# Proceedings
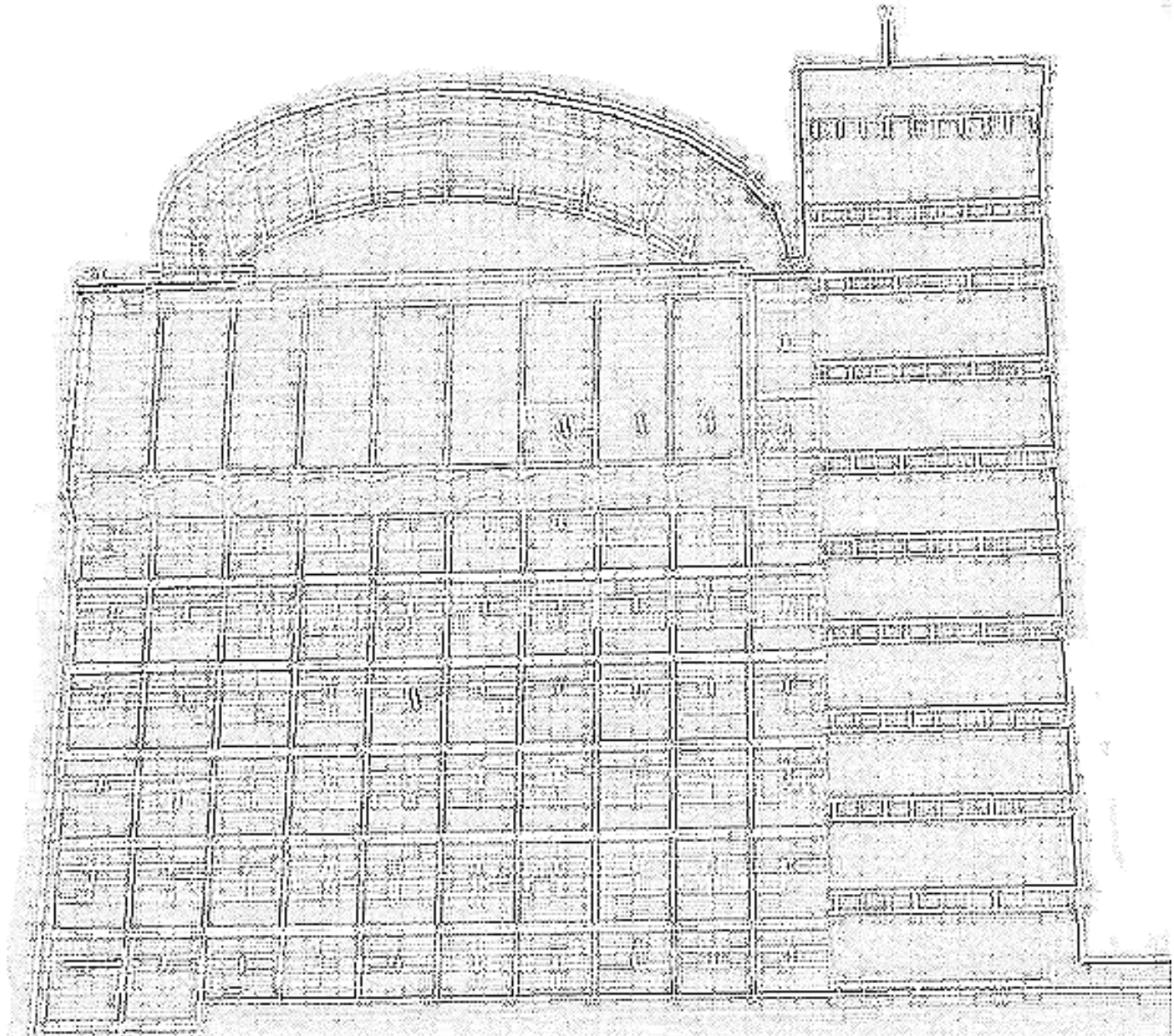# of the
# First Southern African Bioinformatics
# Workshop

University of the Witwatersrand, Johannesburg

28–30 January 2007

# Welcome

It is with great pleasure that we warmly welcome you to the First Southern African Bioinformatics Workshop in the summer of 2007 at the University of the Witwatersrand in Johannesburg. Wits Bioinformatics is pleased to host this inaugural workshop.

As a discipline, bioinformatics is still relatively young in South Africa and we were encouraged by the excellent response to the call for papers. The fact we have been able to attract almost 100 delegates and five international invited speakers is a mark of how bioinformatics has progressed in South Africa. Two of our invited speakers, Alan Christoffels and Janet Kelso, are PhD graduates from the South African National Bioinformatics Institute at the University of the Western Cape and we are delighted with their involvement in this workshop. Eric Rivals joins us from the University of Montpellier and CNRS and Kateryna Makova and Anton Nekrutenko are both from Penn State University.

The scientific programme has been arranged according to themes based on the submitted papers and include: New Computational Techniques; Toolkits, Integration and Services; Gene Identification and Expression; Evolution and Phylogenetics; and Transcriptomics. We have also planned panel discussions to rigorously debate the issues of mathematics teaching and education in preparation for a career in bioinformatics and the role of bioinformatics in high performance computing.

The workshop has been generously funded by the University of the Witwatersrand. The workshop is hosted by Wits Bioinformatics, situated on the 12th and 13th floors of the University Corner Building with magnificent views stretching south over the Nelson Mandela Bridge and north towards the tree-rich suburbs. We hope delegates will have an opportunity to visit us.

We have also had generous support to fund delegates from the national Department of Science and Technology, and the Centre National de la Recherche Scientifique (within the framework of the French support of NEPAD's scientific programme). It is supported by the National Bioinformatics Network (NBN) under whose guidance this and future workshops will be held.

32 papers were submitted for reviewing and 26 were finally accepted. Of these, 9 papers were submitted in the *full research paper* category. The full research papers were reviewed by 3-5 reviewers, and finally 3 of these papers were accepted in the full research paper category. The other papers were all sent to external reviewers and considered by a sub-committee of the programme committee. With the exception of one paper, all papers were reviewed by between 2 and 4 referees. In addition over 40 posters have been submitted and are included in the programme.

This is the first of a series of annual bioinformatics workshops that will take place across southern Africa in the years to come. May you embrace this as an opportunity to get to know like-minded scientists, to make new friends and to enjoy the ambiance of our lovely University. You will be richly rewarded if you take a little time to visit the Origins Centre and reflect on the growth and evolution of our species. We wish you a challenging, informative and enjoyable workshop and look forward to your active participation to make this a truly memorable event.

> Scott Hazelhurst and Michéle Ramsay
> Programme Committee Co-Chairs
> January 2007

## Programme and Organising Committees

—Programme Committee
  —Co-chairs: Scott Hazelhurst & Michèle Ramsay, Wits.
  —Alia BenKahla, Institut Pasteur, Tunis; Judith Bishop, Pretoria; Greg Blatch, Rhodes; Junaid Gamieldien, NBN Central; Winston Hide, SANBI; Jannie Hofmeyr, Stellenbosch; Dan Jacobson, NBN Central; Fourie Joubert, Pretoria; Heikki Lehväslaiho, SANBI; Daniel Masiga, ICIPE, Nairobi; Nicky Mulder, Cape Town; Hugh Murrell, KwaZulu-Natal; Hugh Patterton, Free State.
—Organising committee
  —Scott Hazelhurst (Chair), Jonathan Burke, Pierre Durand, Brenda Lacey-Smith, Khayeni Ndlovu, Andries Oelofse, Michèle Ramsay, Christine Rey

## Reviewers

We thank the following people who were responsible for reviewing the papers.

—Alain Denise, Paris-Sud

—Alia BenKahla, Institut Pasteur, Tunis

—Allen Rodrigo, Auckland

—Anna Kramvis, Wits

—Cathal Seoighe, Cape Town

—Christine Rey, Wits

—Darren Martin, Cape Town

—Dean Goldring, KwaZulu-Natal

—Fourie Joubert, Pretoria

—François Coste, IRISA, Rennes

—Frank Dehne, Carleton

—Hugh Murrell, KwaZulu-Natal

—Ian Sanders, Wits

—Judith Bishop, Pretoria

—Karen Megy, EBI

—Michèle Ramsay, Wits

—Nicky Mulder, Cape Town

—Nir Oren, Aberdeen

—Philip Machanick, Queensland

—Pierre Durand, Wits

—Scott Hazelhurst, Wits

—Steven Orzack, Freshpond Institute

—Tracy McLellan, Wits

—Winston Hide, SANBI/UWC

Contents

# VI  Indices  115

# Author Index  115

# Topic Index  118

# An Efficient Algorithm for Oligonucleotides Selection in a Large EST Databases

Ezekiel Adebiyi

Department of Computer and Information Sciences, Covenant University, PMB 1023, Ota, Nigeria.

ABSTRACT

Identifying unique oligonucleotide (oligo) probe sequences is an important step in PCR and microarray experiments. While there are a growing number of complete and annotated genomes, the largest collection of publicly available genetic sequences are expressed sequence tag (EST) sequences. Furthermore, for many organisms that are important to the society, such as barley, the EST is the major data on the expressed genes in a number of these organisms. For the EST sequences, the unique oligo problem is the selection of oligos each of which appears (exactly) in one EST sequence but does not appear (exactly or approximately, for a given hamming difference $d$) in any other EST sequence.

OligoSpawn, in two phase, has been implemented to efficiently select oligos from ESTs. The notion of a "seed" was used in the construction of OligoSpawn, and its run time is exponential dependent on $q$ (the length of the "seed"). For $q = 11$, it ran on a previous barley dataset of 28MB for 2 hours and 26 minutes using a 1.2GHz AMD machine, but it is very inefficient for large datasets, like the new 43MB barley dataset. We observed this as OligoSpawn, for $q = 11$, runs for about 6 days using a 3.0GHz Pentium IV machine. Furthermore, selection of some important unique oligos (*e.g.*, for which $q = 13$) is unwieldy for OligoSpawn.

In this work, using the suffix tree, we give a careful theoretical characterization of the set of seeds required, and prove a subqradratic time algorithm for extracting these seeds. Using this result, we present an efficient algorithm that takes advantage of the new results, that simplify the solution of the least common ancestor (LCA) problem via the range minimum query (RMQ) problem. The run time of our resulting algorithm is $O(n^3 q d / 4^{2q})$. For $q = 11$ and $q = 13$, our algorithm runs on the new 43MB barley dataset for 4 days using also a 3.0 GHz Pentium IV. As far as we know, our algorithm is the fastest oligonucleotides selector algorithm for large databases of tens of thousands of EST sequences, such as the barley ESTs.