





Introduction to bioinformatics

DNA technology course

Barbera van Schaik Bioinformatics Laboratory Academic Medical Centre (AMC) b.d.vanschaik@amc.uva.nl



Information Management / e-Science

Lost in translation (biology and informatics)



Definitions of bioinformatics

Adapted definition according to Wikipedia

The application of information technology and statistics to the field of molecular biology.

The creation and advancement of databases, algorithms, computational and statistical techniques, and theory to solve formal and practical problems arising from the management, analysis and interpretation of biological data.

What is bioinformatics?

Extraction of biological knowledge from complex data



What is bioinformatics?



... one of the results *might* be a tool you can use

Bioinformatics

*convert data to knowledge *generate new hypotheses



*Design new experiments

And information management......

How does molecule A interact with protein B?



A schematic visual model of oxygen-binding process, showing all four <u>monomers</u> and <u>hemes</u>, and <u>protein chains</u> only as diagramatic coils, to facilitate visualization into the molecule. (<u>http://en.wikipedia.org/wiki/Hemoglobin</u>)

Study migration

Human mtDNA Migrations http://www.mitomap.org/pub/MITOMAP/MitomapFigures/WorldMigrations.pdf

Copyright 2002 © Mitomap.org



+/-, +/+, or -/- = Dde I 10394 / Alu I 10397 * = Rsa | 16329

Mutation rate = 2.2 - 2.9 % / MYRTime estimates are YBP

Study evolution



Which gene(s) causes disease X?





Genetics and population analysis

Gene expression

http://www.oxfordjournals.org/our_journals/bioinformatics/for_authors/scope_guidelines.html

Bioinformatics Groups and institutes

Bioinformatics Laboratory (AMC) Medical Bioinformatics

- Part of the KEBB
- You are welcome if you need bioinformatics expertise
- www.bioinformaticslaboratory.nl

Bioinformatics Laboratory

The Bioinformatics Laboratory was initiated in 1997 to advance biomedical research in the Academic Medical Center (AMC, Amsterdam) and is headed by Antoine van Kampen. The laboratory is part of the department of Clinical Epidemiology, Biostatistics and Bioinformatics (KEBB). The group members represent a broad range of (bio)informatics, biostatistics, and e-bioscience expertise and engage in research, education, and data analysis. Van Kampen is appointed part time at the Biosystems Data Analysis group of the Science Faculty to facilitate collaborations between the AMC, the BDA group, and other groups at the Science Faculty.





Focus

Systems genomics (top-down)

(Statistical) analysis of (multi-)omics data

Systems biology (bottom-up)

Mathematical modelling

Information management

Knowledge bases

Science gateways







Dutch Techcentre for Life Sciences (DTL)

- Support and Education ۲
- http://www.dtls.nl/ •



RECENT & UPCOMING

Events



Courses

Tweets lobs

ECCB

ECCB 2016: Extended submission deadlines

All deadlines for abstract submission to the European Conference on Computational Biology (ECCB) have been extended. The new deadlines are 5 April 2016 (Proceedings & Highlights) and... Continue reading →

WE ARE HIRING!



DTL seeks two Software Engineers Linked Data to provide leadership and innovative software development services. Candidates

PARTNER VIEWS: SURF



AnwarOsseyran, PAC member for SURF: "SURF is a great supporter of life sciences, and believes in good

Other bioinformatics organisations

🔰 Main o News

> Forums Organization What is EMBnet

eLearning
 EMBnet.news (New

Structure © Executive Boa

Education&Tra

Publicity&P

<u>Technical</u> <u>Management</u>

· Specialist No

Achievements

· EMBnet presentations

· How to apply?

Members

o <u>Facts</u> o <u>Role</u>

Community News

Subscribe here to the ISCB Community News

National

- European Bioinformatics Institute (E lacksquare
 - http://www.ebi.ac.uk/
- National Center for Biotechnology Information (NCBI)
 - http://www.ncbi.nlm.nih.gov/
- **EMBnet**

Massachusettes, USA, July 11 - July 13,

2010. The conference will feature an

exciting scientific program including

- http://www.embnet.org/
- International Society for ۲ Computational Biology (ISCB)
 - http://www.iscb.org/

EBI)	EMBL-EBI	EB-ere Search	Enter Text Here	Go Reset Advanc
	Databases Tools Data Resources & To EMBL-BANK = Genomes UniProt = Nucleotid	EBI Groups Training IOIS s = Gene Exp e Sequences = Molecula	Industry About Us Hell pression = Literature = S r Interactions = Taxonomy = P	equence Similarity attern & Motif Searc
	Artray-Express = Protein S Ensembl = Macromo InterPro = Small Mo PDBe	equences Reaction lecular Structures Protein F lecules Enzymes	s & Pattiways = Ontologies = S S amilies = Patent Resources = T = D = W	tructure Analysis ext Mining ownloads leb Services
			uropean Bioinformatic	s Institute
	About the EBI Research PhD Studies Training Industry Support Group & Team Leaders	 User Support EBI Mission People Events at the EBI Genome Campus Events 	Events Plant Bioinformatics 29-31 March 2010 Course full ENBO Practical Course in silico sy reconstruction, analysis and networ 10-13 April 2010 Registration nov	stems biology: netwo k based modelling' v open
S NCBI Resources 🛛 How To 🖓				My NCBI Sign In
National Center for Biotechno	ology Information	Search All Databases	▼ for	Search
Resources NCBI Home All Resources (A-Z) Literature DNA & RNA Proteins Sequence Analysis Genes & Expression Genomes	Nelcome to NCBI The National Center for Biotechnology Information advances science and neatth by providing access to biomedical and genomic information. Are about the NCBI Mission Organization Research R8S PubMed Central Free Full Text. Over 1,500,000 articles		And Popular Resources PubMed PubMed Central PubMed Central BLAST Gene Nucleotide Protein Gene Conserved Domains	
Maps & Markers Domains & Structures	from over 450 journals. Linked to Put and fully searchable.	oMed	Structure PubChem	
Genetics & Medicine	II 1 2 3 4		NCBI News nber and 02 De er News red: New Discovery-orio ed and NCBI Horrepag	o 2009 ented le. T
	Welcome to the EM	Bnet	News - 05 Oc	# 2009
Lister a a a a a a a binnye the combined expert binnye the combined expert of the organization and of this related to bioinformat EMBnet.news, the letter devo nodes. a a a a a a a a a a a a a	roup of collaborating nodes throughou tise of the nodes allows EMBnet to p is members. It provides the visitors wit is members. It provides the visitors wit sc. It also combines the services sted to provide information about wha	It Europe and a number of n cords services to the Europa vide services to the Europa the services of the Europa here soft the EMBret commu available on the nodes a t is happening at the nation:	iodes outside pan molecular inity and new nity and new al and special	



Jan 29: PLoS Computational Biology Article on Live

Coverage of ISMB/ECCR 2009

Bioinformatics

Extraction of biological knowledge from complex data

@AMC: genomics, metabolomics and proteomics data

What is genomics?

The application of **high-throughput** automated technologies to molecular biology.

OR

The experimental study of complete genomes.

Sample storage



DNA microarrays



High throughput sequencing

Ion Proton



Illumina MiSeq



Run: 2-4 hrs Data: 10 GB Reads: 60-80 million Length: 200bp

Run: 4-55 hours Data: 15 GB Reads: 25 million Length: 2x300bp

Illumina HiSeq (VUmc)



Run: 1-3.5 days Data: 1500 GB Reads: 5 billion Length: 2x150bp

2005-now: Next generation sequencing Millions to billions of sequences

DNA sequencing



Sanger sequencing

From DNA to Autoradiograms





Automated Sequencing



ABI 3100 Automated Capillary DNA Sequencer

Electropherogram



Sequencing factories



Whitehead institute

Custom-designed factory-style conveyor belt robots:

perform all functions from purifying DNA from bacterial cultures through setting up and purifying sequencing reactions.

Automated sequencing

2001

The HGP consortium publishes its working draft in *Nature* (15 February), and Celera publishes its draft in *Science* (16 February).

Human genome: 3 billion bases

In total 23 billion bases were sequenced (7.5-fold coverage)

This comprises 23 Gbyte of data

Public Human Genome project -1990-2003

- 3 billion US dollar

Private Celera genome project

-1998 - 2001

- Craig Venter

- 300 million US dollar 2001





Traditional versus high throughput DNA sequencing



http://en.wikipedia.org/

Sanger, one run:

? hours (human genome took 15 years)

1-384 sequences

300-1000 nt per sequence

1 KB - 384 KB data

= 1,000-384,000 bases

Year: 1963 - now



http://www.454.com/

 Roche 454, one run:
 III

 7.5 hours
 1

 1,000,000 sequences
 5 1

 750 nt per sequence
 2x

 35 GB data (including images)
 1...

 1 GB data (excluding images)
 Ye

 = 750,000,000 bases
 Year: 2005 – 2015



http://www.illumina.com

Illumina, one run:

1-3.5 days

5 billion sequences

2x 150 nt per sequence

1.5 TB data

Year: 2008 - now

Data size challenges

Falling fast

In the first few years after the end of the Human Genome Project, the cost of genome sequencing roughly followed Moore's law, which predicts exponential declines in computing costs. After 2007, sequencing costs dropped precipitously.



NextGen Sequencing a Game-Changer



COMPUTING Sequencing rate is higher Sequencing costs lower than Moore's law

STORAGE than data storage

Next Generation Genomics: World Map of High-throughput Sequencers

🖉 Show all platforms 🗖 Illumina GA2 🗖 Illumina HiSeq 🗖 Illumina MiSeq 🗖 Ion Torrent 🗖 PacBio 🗖 Polonator 🗖 Roche/454 🗖 SOLiD 🗖 Service Provider

Southern

Map

Guinea

Map data ©2012 Mapl

Genome projects

- Human genome project (1 genome)
- Exome sequencing (~10 individuals)
- Genome of the Netherlands (770 individuals)
- 1000 genome project (1000 individuals)
- 10K UK project (10,000 individuals)
 - Upgraded to 100K genomes

... many centers have one or more high throughput sequencers

Search

http://omicsmaps.com/

Google search the map

Bringing Order to Chaos with Bioinformatics

Data

plosion

image credit: Digital Vision, PhotoDisc, Matt Ray/EHP

Sign @ Wellcome-Sanger, Cambridge, UK

"oh, shit, that's 320TB!"

ony Cox, The Guardian, Feb 28 2001

Existing hardware



PC





Buy a bigger cluster (centralized model)







Cloud computing



Cloud Computing?



Not only hardware

- Software
- Peopleware
- Sharing

e-Science



More info: <u>http://www.egi.eu/</u> <u>http://www.ebioscience.amc.nl/</u>



Tips and tricks... where to start?

Biological databases Publicly available tools

Publicly available biological databases



Molecular Biology Database Collection 2015

- Published by Nucleic Acids Research (scientific journal)
- Currently ? databases.
- Latest issue new databases
- Coverage is far from exhaustive.



Database access

• Web-interface

- Keyword search
- Browsing / cross-linking
- BLAST
- From other applications
 - Application programming interface (API)
 - Web services
 - In-house / third party software
- Download
 - Full database
 - Specific records
 - Different formats (text, xml, rdf, etc)

Sequence analysis

Function prediction (similarity, sequence search) Localisation (genefinding) Grouping (genes, protein families) Conservation (motifs, functional blocks) SNPs and mutations (variations)

Sequence analysis

Pairwise alignment: in-exact matching of 2 sequences

CTCCTGAGGCAAATCTGTCAGTCCATCCTGGCTGAGTCCTCGCAGTCCCCGGCAGATCTTGAAGAAAAGA

Multiple sequence alignment: in-exact matching

of >2 sequences



Blast output



0.11

Distribution of 196 Blast Hits on the Query Sequence

Blast output - alignments

Score = 115 bits (58), Expect = 1e-22
Identities = 142/170 (83%)
Strand = Plus / Plus

Sbjct: 248 tggacgacctgcccaccgccctgtccgccctgagcgacctgcacgcccacaagctgcgtg 307

Sbjct: 308 tggaccccgtcaacttcaagctcctgagccactgcctgctggtgaccctggcttgccacc 367

```
Score = 58.0 bits (29), Expect = 2e-05
Identities = 52/59 (88%), Gaps = 3/59 (5%)
Strand = Plus / Plus
```

>gi|49420|emb|X57029.1|CAGLOBINM Length = 418
M.auratus mRNA for alpha globin chain

Score = 91.7 bits (46), Expect = 1e-15 Identifies = 151/186 (81%) Done ٠

Galaxy



Tip: Has instruction videos for a few use cases e.g. for analyzing next gen sequence data

https://galaxyproject.org/

Nucleic acid research - annual web server issue

bioinformatics.ca links directory

https://links.bioinformatics.ca/

Search Bioinformatics Links Directory	y Search Directory	Change Text
		Username:
Bioinformatics Links Dire	ectory	Password:
The Bioinformatics Links Directory features cura databases. The links listed in this directory are bioinformatics experts in the field. We also rely users for suggestions. Starting in 2003, we have NAR Webserver issue.	ated links to molecular resources, tools and selected on the basis of recommendations from on input from our community of bioinformatics e also started listing all links contained in the	Log in Create new Request new
Computer Related (82)	DNA (517) This category contains links to useful resources for DNA sequence analyses such as tools for comparative sequence analysis and sequence	Main Page
This category contains links to resources relating to programming languages often used in		Citations
bioinformatics. Other tools of the trade, such as		Acknowled
web development and database resources, are also included here.	assembly. Links to programs for sequence manipulation, primer design, and sequence retrieval	News
	and submission are also listed here.	Suggest a
		NAR Colla

Education (67)

Links to information about the techniques, Links to tools for predicting the expression, materials, people, places, and events of the greater alternative splicing, and regulation of a gene bioinformatics community. Included are current sequence are found here. This section also

Expression (382)

RSS Feed

Support

d



Genetics and population analysis

Gene expression

http://www.oxfordjournals.org/our_journals/bioinformatics/for_authors/scope_guidelines.html

Topics computer exercises



Query public databases Sequence format conversion Automatic RNA — protein translation Primer design Gene finding

Tomorrow 13:30-15:00

https://bioinformatics.amc.nl/