**H3ABioNet Montly Newsletter Issue 3: July 2013**
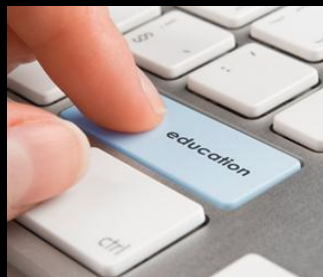
# Foreword

Welcome to the July 2013 newsletter the last official letter for year 1 of the grant! It is amazing to think the first year of the grant is now finished and we need to start working towards all the year 2 milestones. To reflect on year 1, it has been extremely busy and I believe we have had many significant achievements. This was evident when completing the annual report showing completed milestones for year 1 submitted to the NIH in June. The report had an appendix of 50 pages detailing the outcomes of our activities. We have also had challenges and I hope that year 2 will see us ironing these out to ensure that our second year is even more successful than the first one. To get back to the month at hand, July saw the successful completion of the Train-the Trainer Workshop at Icipe, Kenya, which started out as a Masters training course. The course was hosted by Anne Fischer who did an excellent job at coordinating the lives of over 20 participants for 3 weeks! A report on the course and some photos are provided later in this newsletter. This course was followed by a 1 week eBioKit workshop presented by several lecturers, and organised by Prof Erik Bongcam-Rudloff. More details on this course will appear in the next bulletin.

All the working groups have been busy this month revisiting their milestones, so instead of the usual monthly report, this time we show tables of the milestones. Key for Milestones Tables: E= Education and Training Working Group; I = Infrastructure Working Group; R = Research and Tool Development Working Group; U = User Support Working Group; SP = Specific Research Projects; All = All Working Groups.

The Node Assessment working group has progressed well under the excellent leadership of Prof Victor Jongeneel, and the datasets and SOPs are almost ready for launching of the exercise at the next H3ABioNet consortium meeting. In addition to their report, this bulletin provides a feature on the SNAP SNP Analysis Pipeline being developed by the RUBi Node in Grahamstown.

H3ABioNet Central has been undergoing some renovations to create a dedicated H3ABioNet area for the helpdesk staff. Despite the usual delay in completion, these staff have now moved into the new offices. They have in the mean time, been working hard to help nodes complete their hardware orders with Dell. Finally, H3ABioNet has been represented at the GOBLET meeting, which preceded ISMB in Berlin, as well as at a special session at ISMB on "How to run a success bioinformatics network". These are reported on in more detail later.

Prof. Nicky Mulder.

**H3ABioNet Montly Newsletter Issue 3: July 2013**

- **Education and Training**

**&**

The Education and Training Working Group (E&TWG) as with all working groups has been busy in formalising their milestones for the next project period.

| Milestone | Start Date | End Date | WG |
|---|---|---|---|
| Survey of training requirements (part of survey in aim 1) | Dec-12 | Feb-13 | E |
| Gather info on other existing courses | Dec-12 | Feb-13 | E |
| Gather info on training capacity of each node | Dec-12 | Feb-13 | E |
| Set up workshop schedule | Feb-12 | Apr-13 | E |
| Set up course procedures for Cyber-based Learning of Bioinformatics | Dec-12 | Feb-13 | E / I |
| New students for open bursaries recruited and selected | Dec-12 | Feb-13 | E |
| Set up training schedule/plan for each person, retention plan/career development | Feb-13 | Apr-13 | E |
| Education committee develop MSc program curriculum | Dec-12 | Apr-13 | E |
| First MSc program courses run | yr 1 | | E |
| Internship application process and plan | Dec-12 | Jan-13 | E |
| Select interns for year 1 | Jan-13 | Apr-13 | E |
| Set up co-supervision/mentoring programme | Jan-13 | Oct-14 | E |
| Set up template/system for monitoring training quality | Dec-12 | Apr-13 | E / U |
| WORKSHOP: NIH grants management | Nov-12 | Jan-13 | E |
| WORKSHOP: Introductory computing for technicians | TBD, yr 1 | | E |
| WORKSHOP: Train the trainer (coincide with H3A Cons meeting) | TBD, yr 1 | | E |
| WORKSHOP: Data management for researchers, clinicians (coincide with stakeholders meeting) | TBD, yr 1 | | E |
| WORKSHOP: Intro to HPC and Cloud computing | TBD, yr 2 | | E |
| WORKSHOP: Career development for postdocs and mid-career fellows | TBD, yr 2 | | E |
| WORKSHOP: Data analysis with eBioKits (coincide with network meeting) | TBD, yr 2 | | E |
| WORKSHOP: Computing summer school | TBD, yr 2 | | E |
| WORKSHOP: NGS data analysis and Galaxy (coincide with consortium meeting) | TBD, yr 2 | | E |
| Management of course materials | Ongoing | | E |
| Reports on training | Annual | | E |
| 2nd Annual MSc courses run | yr 2 | | E |
| Select interns for year 2 | yr 2 | | E |
| Finalize milestones and deliverables based on project needs | Jan 13 | Mar 13 | E |
| Provide report to MC | Monthly | | E |
| Provide 6 months NIH working group report | Nov 13 | May 14 | E |

**Dr. Nash Oyekanmi.**

Infrastructure

• Infrastructure

# Infrastructure

The following milestones and deliverables have been set for the Infrastructure working group:

| Milestone | Start Date | End Date | WG |
|---|---|---|---|
| Hardware in core centres set up with back-ups in place | Jan-13 | Mar-13 | I / R/ U |
| Set up course procedures for Cyber-based Learning of Bioinformatics | Dec-12 | Feb-13 | E / I |
| eBioKit: development -add new tools and tutorials | Apr-13 | Oct-14 | E / I / U |
| Assess and document computing facilities available at nodes, develop plan to improve where required | Jan 13 | Nov 13 | I |
| Set up mirror sites for major databases | Jan-13 | Jul-13 | I |
| Galaxy/BioMart: data exchanging between BioMarts and Galaxy developed | Nov-14 | Oct-14 | I |
| Ping exercise to assess internet capacity | Feb-13 | Mar-13 | I |
| Engage with vendors for improved internet terms | Feb-13 | Oct-14 | I |
| Hardware centres - installation & testing of software | Mar-13 | May-13 | I |
| Investigate Computing technologies (Cloud, GRID, HPC) | Mar-13 | Oct-14 | I |
| Work with HPCs to get data and tools available | Mar-13 | Oct-14 | I |
| Investigate data storage on Cloud and GRID | Mar-13 | Mar-14 | I |
| Hardware centres - installation & testing of software | Mar-13 | May-13 | I |
| Galaxy: user requirements defined | Dec-12 | Oct-13 | I / R |
| Galaxy: required tools all integrated into each installation | Apr-13 | Apr-14 | I / R |
| Galaxy installed at CHPC and core hardware centers | Feb-13 | Jun-13 | I / R |
| Investigate other solutions, e.g. Chipster | Apr-13 | Oct-14 | I / R / U |
| Develop/customize GUI for BioMart | Oct-13 | Apr-14 | I / R / U |
| Data submission tools developed | Apr-14 | Oct-15 | I / R / U |
| Galaxy installed at CHPC and core hardware centers | Feb-13 | Jun-13 | I / U |
| eBioKit synchronization system set up and running | Apr-13 | Oct-14 | I / U |
| Work with HPCs to get data and tools available | Mar-13 | Oct-14 | I / U |
| Central server for eBioKit set up | Jan-13 | Apr-13 | I / U / E |
| Guidelines and SOPs for data management developed | Oct-13 | Oct-14 | I / U / R |
| Guidelines and SOPs for data management developed | Oct-13 | Oct-14 | R / I |
| Finalize milestones and deliverables based on project needs | Jan 13 | Mar 13 | I |
| Provide report to MC | Monthly | | I |
| Provide 6 months NIH working group report | Nov 13 | May 14 | I |

Dr. Alia Benkahla.                    Prof. Scott Hazelhurst.

# Research and Tool Development

- **Research and Tool Development**

The research and tool development working group's goals are as follows:

| Milestone | Start Date | End Date | WG |
|---|---|---|---|
| Hardware in core centres set up with back-ups in place | Jan-13 | Mar-13 | R / U / I |
| Galaxy: user requirements defined | Dec-12 | Oct-13 | R / I |
| Galaxy: required tools all integrated into each installation | Apr-13 | Apr-14 | R / I |
| Galaxy installed at CHPC and core hardware centers | Feb-13 | Jun-13 | R / I |
| Investigate other solutions, e.g. Chipster | Apr-13 | Oct-14 | R / U / I |
| Develop/customize GUI for BioMart | Oct-13 | Apr-14 | R / U / I |
| Data submission tools developed | Apr-14 | Oct-15 | R / U / I |
| Guidelines and SOPs for data management developed | Oct-13 | Oct-14 | R / U / I |
| Set up research subcommittee, meeting of personnel on joint projects | Dec-12 | Quarterly | R |
| Develop Galaxy workflows for H3Africa projects (Specifications) | Oct-12 | Apr-13 | R / SP |
| Develop genome assembly pipeline for NGS (Specifications) | Oct-12 | Apr-13 | R / SP /N |
| Develop DAS-based visualization tool (Specifications) | Oct-12 | Apr-13 | R / SP |
| Develop Admixture mapping tool (Specifications) | Oct-12 | Apr-13 | R / SP |
| Develop Structural and functional SNP-calling tools/pipelines (Specifications) | Oct-12 | Apr-13 | R / SP |
| Develop Grid-based tool for data storage and sharing (Specifications) | Oct-12 | Apr-13 | R / SP |
| Develop Recombination tool (Specifications) | Oct-12 | Apr-13 | R / SP |
| Patient DB: guidelines for patient data storage and database development started | Jan-13 | Apr-13 | R / SP |
| Patient DB: existing ontologies, CVs sourced | Jan-13 | Apr-13 | R / SP |
| Generic framework for building a database for clinical data developed | Apr-13 | Apr-14 | R / SP |
| Develop pipelines for processing GWAS data into BioMart | Apr-13 | Oct-14 | R / SP |
| BioMart implementation plan developed | Dec-12 | Jan-13 | R / U |
| BioMart: all public data in a central shared BioMart | Jan-13 | Ongoing | R / U |
| Guidelines and SOPs for data management developed | Oct-13 | Oct-14 | R / U |
| Evaluate user support systems | Oct-12 | Nov-12 | U / R |
| Finalize milestones and deliverables based on project needs | Jan 13 | Mar 13 | R |
| Provide report to MC | Monthly | | R |
| Provide 6 months NIH working group report | Nov 13 | May 14 | R |

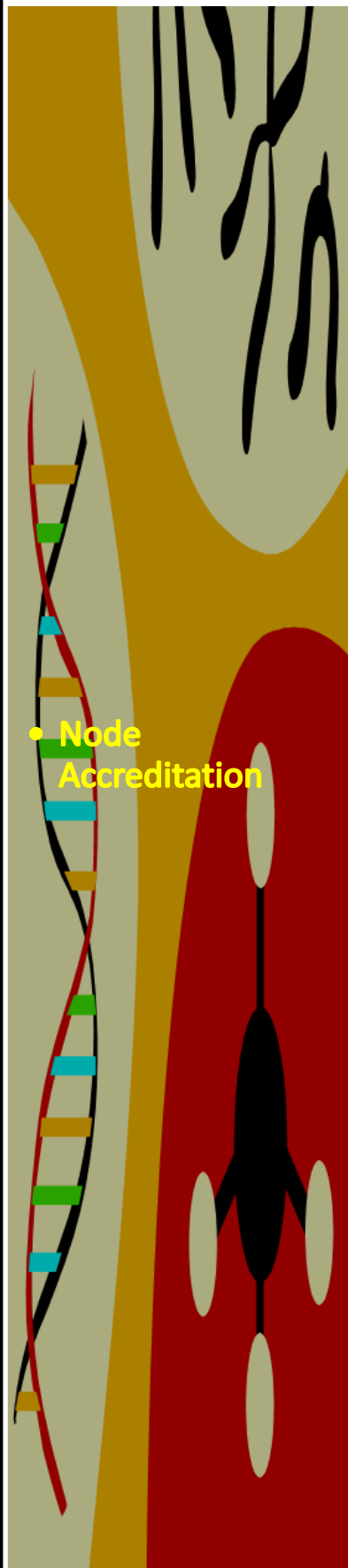Dr. Julie Makani.            Prof. Ezekiel Adebiyi.

- User Support

# User Support

The user support working group will work on the following set of milestones for this period:

| Milestone | Start Date | End Date | WG |
|---|---|---|---|
| Hardware in core centres set up with back-ups in place | Jan-13 | Mar-13 | U / R / I |
| Investigate other solutions, e.g. Chipster | Apr-13 | Oct-14 | U / R / I |
| Develop/customize GUI for BioMart | Oct-13 | Apr-14 | U / R / I |
| Data submission tools developed | Apr-14 | Oct-15 | U / R / I |
| Guidelines and SOPs for data management developed | Oct-13 | Oct-14 | U / R / I |
| BioMart implementation plan developed | Dec-12 | Jan-13 | U / R |
| BioMart: all public data in a central shared BioMart | Jan-13 | Ongoing | U / R |
| Guidelines and SOPs for data management developed | Oct-13 | Oct-14 | U / R |
| Evaluate user support systems | Oct-12 | Nov-12 | U / R |
| eBioKit: development -add new tools and tutorials | Apr-13 | Oct-14 | U / I / E |
| Galaxy installed at CHPC and core hardware centers | Feb-13 | Jun-13 | U / I |
| eBioKit synchronization system set up and running | Apr-13 | Oct-14 | U / I |
| Work with HPCs to get data and tools available | Mar-13 | Oct-14 | U / I |
| Central server for eBioKit set up | Jan-13 | Apr-13 | U / I / E |
| Co – opt more expertise for helpdesk | Jan-13 | Ongoing | U / ALL |
| Set up licensing agreements on commercial software (Ingenuity, Visual Analytics) | Jan-13 | Mar-13 | U / ALL |
| Set up template/system for monitoring training quality | Dec-12 | Apr-13 | U / E |
| Galaxy: user requirements defined | Dec-12 | Oct-13 | U |
| Audit Helpdesk  (once every 6 months) | May-13 | Ongoing | U |
| Bioinformatics help desk FAQ set up | Oct-13 | Ongoing | U |
| Decide type of user support system to implement | Nov-12 | Dec-12 | U / C |
| Bioinformatics help desk & tracking system set up | Dec-12 | Feb-13 | U / C |
| Create common framework for help desk | Jan-13 | Mar-13 | U / C |
| Finalize milestones and deliverables based on project needs | Jan 13 | Mar 13 | R |
| Provide report to MC | Monthly | | R |
| Provide 6 months NIH working group report | Nov 13 | May 14 | R |

Dr.  Judit Kumuthini.          Dr. Jonathan Kayondo.

- **Node Accreditation**

# Node Accreditation

## Generation of synthetic datasets

The NAWG will be using synthetic datasets rather than real ones, as they can be easily produced in unlimited numbers, are impossible to match to known data, and can be made to all contain the same density of information (frequency of variants, types of variation, etc).

To produce artificial NGS datasets, we have been using the following procedure:

- We use as input a single human chromosome of average size and gene density rather than a full genome; the input will be selected from one of the African genomes from the 1000 Genomes collection.
- First, we generate a list of random variants relative to the reference African genome. The Genome Smasher scripts from Steve Hart at the Mayo Clinic (https://code.google.com/p/genome-smasher/) generates a VCF file with random SNPs or CNVs from an input FASTA file.
- The reference FASTA and VCF files are processed by Genome Smasher to produce an output FASTA file that contains the variants described in the VCF file.
- The new FASTA file containing the variants is then converted to simulated exome reads using the Wessim package written by Sangwoo Kim at UCSD (http://sak042.github.io/Wessim/).

Wessim takes as input a genome sequence (in our case the mutated FASTA above) and information about the target regions (exons) as a capture probe set. Various parameters describing the exon capture and sequencing steps can be specified. The output is a FASTQ file with simulated Illumina reads.

We have successfully installed and tested both Genome Smasher and Wessim on our local cluster at the University of Illinois HPC. Wessim requires a number of programs and libraries to be pre-installed: pysam library, numpy library, gfServer and gfClient, faToTwoBit, samtools and GemSim error models, and needs to find them via the relevant paths on the installation / server. Genome Smasher scripts can be downloaded as a gzipped tar file (careful, the extension is .gz rather than .tar.gz) and only requires Perl.

The VCF files generated by Genome Smasher contain the "true" variation documented by the simulated reads and thus, can be compared directly to the VCF files generated by the candidate nodes. This will preclude any bias in the interpretation of the data.

Since we are using an African genome as a reference rather than the standard human reference assembly (GRCh37), the variants contained within this African genome, which are documented in dbSNP, will have to be identified by the candidates. The additional variants introduced by Genome Smasher will be random and not referenced in dbSNP. However, they will appear as "private" variants enabling the assessment panel to determine how successfully these variants were detected.

Dr. Victor Jongeneel.

# H3ABioNet Research Focus - The SNP Analysis Pipeline (SNAP)

SNAP is a web-based workflow management system that is currently under development in the Rhodes University Bioinformatics (RUBi) research group in Grahamstown, South Africa (http://rubi.ru.ac.za/). It will provide a range of bioinformatics tools that will initially be targeted at homology modelling and SNP analysis. Being a workflow management system, SNAP will provide users with the ability to string multiple tools together in complex, non-linear pipelines. Users can then save these pipelines/workflows for future use. At launch, a number of predefined pipelines will also be provided for users, who do not wish to build their own.

SNAP will initially be launched with several tools, organized into categories. With the initial goal being homology modelling and SNP analysis, the proposed initial categories are currently Template Selection, Sequence Alignment, Homology Modelling, Model Validation, Energy Analysis, GWAS, Visualization, and Data Manipulation. Amongst other things, these tools will allow proteins with unknown structures to be modelled, SNPs to be mapped to a location on the protein structure, and a stability analysis to be carried out to determine whether an amino acid change caused by a SNP might destabilize the protein. Homology modelling has also been found to be useful in drug discovery.

Over and above the tools that will be provided at launch, SNAP has been designed with the goal of allowing admin users to add additional tools and programs, as well as categories of programs, by simply filling in a few forms on a web page. This data will be stored in a MySQL database. Based on the input to these forms, HTML pages will be generated as interfaces to the tools. This ability will give SNAP an unprecedented level of flexibility and customizability, allowing users without much technical experience to add tools without difficulty.

SNAP will provide access to its tools in two ways. Firstly, access will be provided through a web interface for the general user. And secondly, programmatic access will be provided via a RESTful web API.

SNAP will also provide a social networking aspect that will allow users to collaborate with one another, sharing data and results. Data can be uploaded by one user and shared with his/her peers or be made public for all users on the server. This will also hopefully help to reduce the amount of data uploaded to and stored on the server. Users will also be able to share results with their peers as well as comment on results that have been shared with them, fostering an environment for collaboration.

Although the tools that SNAP will initially provide are focused on homology modelling and SNP analysis, they actually have a wide range of bioinformatics applications. Add this to the ease at which tools can be added to the system, and it is easy to see SNAP becoming a general bioinformatics workflow management system, where users will be able to satisfy all of their bioinformatics needs.
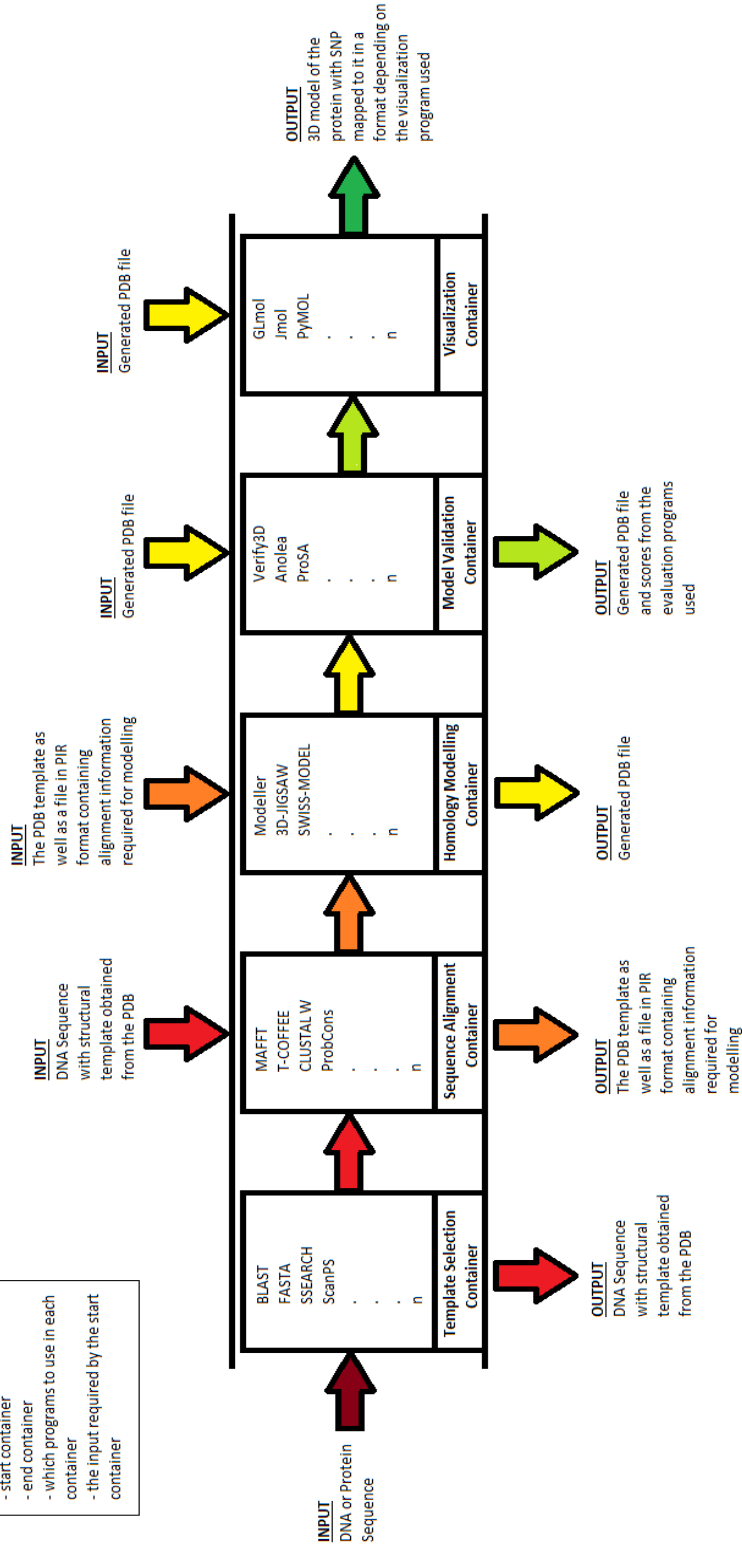
- **H3ABioNet Research Focus**

**H3ABioNet Montly Newsletter  Issue 3: July 2013**

# H3ABioNet Research Focus - The SNP Analysis Pipeline (SNAP)

- **H3ABioNet Research Focus**



**Schematic representation of the SNAP platform being developed.**

# H3ABioNet Research Focus - The SNP Analysis Pipeline (SNAP)

SNAP is being developed using Mono, a cross platform, open source .NET framework, and makes use of a MySQL database. It will be hosted over a cluster at Rhodes University. It is being developed by MSc student and software developer, David Brown, under the supervision of Dr. Özlem Tastan Bishop.

## How can You get Involved?

SNAP can be used to analyse other data types as well as SNPs. We are always looking to build collaborations and would also like to test SNAP with novel SNP data. If you have a novel SNP dataset that you would like to characterise structurally, any ideas for what types of bioinformatics analysis that you would like to see integrated into SNAP for collaborations then please do contact me at: ozlem.tastanbishop@gmail.com

Dr. Özlem Tastan Bishop.

- **H3ABioNet Research Focus**

# Train the Trainer (ToT) Bioinformatics Workshop - ICIPE 2013

The H3ABioNet Train the Trainer Bioinformatics workshop was held for 3 weeks at the International Center for Insect Physiology and Ecology (ICIPE) between 8th - 26th July. 22 participants from 14 different African countries were intensively taught a variety of bioinformatics topics ranging from Python programming, Biostatistics and R, Unix, population genetics, genome wide association studies, next generation analysis and the software programmes associated with them such as Plink, BWA etc. The Vidyo equipment used by the National Biotechnology Development Agency (NABDA) Node was transported to the ICIPE venue by Dr. Nash Oyekanmi and their consultant Carlos Lijeron. The Vidyo equipment enabled the lectures to be live streamed to participating classrooms in Tunisia and Nigeria for other scientists to follow the training and exercises with the last hour of each workshop day open to both the classrooms to ask questions.

The majority of participants found the workshop to be informative and the teaching to be of a high quality considering the amount of material covered in a relatively short period. The participants of the ToT Bioinformatics workshop will provide training and support in bioinformatics at their home institutions and it is hoped that they will form the nucleus of bioinformatics trainers for Africa.



H3 ABionet Train-The-Trainer Bioinformatics Course Participants

Left to right, Back Row: Ben Kulohoma, Jean-Baka Domelevo Entfellner
Left to right, Middle Row: Anne Fischer, Oumaru Samna, Benjamin Kumwenda, Anita Ghansah, Sunji Nadoma.
Left to right, Front Row: Mutawakil Saad, Moise Elegbe, Majdi Nagara, Deogratius Ssemwanga, Samar Kassim, Lerato Magosi, Angela Makolo, Bruno Mmbando, Samson Pandam Salifu, Sylvester Lyantagaye, Bonface Nganyi, Morne Du Plessis, Twaha Mlwilo, Emmanuel Oga.
Seated Front Row: Dr. Nash Oyekanmi.

- Train the Trainer (ToT) Bioinformatics Workshop - ICIPE 2013

**H3ABioNet Montly Newsletter Issue 3: July 2013**

## Global Organization for Bioinformatics Learning, Education and Training (GOBLET) July 2013 meeting feedback

The African Society for Bioinformatics and Computational Biology and the CPGR as members of GOBLET (http://mygoblet.org) attended its third meeting that was held in Berlin. Dr. Judit Kumuthini represented the CPGR and ASBCB as proxy for Prof. Nicola Mulder. On 19th of July, at Berlin ISMB/ECCB 2013, the GOBLET meeting commenced with a brief 'tour de table', during which participants introduced themselves and their organisations.

GOBLET is a legally registered foundation whose mission is to provide a global, sustainable support structure for bioinformatics trainers and trainees; facilitate capacity development in bioinformatics in all countries; contribute to the development of standards and guidelines for bioinformatics education and training; act as a hub for fund gathering; reach out to school teachers, bridge the gap to the next generation of bioinformaticians and to foster the international community of Bioinformatics, Biocomputing, Biocuration and Computational Biology (B$^3$CB) trainers.

The meeting in Berlin proper began with a review of 12 tasks, most of which were now complete or in progress. These included writing and circulating the kick-off meeting report; establishing a bank account and sending out the first invoices; progressing the website based on feedback from the kick-off; submitting a proposal for funds to support a meeting in Toronto.

The principal business goals of the meeting were:

i) To provide an update on work carried out since the KO, including reports from three task-force rapporteurs.

ii) To continue previous discussions and finally agree a framework for GOBLET's fee/benefit structure.

iii) To agree nomination and election processes for its future Executive Board members and Committee Chairs.

In terms of outputs, a paper had recently been accepted in *Briefings in Bioinformatics:* Best Practices in Bioinformatics Training for Life Scientists, by Via *et al*. Although not a GOBLET article *per se*, it introduces GOBLET as the natural evolution of the BTN. In addition, GOBLET also worked closely with the ISCB in order to create an education poster track for ISMB 2013 with 20 posters being accepted for International Conference on Intelligent Systems for Molecular Biology (ISMB) and the European Conference on Computational Biology (ECCB).

- **GOBLET July 2013 Meeting Feedback**



**Attendees of the GOBLET meeting hosted by the ISCB at the Berlin Hilton.**

**Dr. Judit Kumuthini.**

- **ISMB/ECCB July 2013 Meeting Feedback**

# ISMB/ECCB 2013
# Meeting Feedback

The 21[st] annual International Conference on Intelligent Systems for Molecular Biology (ISMB 2013) was jointly held with the 12[th] Annual European Conference on Computational Biology (ECCB), under the name of ISMB/ECCB 2013 (http://www.iscb.org/ismbeccb2013). The conference took place in Berlin, Germany between July from July 21-23, 2013, where over 1500 attendees participated.

The main conference included special session tracks for experimental biology and related inter-disciplines, highlight tracks for work published in journals frequented, technology track and late breaking research track for the presentation of submitted abstracts of research in progress. One of the highlights of the conference was keynotes delivered by 6 researchers of the highest international esteem who informed the community of historical perspectives and landmark advances in computational and experimental research, and injected new directions into the field of computational molecular biology (http://www.iscb.org/ismbeccb2013-program/ismbeccb2013-keynotes/ismbeccb2013-all-keynotes).

Although the ISMB/ECCB 2013 conference provided a multidisciplinary forum for disseminating the latest developments in bioinformatics and advanced computational biology, it allowed seniors and junior researchers and students to have the most opportunities of finding collaborators, jobs, mentor-ships.This is a successful of education training in bringing world class researchers together, sharing their experiences with young and junior trainees. However the number of members of H3ABioNet attending this conference was underestimated, roughly less than 10.

Prof. Nicky Mulder was invited to give a talk on H3ABioNet at a special workshop at ISMB on "How to run a successful bioinformatics network". The session included talks on some existing networks based around specific projects or geographically linked sites, and a panel discussion on how to run a bioinformatics network. I presented an overview of H3ABioNet, some early successes of the network and some comments on what has the potential to make it successful.

The CPGR presented a poster on KTP (Knowledge Transfer Programme) and took this as an opportunity for Dr. Kumuthini to lead a separate meeting with scientific advisory committee / reviewer committee meeting to discuss the way forward and the KTP launch.

Dr. Judit Kumuthini.

**H3ABioNet Montly Newsletter Issue 3: July 2013**

- Important Dates

# Important Dates

- 30th August – deadline for registering for H3ABioNet Consortium meeting

- 9th September, 2013 – deadline for registering for the H3Africa Consortium meeting

- 31st August, 2013 – deadline for Early bird registration to the South African Society for Human Genetics Conference, Johannesburg, South Africa ( http://www.sashg2013.co.za/)

- 5th September, 2013 – deadline to fill in NetCapDB details – all Nodes

- 1st -3rd of October, 2013 - H3ABioNet Annual Meeting, Johannesburg, South Africa

- 3rd - 6th of October, 2013 - H3Africa Consortium Meeting, Johannesburg, South Africa

- 6th - 9th of October 2013, South African Society for Human Genetics Conference, Johannesburg, South Africa ( http://www.sashg2013.co.za/)