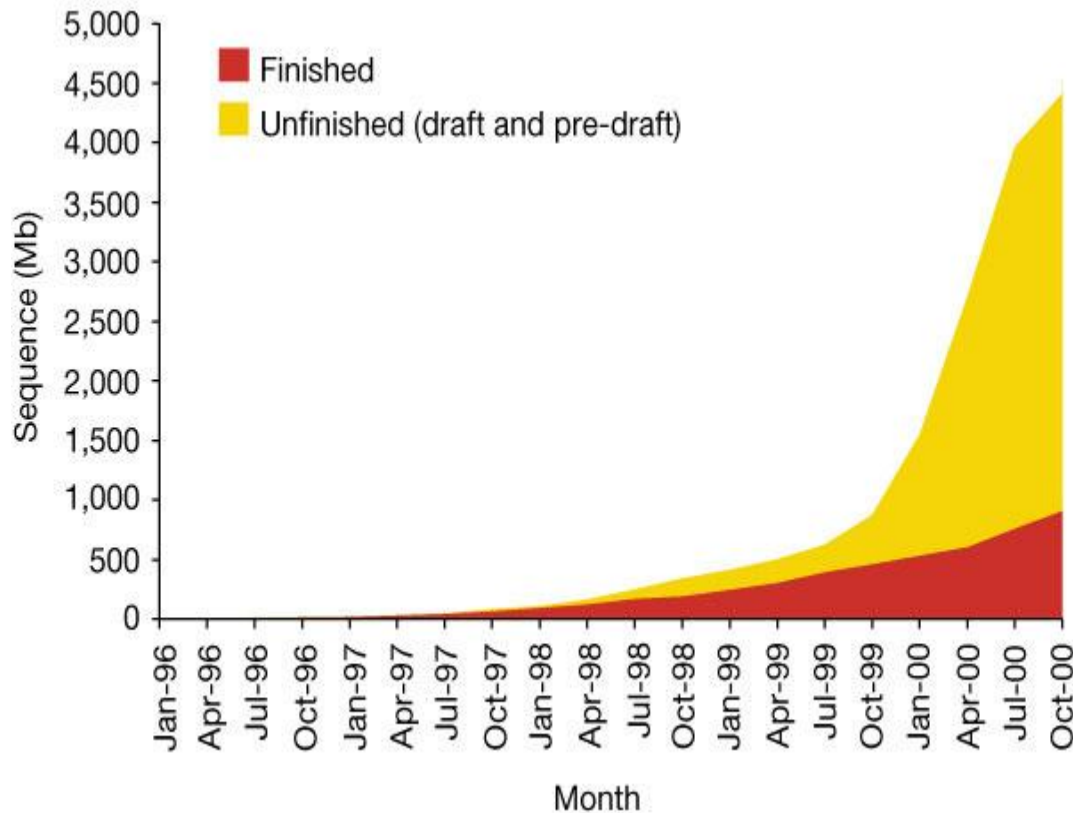
Pan African Bioinformatics Network for H3Africa



Why do we need IT Infrastructure?

- Collection and storage of biological information
- Manipulation of biological information
- Small- and large-scale biological analyses
- New laboratory technologies, move towards big data generation
- Example is genome sequencing

Sequencing the human genome



International Human Genome Sequencing Consortium 2001. Nature 409, 860 – 921.

First human genome took ~5 years and cost ~\$3 billion

Now, can sequence in a few weeks for ~\$5,000

BUT: doesn't consider cost and time for data analysis!

human genome, sequenced at 30x coverage = ~ 1 billion raw reads of about 100bp = ~250Gb of raw data, when processed = $\times 10-15$



H3ABioNet

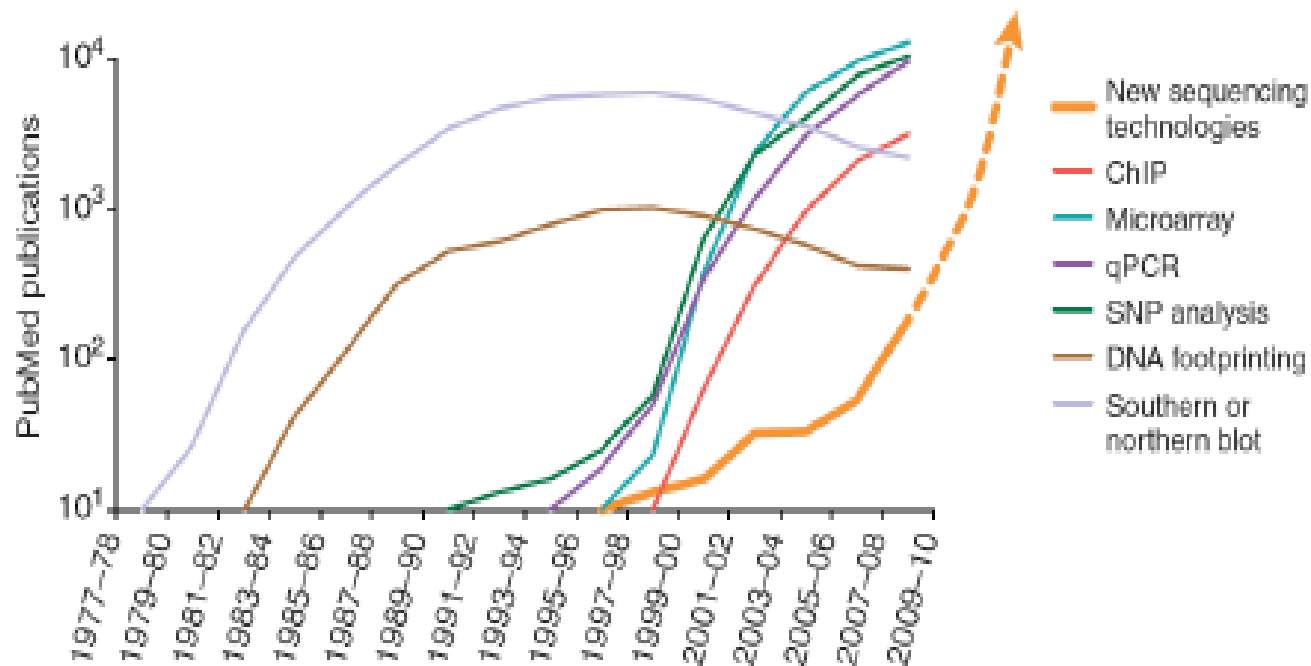
Pan African Bioinformatics Network for H3Africa



Computational Biology @ UCT



Technologies used based on publications



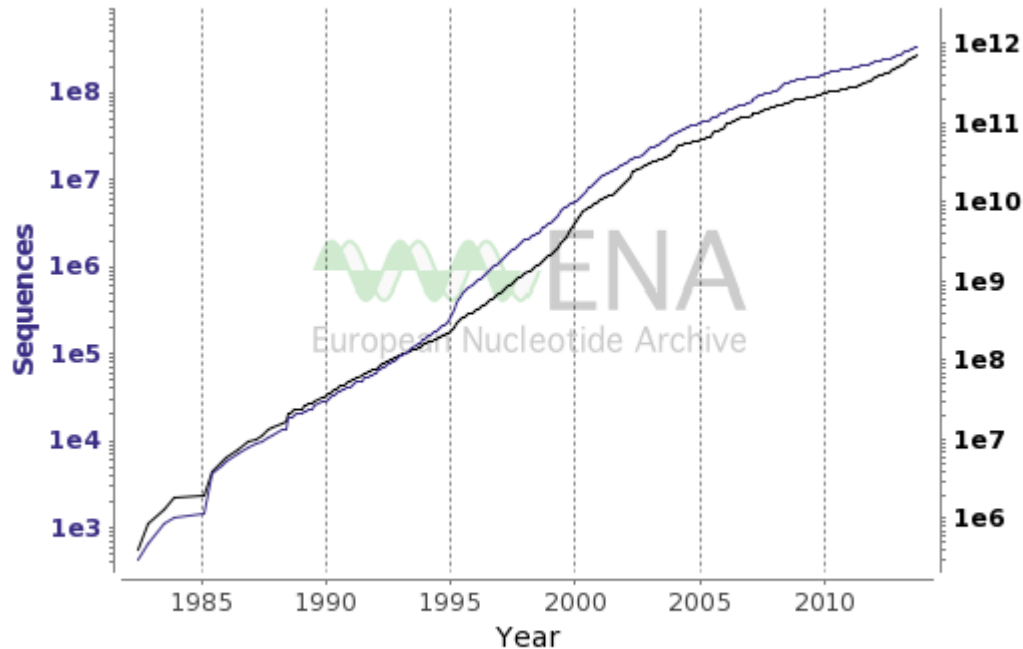
H3ABioNet

Pan African Bioinformatics Network for H3Africa

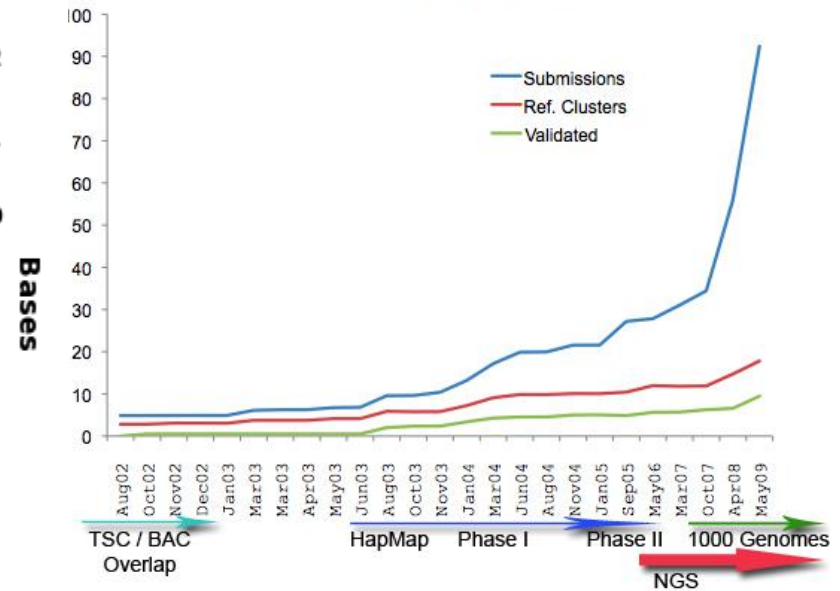


Increase in biological data

EMBL-Bank Growth
07-Oct-2013



Growth of dbSNP, 2002-2009



<http://www.ebi.ac.uk/ena/about/statistics> -Last release of ENA >330 mill seqs, 1.6TB

Biological data growth has exceeded Moore's law



H3ABioNet

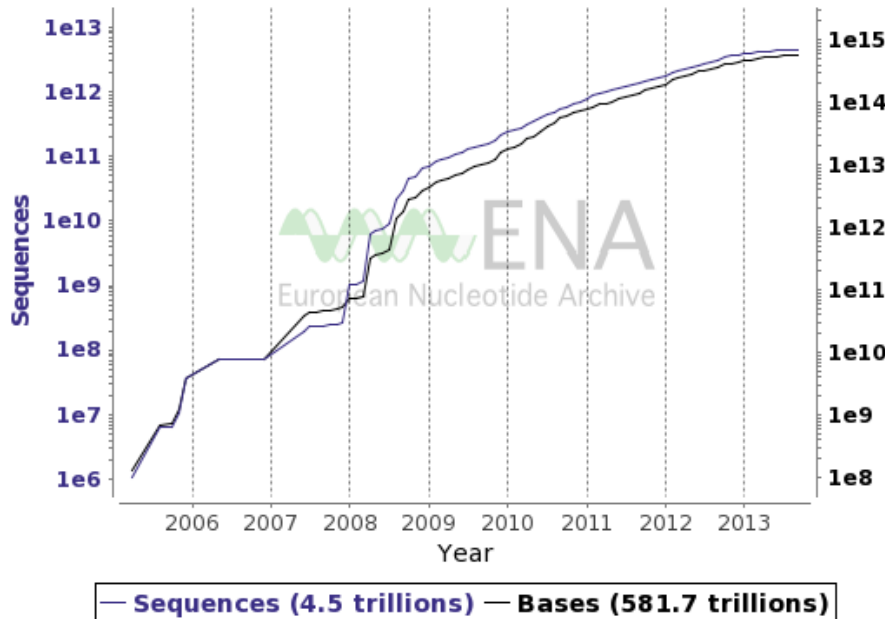
Pan African Bioinformatics Network for H3Africa



New types of data

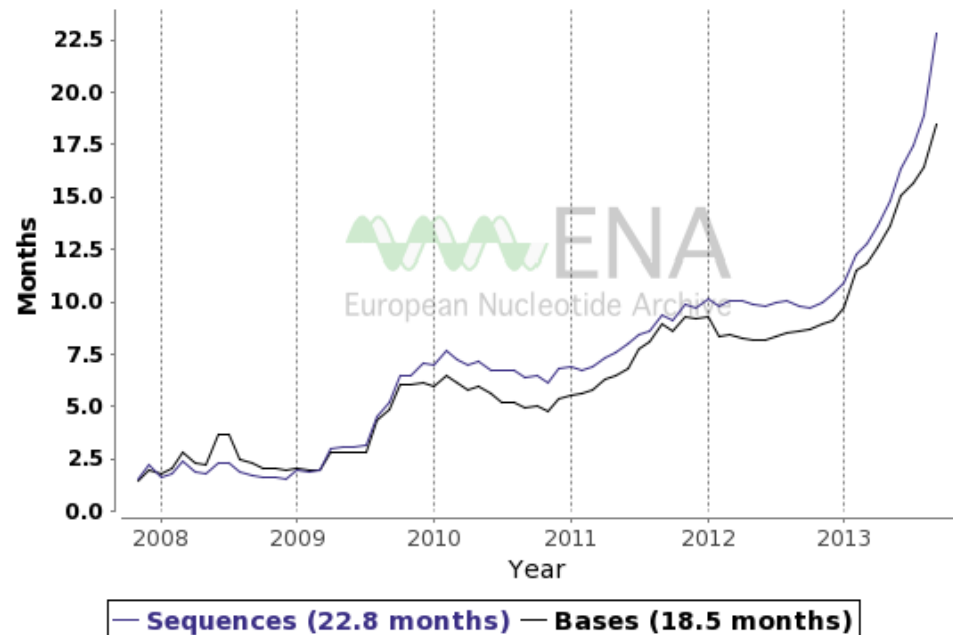
Sequence Read Archive (SRA) Growth

07-Oct-2013



Sequence Read Archive (SRA) Doubling Time

07-Oct-2013



<http://www.ebi.ac.uk/ena/about/statistics>

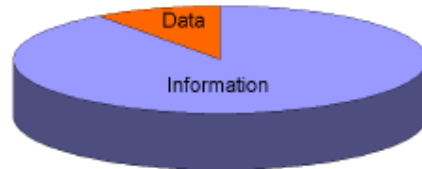


H3ABioNet

Pan African Bioinformatics Network for H3Africa



Data versus information



Genes in the DNA...

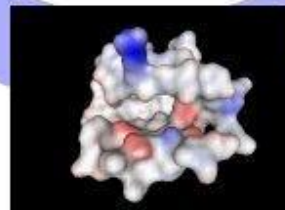


...code for proteins...

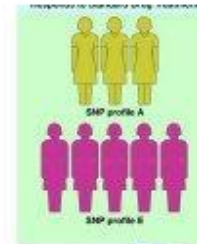
```
>protein kinase  
acctgttgatggcgacagggactgt  
atgctgatctatgctgatgcatgcatg  
ctgactactgatgtggggctattga  
cttgatgctatc....
```

...whose structure accounts for function...

From genotype to phenotype.



...produces the final phenotype



...plus the environment...



H3AB
Pan African



What infrastructure is required?

- Data
 - Storage and management –flatfile archive versus mineable database
 - Processing –single nodes, HPC, Cloud
 - Interpretation –analysis and visualization tools
- Learning how to work with it
 - IT/technical skills –developing tools and data management
 - Bioinformatics skills –using tools and interpreting results



H3ABioNet

Pan African Bioinformatics Network for H3Africa



Data requirements for genomes

- Raw files for NG sequencing are large and increase 5-10 fold after analysis
- Data needs to be stored and processed
- Raw data:
 - Storage requirements are large –some people are working on compression formats for NGS data
 - Data needs to be backed up
 - Also need access to public data
- Data needs to be processed, and made available in a user-friendly format
 - Need compute hardware for QC, alignment, imputation, and tool/visualization development

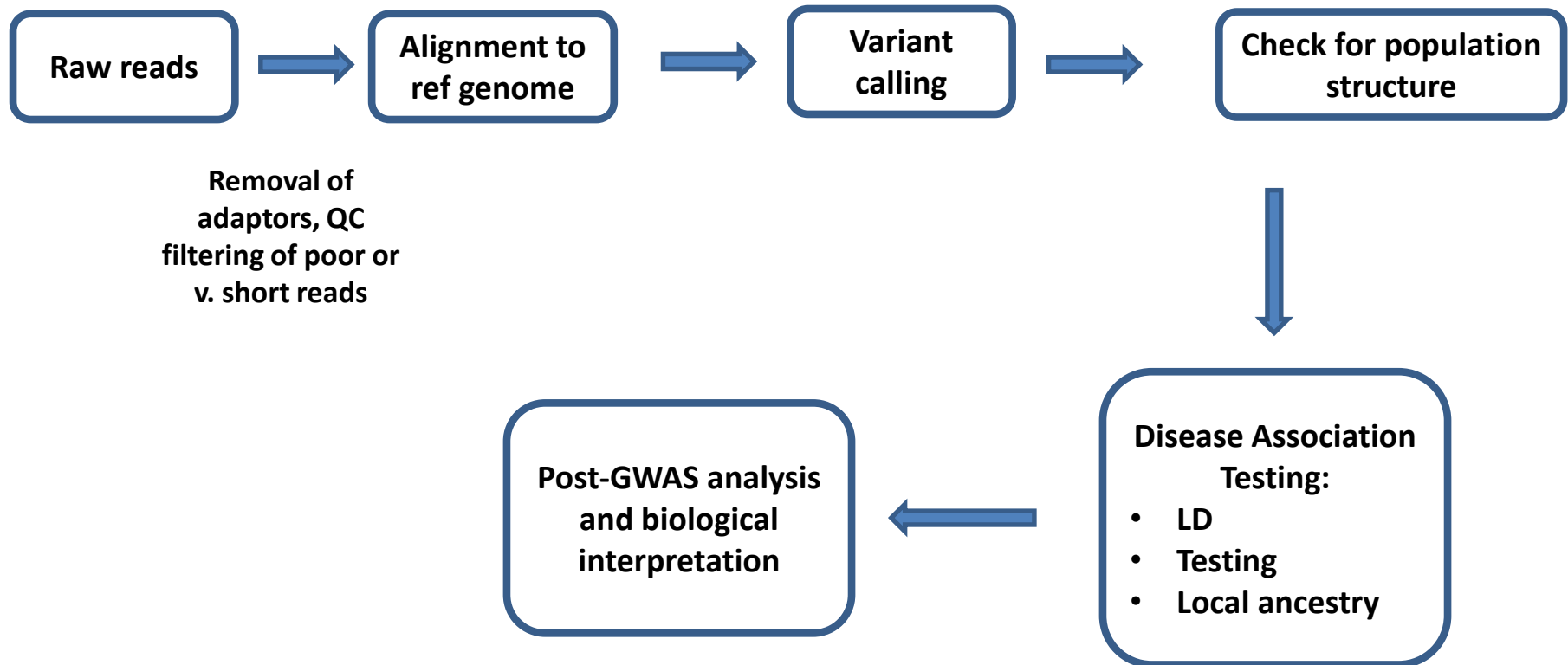


H3ABioNet

Pan African Bioinformatics Network for H3Africa



Workflow for genotyping by NGS

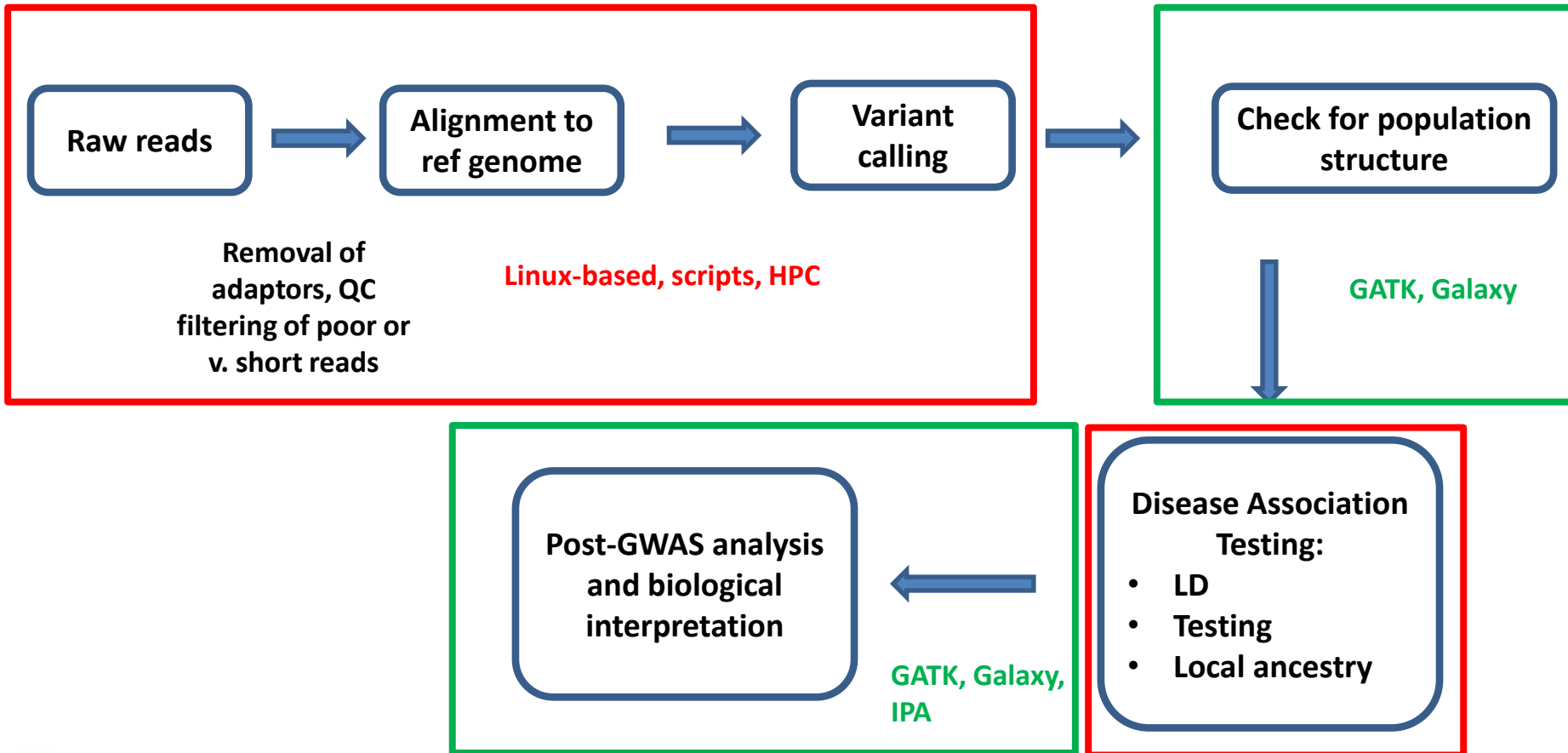


H3ABioNet

Pan African Bioinformatics Network for H3Africa



Workflow for genotyping by NGS -IT

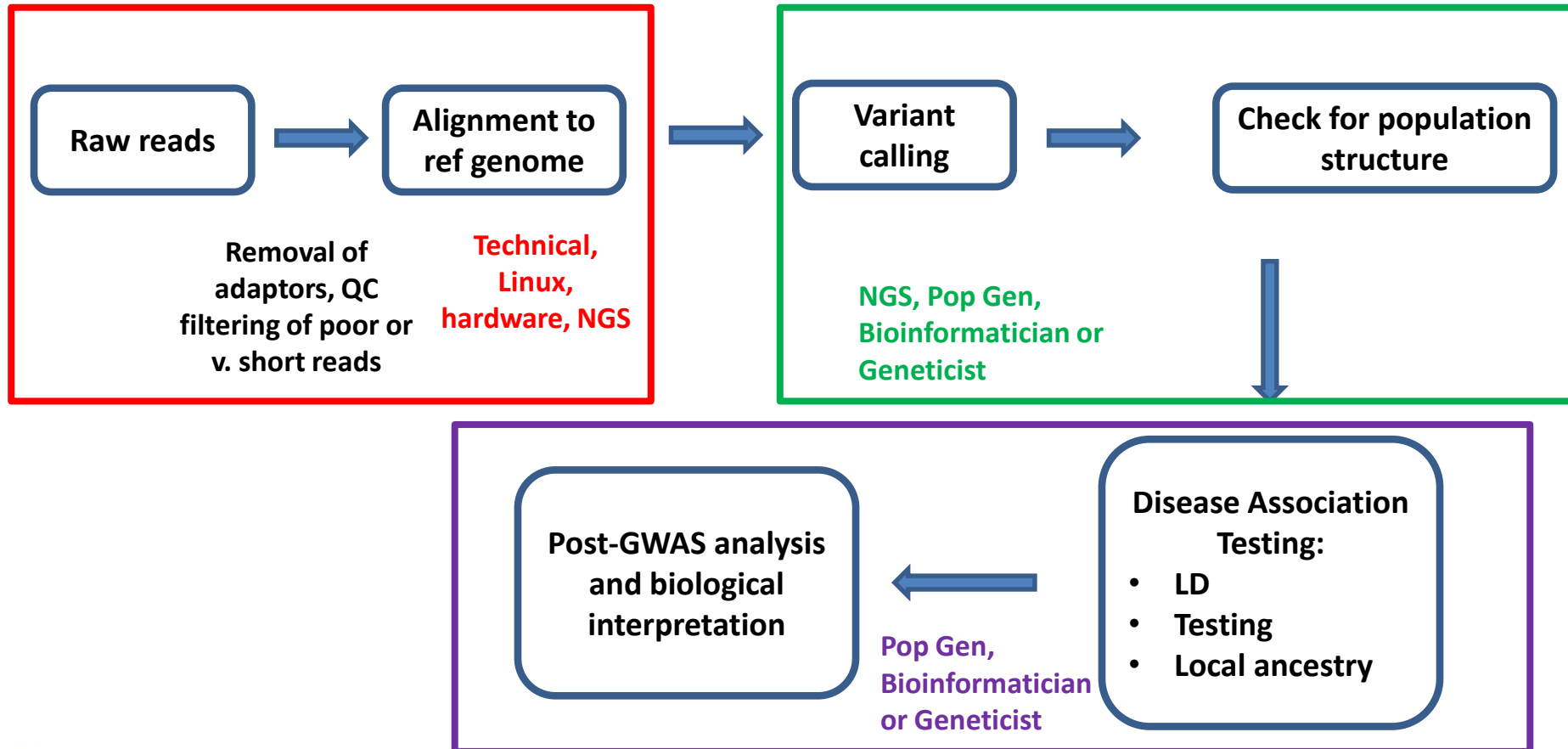


H3ABioNet

Pan African Bioinformatics Network for H3Africa



Workflow for genotyping by NGS -HCD



H3ABioNet

Pan African Bioinformatics Network for H3Africa



Computational Biology @ UCT



Human Capacity Development

Training of 2 groups of scientists:

- Bioinformaticians
 - Programming, data management
 - Data processing
 - Biostatistics
 - Specialised skills (NGS, GWAS, PopGen)
- Researchers –data analysers
 - Basic file manipulation (Galaxy)
 - Biostatistics
 - Specialised skills (NGS, GWAS, PopGen)

How we are developing infrastructure

- Work on national & Africa-wide initiatives
- Local:
 - Support researchers at UCT with data analysis and training for this
 - Working with DST on a national Bioinformatics support platform
- Africa-wide –H3ABioNet



H3ABioNet

Pan African Bioinformatics Network for H3Africa



H3ABioNet Project Goal

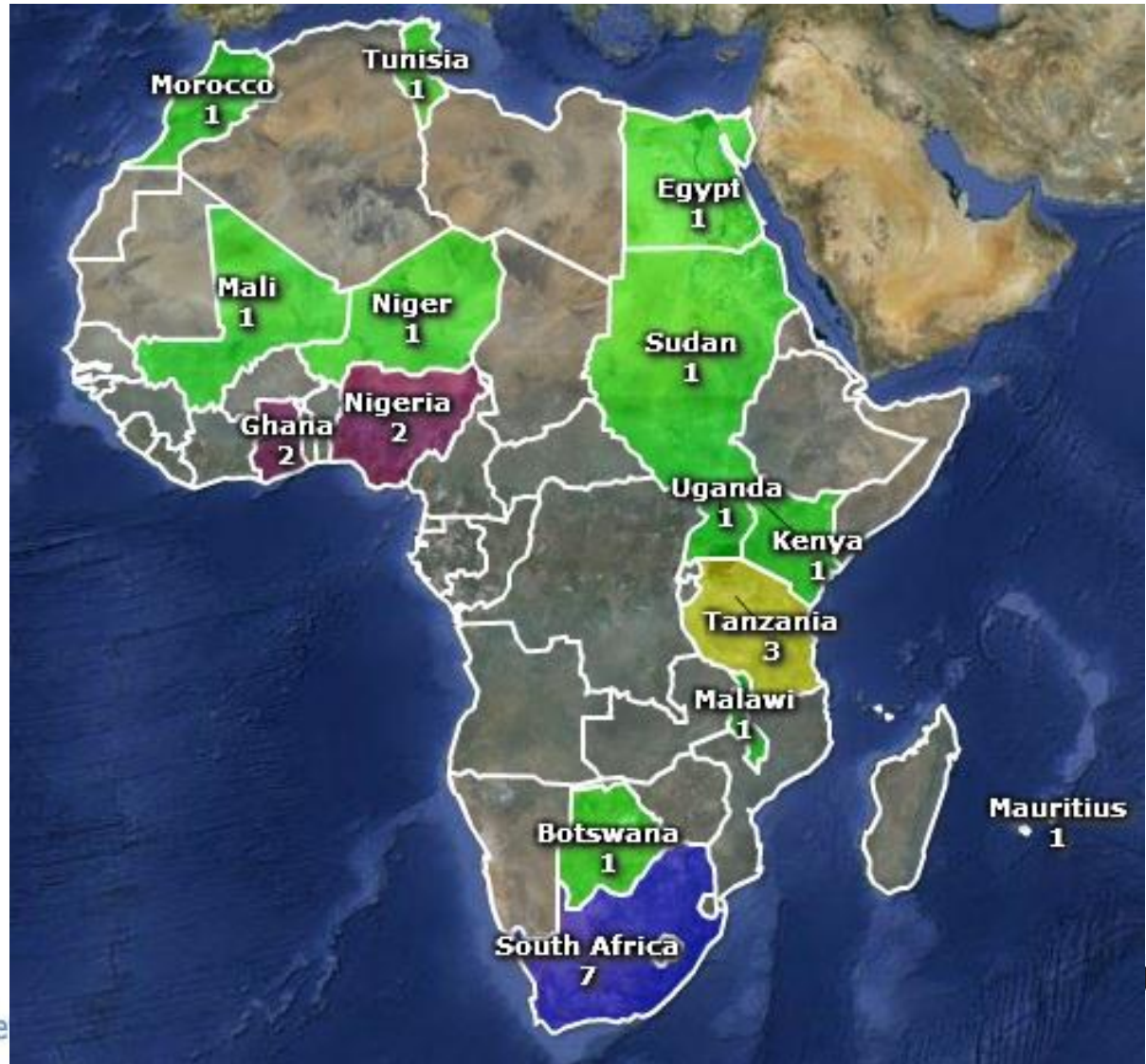
- To build H3ABioNet, a sustainable African Bioinformatics Network, to provide bioinformatics infrastructure and support for the H3Africa consortium.
- H3Africa (Human Heredity and Health in Africa) is an initiative to develop genomics research in Africa on diseases of importance in Africa



Partner institutions

Administrative hub at UCT

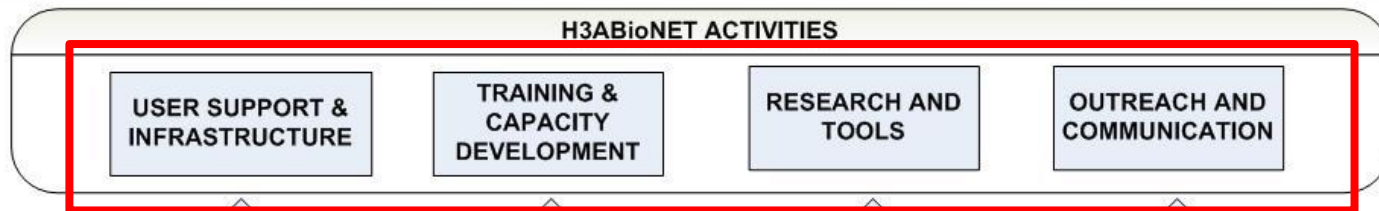
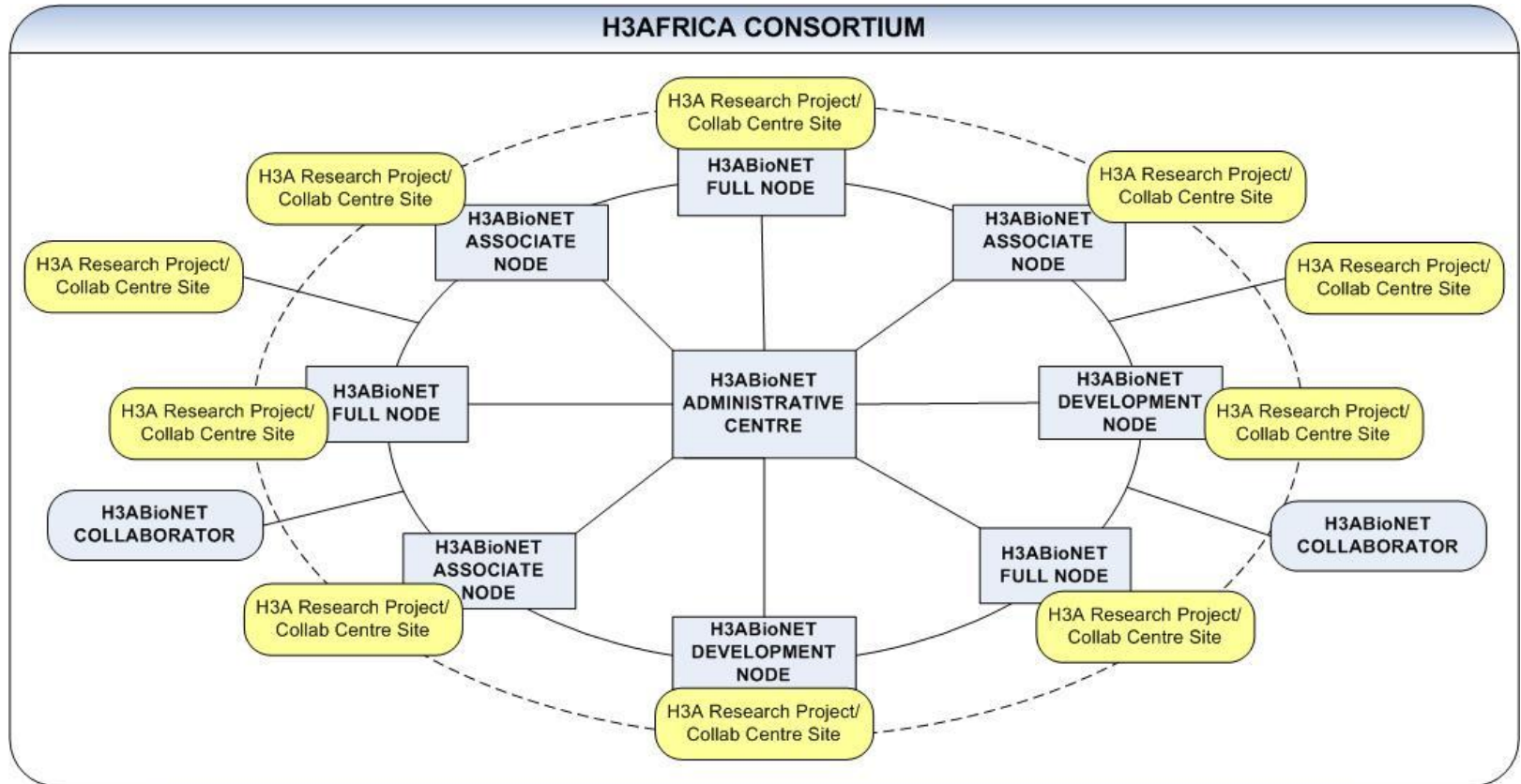
34 partner institutions, 32
in 15 African countries, 2 in
USA



H3ABioNet

Pan African Bioinformatics Network

Network structure



OTHER INTERACTIONS

CHPC

EBI training network

EMBNet

ABioNET



Pan African Bioinformatics Network for H3Africa

CDIO
Computational Biology @ UCT



Infrastructure development & support

- Node server purchases
- Investigating access to HPC, Cloud
- Internet connectivity measurement
- Communication structure established: website, surveys and mailing list setup
- Set up help desk

Hardware Infrastructure development

- Process followed:
 - 34 nodes were surveyed via email to determine existing infrastructure and expected needs
 - 18 nodes requested funding to purchase infrastructure
 - One on one skype follow-up calls with nodes who requested assistance
 - 12 nodes requested assistance
 - 10 bought the recommended hardware
 - 1 bought HP based on recommended hardware
 - 1 bought own configuration

Hardware Infrastructure development

- 3 server builds developed based on survey feedback
 - Option A – Recommended for nodes with an existing infrastructure (single server with max cores and high RAM)
 - Option B – Recommended for nodes needing two physical servers (smaller HPC server and one database server)
 - Option C – Recommended for nodes with no infrastructure (Smaller server but included a rack, switch and UPS)
- Hardware Vendor
 - Approached Dell, IBM and HP
 - Dell was the preferred provider as they have a presence across Africa
 - Provided a 5 year next business day warranty / support to nodes
 - Offered best value for money
 - Negotiated a discount of up to 45% per hardware device

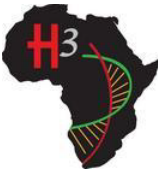
Iperf tests



H3ABioNet

Pan African Bioinformatics Network for H3Africa

CBIO 
Computational Biology @ UCT



H3ABioNet help desk



H3ABioNet

Pan African Bioinformatics Network for H3Africa

[Home](#)[About](#)[Consortium Members](#)[Training and Education](#)[User Support](#)[Members](#)

[Contacts](#)
[Events](#)
[Links](#)
[iAnn](#)

Help desk - dashboard

Helpdesk

[+ Submit new issue](#)[View submitted issues](#) [Edit your contact information](#)

Latest News

- [Second Meeting of the H3Africa Consortium, Accra Ghana](#)

<http://www.h3abionet.org/helpdesk>



National Human
Genome Research
Institute

Computational Biology @ CC



H3ABioNet help desk

Help desk - dashboard

Helpdesk

New Ticket



Contact Information

User Name: *
E-Mail: *
Department: *
Location:
Phone:

Classification

Category: *
Status:
Priority:
Assigned To:
Time Spent:

TicketInformation

Title:

Description

B I U

Notes

Enter Additional Notes

B I U

Solution



[Home](#) [About](#) [Consortium Men](#)

Help desk - dashboard

Helpdesk

[+ Submit new issue](#)

[View submitted issues](#)

[View Issue #](#)

[Edit your contact inform](#)

[Contacts](#)
[Events](#)
[Links](#)
[iAnn](#)

Latest News

- [Second Meeting of the H3Africa Consortium, Accra Ghana](#)

[Search](#)

[Helpdesk user guide](#)

Events

[Second Meeting of the H3Africa Consortium, Accra Ghana](#)

[Helpdesk](#)

[iAnn](#)

H3ABioNet help desk



Help desk - dashboard

Helpdesk

New Ticket



Contact Information

User Name: *
E-Mail: *
Department: *
Location:
Phone:

Classification

Category: *
Status:
Priority:
Assigned To:
Time Spent:

TicketInformation

Title:
Description:

B I U

Notes

Enter Additional Notes

B I U

Solution

[Helpdesk user guide](#)

Events

[Second Meeting of the H3Africa Consortium, Accra Ghana](#)

[Helpdesk](#)

[iAnn](#)

[Home](#) [About](#) [Consortium Men](#)

Help desk - dashboard

Helpdesk

[+ Submit new issue](#)

[View submitted issues](#)

[View Issue #](#)

[Edit your contact inform](#)

[Contacts](#)
[Events](#)
[Links](#)
[iAnn](#)

Latest News

- [Second Meeting of the H3Africa Consortium, Accra Ghana](#)

Categories:

- General Project Administration
- Technical / System Administration
- Website / Mailing List
- Analysis - Genotyping arrays
- Data Management (storage, etc)
- Software Development / Programming
- Analysis - NGS data
- Analysis - Other
- Biostatistics
- Other
- NetCapDB
- Software license request

Event registration

- Field validation
- Email notifications
- Admin backend
- Excel export
- 10 events, 461 registrations

Django administration Welcome, ayton. Change password / Log out

Home > H3aghb > Applicants

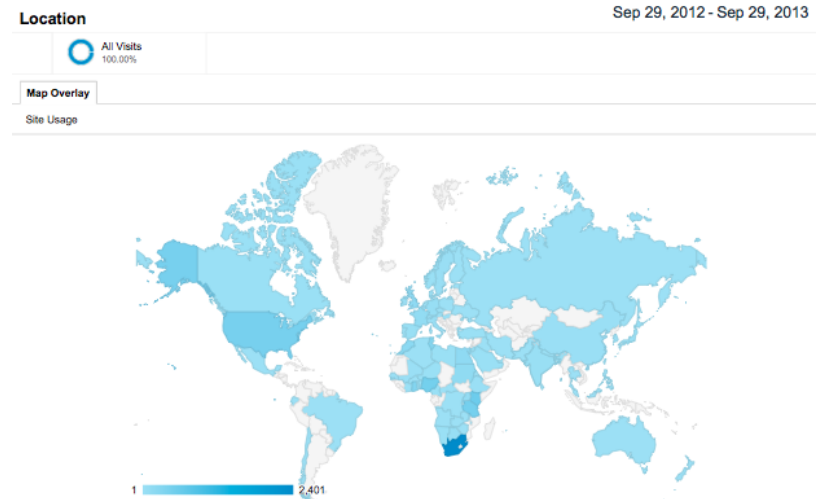
Select applicant to change Add applicant. +

Action: [-----] Go 0 of 100 selected

	First name	Last name	Research group or funding agency to which you belong	VISA required
<input type="checkbox"/> Created	ANANYO	CHOU DHURY	H3A project (Pis-Michele Ramsay, Osman Sankoh)	NO
<input type="checkbox"/> Sept. 26, 2013, 4:39 p.m.	Matti	Kimberg	other (see below)	NO
<input type="checkbox"/> Sept. 26, 2013, 10:24 a.m.	John	Walters	other (see below)	NO
<input type="checkbox"/> Sept. 26, 2013, 9:26 a.m.	Izak	Storm	other (see below)	NO
<input type="checkbox"/> Sept. 25, 2013, 8:28 a.m.	Anne	Arens	other (see below)	NO
<input type="checkbox"/> Sept. 23, 2013, 2:18 p.m.	Reinhard	Eckloff	other (see below)	NO
<input type="checkbox"/> Sept. 23, 2013, 2:07 p.m.	Paul	Mola	other (see below)	NO
<input type="checkbox"/> Sept. 23, 2013, 9:16 a.m.	Dissou	Affolabi	H3A project (P.Dissou Affolabi)	NO
<input type="checkbox"/> Sept. 19, 2013, 11:19 p.m.	Ruth	Chadwick	Funding agency: Wellcome Trust	NO
<input type="checkbox"/> Sept. 19, 2013, 12:01 p.m.	Alhass	Amosh	H3A contact (Dr Alhass Amosh)	YES

Website monitoring

- Google analytics
 - traffic flow
 - traffic sources
- Uptime monitored with pingdom
- Amazon failover server (manual switch)



Country / Territory	Visits	Pages / Visit	Avg. Visit Duration	% New Visits	Bounce Rate
	7,044 % of Total: 100.00% (7,044)	3.67 Site Avg: 3.67 (0.00%)	00:04:10 Site Avg: 00:04:10 (0.00%)	34.64% Site Avg: 29.40% (17.82%)	52.81% Site Avg: 52.81% (0.00%)
1. South Africa	2,401	6.08	00:07:08	18.70%	35.61%
2. Kenya	572	2.13	00:02:16	46.50%	61.19%
3. Nigeria	477	2.26	00:03:26	43.19%	57.65%
4. Tanzania	463	1.91	00:02:23	16.20%	73.87%
5. United States	458	3.16	00:02:21	59.83%	49.78%
6. (not set)	428	2.72	00:03:00	25.47%	70.09%
7. Ghana	311	2.29	00:03:08	26.37%	67.85%
8. Tunisia	259	2.31	00:02:33	47.88%	61.78%
9. Morocco	210	3.58	00:03:21	55.24%	50.48%
10. Egypt	165	2.42	00:02:37	49.70%	50.91%

Rows 1 - 10 of 85

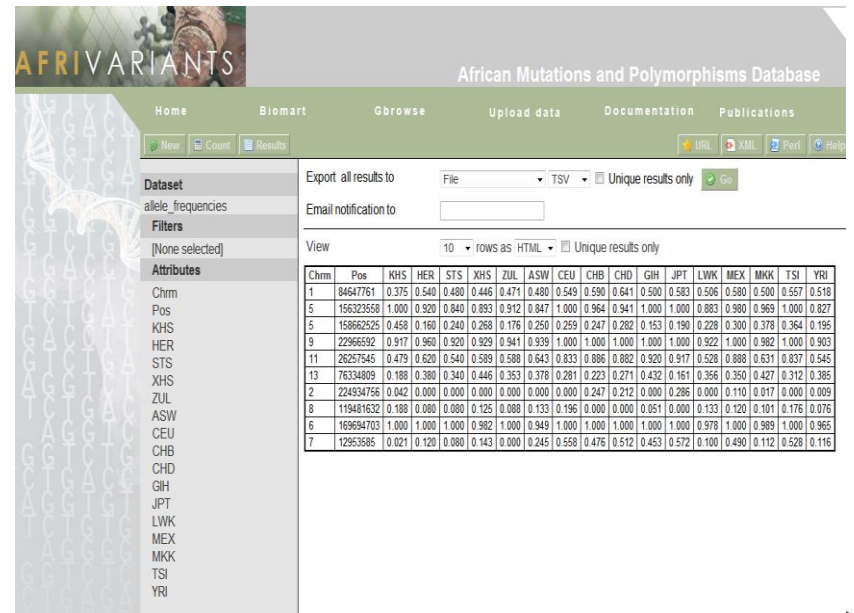
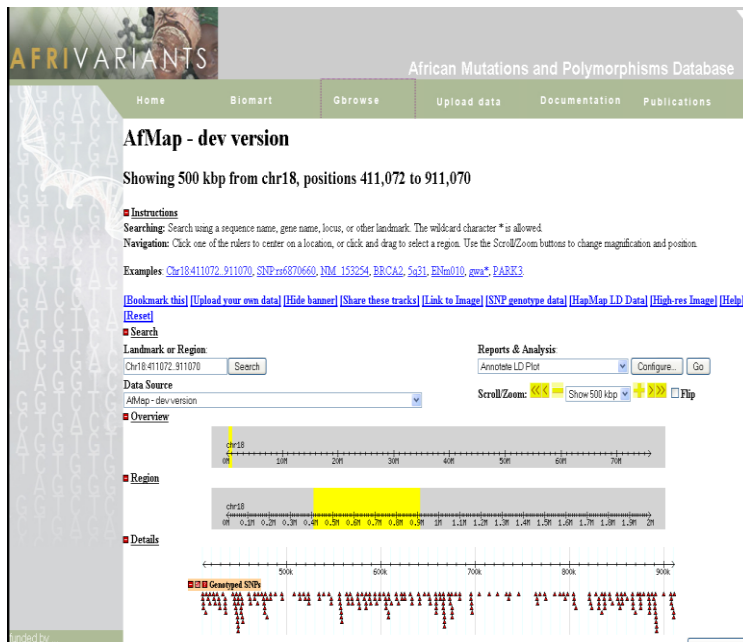


Research and tool development

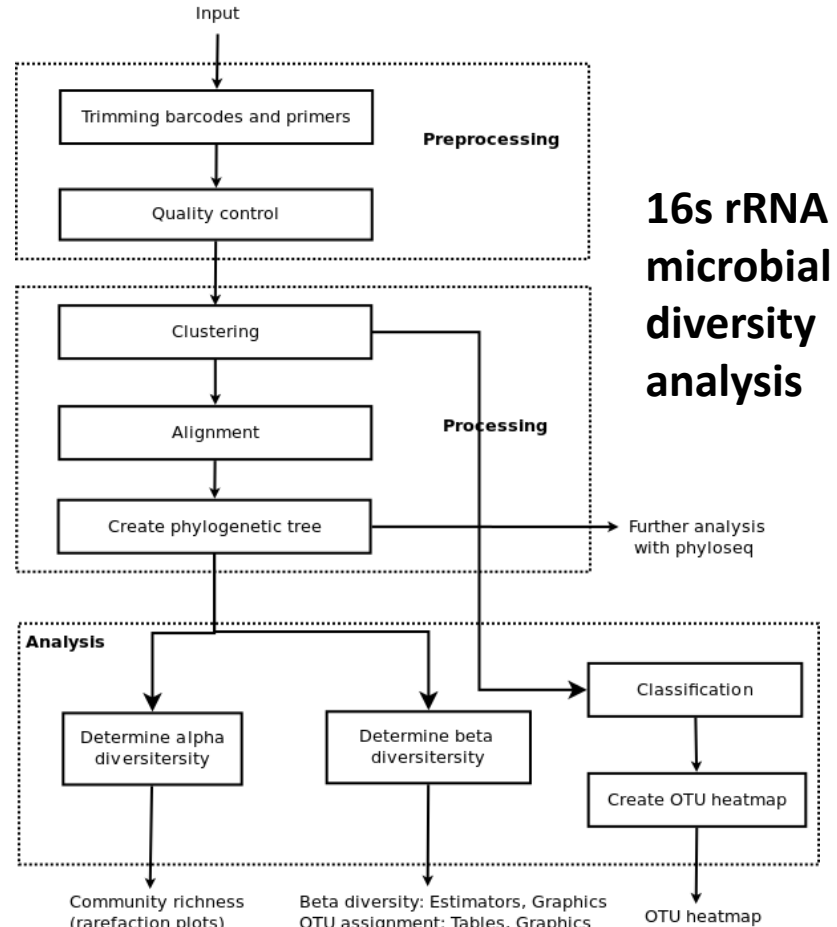
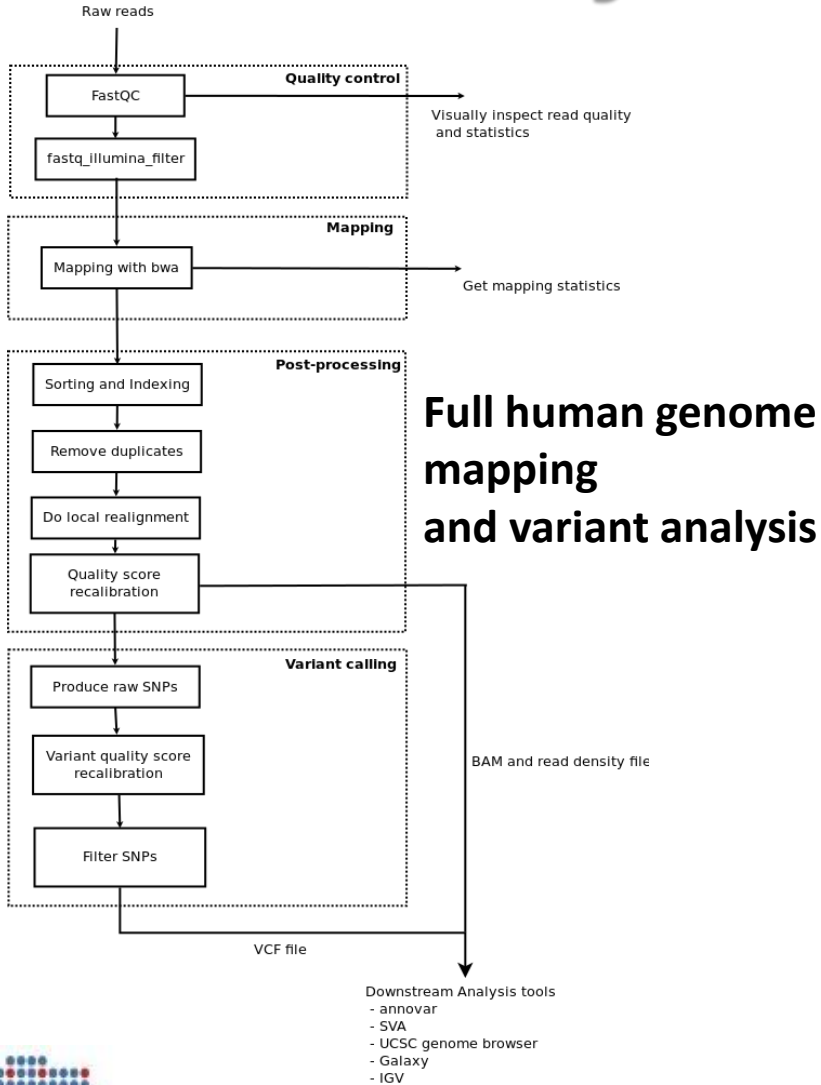
- Data management and storage
 - Data integration platforms, e.g. BioMart
- Data analysis tools
 - Developing pipelines in Galaxy
 - Further developing eBioKits
 - Visualization tools

BioMart database

- Converting relevant data from flatfiles into mineable format
- Includes user-friendly GUI



Analysis workflows

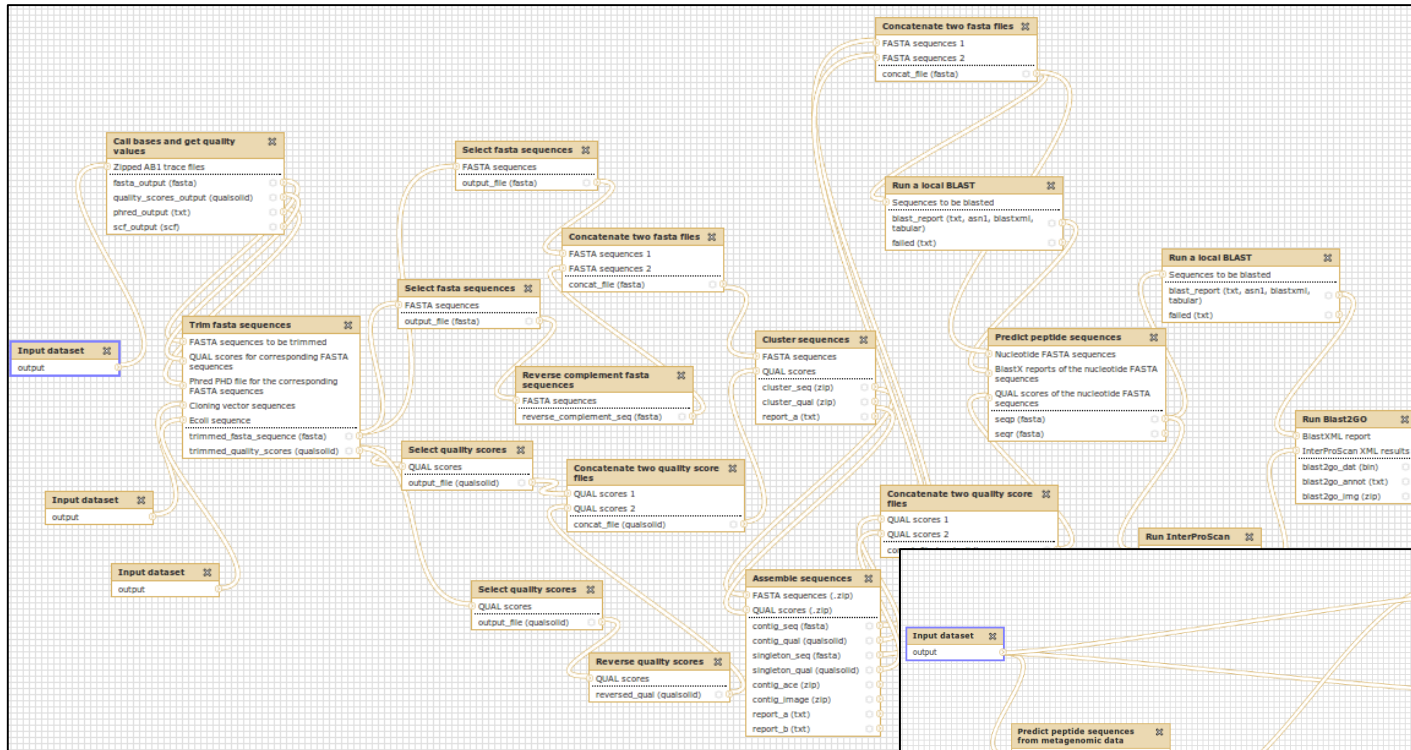


H3ABioNet

Pan African Bioinformatics Network for H3Africa

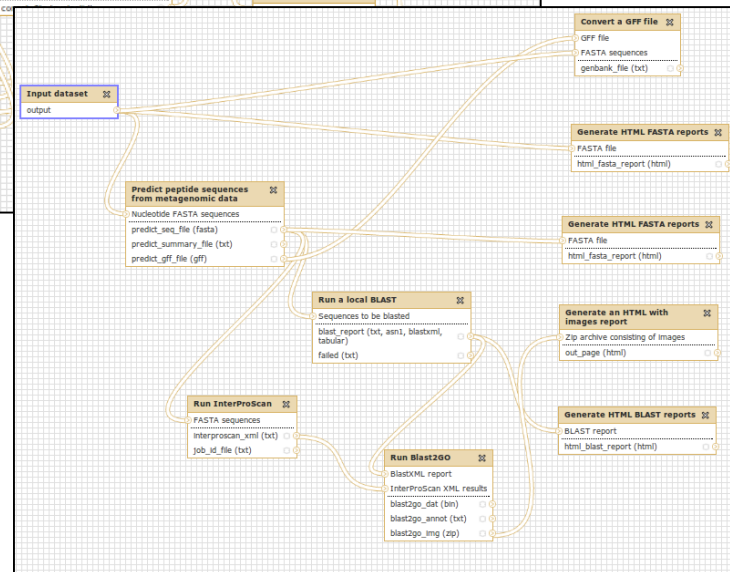


Galaxy workflows



EST assembly and annotation pipeline

Metagenomics functional annotation pipeline



H3ABioNet

Pan African Bioinformatics Network for H3Africa



Computational Biology @ UCT



eBiokits

- Standalone hardware without the need for internet
- Includes databases, tools, tutorials
- Based on Mac hard drives



H3ABioNet

Pan African Bioinformatics Network for H3Africa



Visualization tools

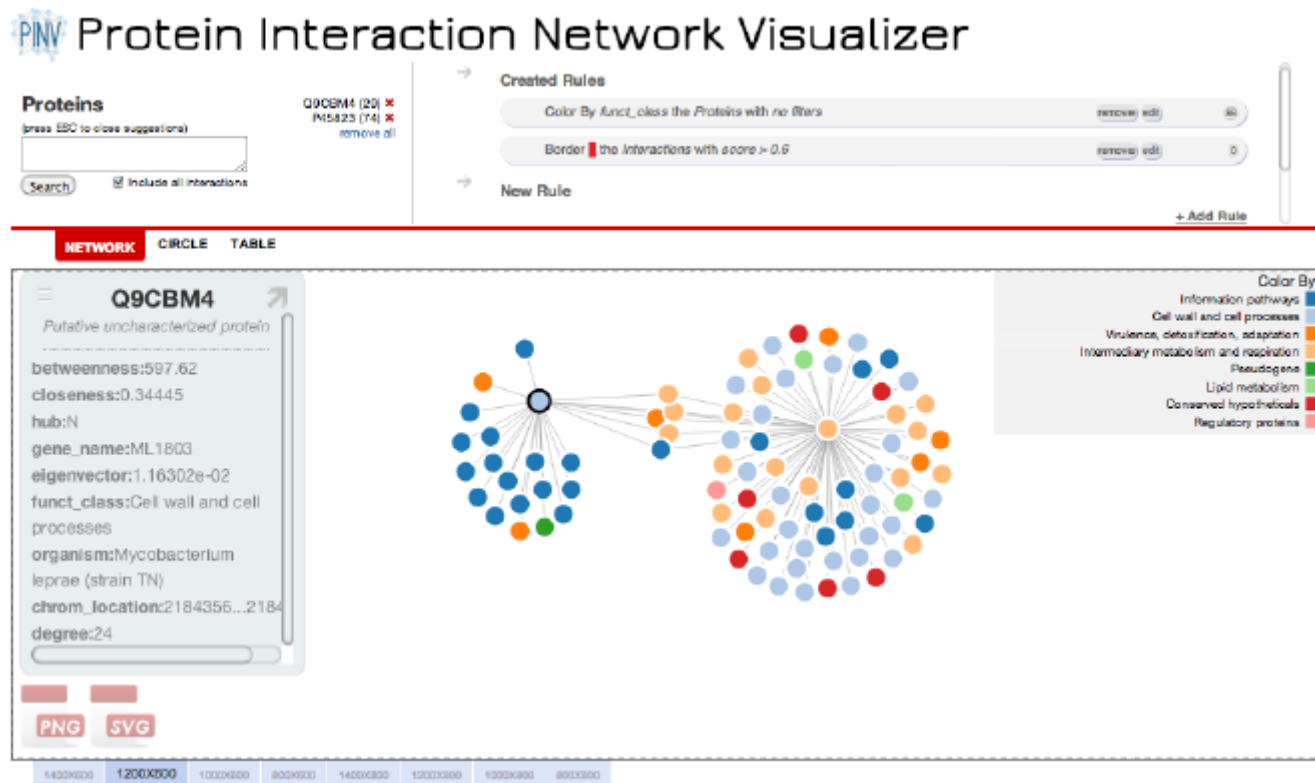
- DAS-based tools –avoid integrating data locally

The screenshot displays the Dasty2 web interface for protein P05067. The browser address bar shows the URL: `http://192.168.2.146/dasty2_release/biosapiens.html?q=P05067&label=BioSapiens&t=3`. The interface includes a search bar with the protein ID and registry label, a 'CHECKING' progress bar at 100%, and a 'QUERY INFORMATION' section showing the sequence ID and length. The 'SEQUENCE' section displays the amino acid sequence in orange text. Below this, the 'FEATURES DETAILS' section shows 'SELECTED DAS SOURCES' with checkboxes for various databases like signalp, netphos, and uniprot. The 'MANIPULATION OPTIONS' section allows for sorting and zooming. The 'POSITIONAL FEATURES' section at the bottom provides a visual representation of the protein's features, including signal peptides and mature protein regions, with corresponding DAS sources listed on the right.



Visualization tools

- Interactive web-based network visualization tool



H3ABioNet training plan

- Short specialised courses
- Research driven workshops
- Regional workshops
- Career development workshops
- Internships, KTP
- Train the trainer program
- Mentorship
- Postgrad degree curriculum development
- Facility for live-streaming to other classrooms



H3ABioNet

Pan African Bioinformatics Network for H3Africa



Training activities

- Organized training courses:
 - Technical course (Linux, Cloud, data security, HPC)
 - Train-the-trainer course (Biostats, GWAS, NGS)
 - eBiokit course (NGS)
- Courses live streamed to 2 class rooms
- Set up online application system
- Set up online evaluation system
- Set up training course website
- Developed training course planning document



H3ABioNet

Pan African Bioinformatics Network for H3Africa



Summary

- Bioinformatics is required to convert data into information
- As biology goes more high-throughput so the need for mathematics, statistics and computing increases
- As well as ensuring adequate IT infrastructure is in place, we need to train people at different levels on the use of it
- If we can build these in Africa we will stop the movement of data off the continent!



H3ABioNet

Pan African Bioinformatics Network for H3Africa

