

#### POZNAŃ -> MONTREAL ->BETHESDA -> STATE COLLEGE -> MÜNSTER



2

#### TOPICS TO BE COVERED IN THIS COURSE

1

- Introduction to bioinformatics from the evolutionary perspective. [WM]
- Next Generation Sequencing. [WM]
- Sequence alignment and similarity search. [WM]
- Sequence assembly and gene prediction. [WM]
- Principles of heredity. Mutations, substitutions and polymorphisms. [CA]
- Distances and models. Synonymous and non-synonymous substitutions. Basics of the neutral theory. [CA]

3

• Phylogenetic inference. [CA]

# HANDS ON COMPUTER LAB

#### Computer Lab B, Schlossplatz 2b

- Nucleotide sequence analyses [November 10 12]
- Phylogenetic inference [December 8 10]
- Registration after November 1st at

http://bioinformatics.uni-muenster.de/teaching/register.pl



4

## CONTACT

- Prof. Claudia Acquisti <u>claudia.acquisti@uni-muenster.de</u>
- Prof. Wojciech Makałowski wojmak@uni-muenster.de
- Norbert Grundmann <u>ngrundma@uni-muenster.de</u> (lab coordinator)
- http://www.bioinformatics.uni-muenster.de/teaching/ courses-2013/bioinf1/index.hbi
- office hours see the web site







#### GROWTH OF BIOMEDICAL INFORMATION - GENBANK



8



9

#### IMPROVING TECHNOLOGY



10





# IMPROVING TECHNOLOGY



12 bioinfo1\_1\_2015 - October 22, 2015

# GETTING SEQUENCES



13

#### **READING** $\neq$ **UNDERSTANDING**

Carmina qui quondam studio florente peregi, flebilis heu maestos cogor inire modos.

Ecce mihi lacerae dictant scribenda Camenae et ueris elegi fletibus ora rigant.

14

### **READING** $\neq$ **UNDERSTANDING**

We shall best understand the probable course of natural selection by taking the case of a country undergoing some physical change. If the country were open were open on its borders, new forms would certainly immigrate, and this also would bla, bla bla become extinct inhabitants.

Charles Darwin - The Origin of Species

#### **READING** $\neq$ **UNDERSTANDING**

We shall best understand the probable course of by taking the case of a country undergoing some physical change. If the country were open were open on its borders, new forms would certainly immigrate, and this also would bla, bla bla become extinct inhabitants.

Charles Darwin - The Origin of Species

Boethius, Consolatio Philosot

15

#### 16

# CHALLENGE: HOW FROM THIS...







- & very smart
- Slow
- & error prone
- doesn't like repetitive tasks
- not so smart (stupid)
- & extremely fast
- & very accurate
- doesn't understand human languages;



# <text>

20

# EXAMPLE TASK: PUT SHOES ON!



A human just understands an order and often executes it automatically even without thinking

A computer needs detailed instruction (an algorithm)



#### PUT SHOES ON! INSTRUCTION FOR A COMPUTER

- 1. Find two the same shoes
- 2. Check if you have left and right shoe
- 3. Check if they are of the same size
- 4. Check if this is the right size
- 5. Put the left shoe on
- 6. Put the right shoe on
- 7. Tie the laces



22

# THE ORIGIN OF THE FIELD

21



Paulien Hogeweg coined the term *bioinformatica* to define "the study of informatic processes in biotic systems". <sub>Hege R, Hegereg P (1970) Bioinformatice cer</sub>

rkconcept. Kameleon 1(6): 28–29. (In Dutch.) Leiden: Leidse Biologen Club.

... but its origin can be tracked back many decades earlier.



# BIOINFORMATICS EMERGED AS AN INTERSECTION BETWEEN DIFFERENT DISCIPLINES

Information

technology

#### 24 bioinfo1\_1\_2015 - October 22, 2015

Molecular

evolution

# BIOINFORMATICS -DEFINITION

- Research, development, or application of computational tools and approaches for expanding the use of biological data, including those to acquire, store, organize, archive, analyze, or visualize such data.
- Its goal is to enable biological discovery based on existing information or in other words transform biological data into information and eventually into knowledge.

25

#### ROLE OF BIOINFORMATICS IN MODERN BIOLOGY

- molecular biology
- molecular evolution
- genomics
- system biology
- protein engineering
- drug design
- personalized medicine
- biogeography



26





27

28





30 bioinfo1\_1\_2015 - October 22, 2015



# **BIOLOGICAL DATABASES**

- organized sets of large amount of data, usually coupled with a software that enables data search, information extraction, and data update
- databases should be characterized by
  - easy data access
  - the possibility to extract only the information that is desirable

32

# INFORMATION IN DATABASES

- Databases and resources may contain many different kinds of information. Each item of entry is typically called an entry. Regardless of the type of resource, each entry comprises two main parts, each broken into one or more fields
- Descriptive information Annotation
  - Description
  - Literature references
- The raw data sequence or observations
- The most valuable information is frequently the annotation with the raw data providing a scaffold to organize this curated information.

33

## HISTORICAL (?) LOOK AT DATABASES

- Early systems were file based
- One entry one file
- Lookup based on computer system functions such as grep
- Drawbacks to file-based systems
  - Concurrency
  - No way to check consistency
    - Are values appropriate for fields?
    - Have you updated all necessary information?
  - Unable to limit queries to specific fields
  - Queries and especially updates may be slow and require special programming skills

34

#### GENBANK RECORD

LOCUS	AF062069 3808 bp mRNA INV 02-MAR-2000
DEFINITION	Limulus polyphemus myosin III mRNA, complete cds.
ACCESSION	AF062069
VERSION	AF062069.2 GI:7144484
KEYWORDS	
SOURCE	Atlantic horseshoe crab.
ORGANISM	Limulus polyphemus
	Eukaryota; Metazoa; Arthropoda; Chelicerata; Merostomata;
	Xiphosura; Limulidae; Limulus.
REFERENCE	1 (bases 1 to 3808)
AUTHORS	Battelle,BA., Andrews,A.W., Calman,B.G., Sellers,J.R.,
	Greenberg,R.M. and Smith,W.C.
TITLE	A myosin III from Limulus eyes is a clock-regulated phosphoprotein
JOURNAL	J. Neurosci. (1998) In press
REFERENCE	2 (bases 1 to 3808)
AUTHORS	Battelle,BA., Andrews,A.W., Calman,B.G., Sellers,J.R.,
	Greenberg, R.M. and Smith, W.C.
TITLE	Direct Submission
JOURNAL	Submitted (29-APR-1998) Whitney Laboratory, University of Florida,
	9505 Ocean Shore Bivd., St. Augustine, FL 32086, USA
REFERENCE	3 (Dases I to 3808)
AUTHORS	Battelle, BA., Andrews, A.W., Calman, B.G., Sellers, J.K.,
	Greenberg, R.M. and Smith, W.C.
TITLE	Direct Submission
JOORNAL	Submitted (02-MAR-2000) Whitney Laboratory, University of Fiorida,
DEMADY	Sourcean Shore Bivd., St. Augustine, FL 52086, USA
COMPANY	On Mar 2, 2000 this sequence version replaced si 2122700
COMPANY	on mar 2, 2000 unto sequence version replaced g1:3132/00.

# GENBANK RECORD

FEATURES Location/Qualifiers			
source	13808		
	/organism="Limulus polyphemus"		
	/db xref="taxon:6850"		
	/tissue type="lateral eye"		
CDS	2583302		
	/note="N-terminal protein kinase domain; C-terminal myosin		
	heavy chain head; substrate for PKA"		
	/codon start=1		
	/product="myosin III"		
	/protein id="AAC16332.2"		
	/db xref="GI:7144485"		
	/translation="MEYKCISEHLPFETLPDPGDRFEVQELVGTGTYATVYSAIDKQA		
	NKKVALKIIGHIAENLLDIETEYRIYKAVNGIQFFPEFRGAFFKRGERESDNEVWLGI		
	EFLEEGTAADLLATHRRFGIHLKEDLIALIIKEVVRAVQYLHENSIIHRDIRAANIMF		
	SKEGYVKLIDFGLSASVKNTNGKAOSSVGSPYWMAPEVISCDCLOEPYNYTCDVWSIG		
	ITAIELADTVPSLSDIHALRAMFRINRNPPPSVKRETRWSETLKDFISECLVKNPEYR		
	PCIQEIPQHPFLAQVEGKEDQLRSELVDILKKNPGEKLRNKPYNVTFKNGHLKTISGQ		
BASE COUNT	1201 a 689 c 782 g 1136 t		
ORIGIN			
1 tcgac	atotg tggtcgcttt ttttagtaat aaaaaattgt attatgacgt cctatctgtt		
3781 aaga	tacagt aactagggaa aaaaaaaa		
11			

# MODERN RESOURCES

- Relational Database Management Systems (RDBMS)
  - Introduced in the 1970s
  - Off-the-shelf software (commercial and open source)
    - Oracle, DB2, MySQL
  - High level declarative language SQL
  - Concurrency
  - Transaction control
  - Consistency



37

# <section-header><complex-block>

38

## CRITICAL ISSUES FOR BIOLOGICAL DATABASES

- Annotation
  - Correctness
  - Consistency
- Quality
- Archival Quality
- Updates
  - Raw data
  - Annotation



# ANNOTATION

CRITICAL ISSUES

- Correctness many genes are annotated primarily based on sequence comparisons. Annotation is copied from a similar sequence to a novel sequence. This may cause some problems
  - Comparison may have been done when the data were less complete
  - If sequence is incorrectly annotated, this error propagates through the database

40

## CRITICAL ISSUES ANNOTATION QUALITY

39

- Who supplies the annotation? An expert, or a non-expert at the database
- Many databases have defined groups of "experts" to help annotated genes or gene families, but there is no peer-review of information in databases
- What is the vocabulary?



#### CRITICAL ISSUES ARCHIVAL QUALITY

- Databases have been torn between trying to be archival – to simply report information as experts publish it (*primary databases*), or curated – to provide the best editorially reviewed data on a topic (*secondary DB*).
- Can the same entry be recovered later?
  - Accession numbers are more stable than entry or locus names
  - Many databases do not note that there have been changes to the data! What you retrieve today may be different than yesterday

## CRITICAL ISSUES UPDATES

- How often are updates done? Major databases take direct submissions.
- Generally, only the original submitter can change an entry, even if you can prove it is wrong. This is tied to the question of archival versus curated.
- How is annotation updated as more knowledge is available? Who decides?



#### SECONDARY (SPECIALIZED) DATABASES

- Boom of biological databases
- Every year first issue of *Nucleic Acids Research* dedicated to biological databases
  - http://nar.oxfordjournals.org/content/43/D1.toc
  - this year's database issue includes 1509 databases
  - the first collection published in 1993 contained description of 24 databases

44

43



#### 45

## EVOLUTIONARY BASIS OF BIOINFORMATICS



46



#### HOMOLOGS AT THE MOLECULAR LEVEL

COW	ATGACTARCATTCGRAAGTCCCACCACTARTARARATTGTARAC
sheep	ATGATCAACATCCGAAAAAACCCACCCACTAATAAAAATTGTAAAC
goat	ATGACCARCATCCGAAAGACCCACCATTAATAAAAATTGTAAAC
horse	ATGACAAACATCCGGAAATCTCACCCACTAATTAAAATCATCAAT
donkey	ATGACAAACATCCGAAAATCCCACCCGCTAATTAAAATCATCAAT
ostrich	ATGGCCCCCAACATTCGAAAATCGCACCCCCTGCTCAAAATTATCAAC
emu	ATGGCCCCTAACATCCGAAAATCCCACCCTCTACTCAAAATCATCAAC
turkey	ATGGCACCCARTATCCGARAATCACACCCCCTATTARAAACAATCAAC

Two sequences that share common ancestry. Significant sequence similarity usually suggests homology, however sequence similarity may occur also by chance and some homologous sequences may diverge beyond detectable similarity.

#### HOMOLOGS: ORTHOLOGS AND PARALOGS

**ORTHOLOGS.** Genes or sequences that result from a speciation event followed by a sequence divergence. Such genes may not exist side by side in the same genome. The last common ancestor of two orthologous sequences existed just before speciation event.



#### HOMOLOGS: ORTHOLOGS AND PARALOGS

PARALOGS. Genes or sequences that resulted from duplication of genetic material followed by a sequence divergence. Such genes may descend and diverge while existing side by side in the same genome. If speciation occurs after gene duplication, then two paralogous genes may exist in two different genomes. The last common ancestor of two paralogous sequences existed just before duplication event.



50

### EVOLUTIONARY BASIS OF BIOINFORMATICS

49



#### HOMOLOGS: ORTHOLOGS AND PARALOGS

Compared Genes	Relation	Time of last comm. ancestor	Evolutionary event at the time of last common ancestor	Presence i the same species
A - B	paralogy	t <sub>1</sub>	gene duplication	yes
A1 - A2	orthology	t2	speciation	no
A1 - B1	paralogy	t1	gene duplication	yes
A1 - B2	paralogy	t <sub>1</sub>	gene duplication	no
A1 - B3	paralogy	tı	gene duplication	no
A2 - A1	orthology	12	speciation	no
A2 - B1	paralogy	t <sub>1</sub>	gene duplication	no
A2 - B2	paralogy	t <sub>1</sub>	gene duplication	yes
A2 - B3	paralogy	t <sub>1</sub>	gene duplication	yes
B1 - A1	paralogy	t <sub>1</sub>	gene duplication	yes
B1 - A2	paralogy	t <sub>1</sub>	gene duplication	no
B1 - B2	orthology	12	speciation	no
B1 - B3	orthology	12	speciation	no
B2 - A1	paralogy	t <sub>1</sub>	gene duplication	no
B2 - A2	paralogy	t <sub>1</sub>	gene duplication	yes
B2 - B1	orthology	t2	speciation	no
B2 - B3	paralogy	13	gene duplication	yes
B3 - A1	paralogy	t <sub>1</sub>	gene duplication	yes
B3 - A2	paralogy	tı	gene duplication	no
B3 - B1	orthology	12	speciation	no
B3 - B2	paralogy	ts .	gene duplication	yes









#### **COMPARATIVE GENOMICS**





15 000 victims of thalidomide

What is true for mouse is not necessarily true for human...

56



#### DID THE FLORIDA DENTIST INFECT HIS PATIENTS WITH HIV?



57



## THE MYSTERY OF THE CHILEAN BLOB

#### >Chilean\_Blob





# THE MYSTERY OF THE CHILEAN BLOB

> emb AJ277029.2  D Physeter macrocephalus mitochondrial genome Length=16428						
<pre>Score = 1074 bits (581), Expect = 0.0 Identitics = 585/587 (99%), Gaps = 0/587 (0%) Strand=Pius/Pius</pre>						
Query Sbjct	1 4400	TANTACTANCTATATCCCTACTCCCATTCTCATCGGGGGTTGAGGAGGACTANACCAGA	60 4459			
Query Sbjct	61 4460	CTCAACTCCGAAAAATTATAGCTTACTCATCAATCGCCCACATAGGATGAATAACCACAA	120 4519			
Query	121	TCCTRCCCTRCRATRACRACCTATARACCCTACTARATCTRTGTCRCRATRACCT	180			
Query	181	TCACCATATTCATACTATTTATCCAAAACTCAACCACACCAC	240			
Query	241	TURCURTATION AND ANTITATO ANALYCANAL TO ACCARACIAL CARACIAN CONTROL CARACI	300			
Sbjot Query	4640 301	CATGAAACAAAACACCCATTACCACAACCCTTACCATACTTACCATACTTACCATAGGGG GCCTCCCACCACTTCGGGCTTTATCCCCAAATGAATAATTATTCAAGAACTAACAAAAA	4699 360			
Sbjct Query	4700 361	OCCYCCCACCACTCTCGOGCTTTATCCCCCAAATGAATAATTATTCAAGAACTAACAAAAA ACGAAACCCTCATCATACCAACCTTCATAGCCACCACCACCACCACCTCCTCCTCCT	4759 420			
Sbjet	4760	ACGAAGCCCTCATCATACCAACCTTCATAGCCACCACAGCATTACTCAACCTCTACTTCT	4819			
Sbjct	4820	ATATACGCCTCACCTACTCAACAGCACTAACCCTATTCCCCTCCCACAAATAACATAAAAAA	4879			
Query	481	TAAAATGACAATTCTACCCCCACAAAACGAATAACCCTCCTGCCAACAGCAATTGTAATAT	540			
Query	541	CAACAATACTCCTACCCCTTACACCAATACTCTCCCACCCTATTAT	4555			
Sbjct	4940	CAACAATACTCCTACCCCTTACACCAATACTCTCCACCCTATTAT				



62

# **BIOINFORMATICS CREED**

- Remember about biology
- Do not trust the data
- Use comparative approach
- Use statistics
- Know the limits
- Remember about biology!!!

