

Information Systems in the Life Sciences (ISiLS)

Fons J. Verbeek
Imaging & BioInformatics, LIACS

ISiLS #1, FJV

1

Lectures & Supervision

- Dr. Erwin Bakker
- Dr. Nies Huijsmans
- Dr. Fons Verbeek (coordination)

ISiLS #1, FJV

2

Unraveling the complexity of live,
the ultimate challenge

ISiLS #1, FJV

3

Contents

- Introduction
- Organization of the Course
- Data
- Informatics
- Processes
- Contents
- Questions

ISiLS #1, FJV

4

Organization (1)

- Seminarium
 - Limited number of participants
 - Taking the course is participating in the course
 - Attending the course
- Introductory Lectures (2 series)
 - Bakker
 - Huijsmans
 - Verbeek

ISiLS #1, FJV

5

Organization (2)

- Paper presentations
 - Schedule
 - Groups, depending on # people attending
 - Deadline
- Paper writing
 - End of the course
 - Discuss with course administration
 - Subjects equally divided over participants
 - Deadline

ISiLS #1, FJV

6

Information Systems in LS

- Fons Verbeek
 - From signals to systems
 - Virtual Cell
 - Virtual organism
- Nies Huijsmans
 - Molecular modeling, shape
 - Content bases search, image
- Erwin Bakker
 - Data integration MicroArray (DIAL)
 - Gene Browsers
- *Other*

ISLS #1, FJV

7

Information Systems in LS

- Closely related to BioInformatics
 - What is BioInformatics?
- Information Retrieval
 - How is the information structured
 - Sequences
 - Graphics
 - Images
 - Other ...
 - Learn basic components

ISLS #1, FJV

8

BioInformatics

BIO-Informatics ?

$$(2b) \vee \neg(2b)$$

ISLS #1, FJV

9

Bio-Informatics: $(2b) \vee \neg(2b)$

- In the early 1980 Super-Computer scientists (Hwa Lim) realized the potential of combining of biology and computer science:
CompBio (“that is not a word ...”)
- More whimsical: **Bioinformatique**
- This changed to: Bio-informatics
- Using - or / was troublesome: Bioinformatics

ISLS #1, FJV

10

Interpretations of BioInformatics

- Political interpretation: definition is under debate.
- Narrow interpretation: the information science techniques needed to support genome analysis.
- Broader interpretation: synonymous with computational biology or computational molecular biology.

ISLS #1, FJV

11

Definitions of BioInformatics (1)

Bioinformatics:

Research, development or application of computation tools and approaches for expanding the use of biological, medical behavioral or health data, including those to acquire, store, organize, archive, analyze or visualize such data.

after:

NIH Biomedical information Science & Technology Initiative Consortium

ISLS #1, FJV

12

Definitions of BioInformatics (2)

Computational Biology:

The development and application of data-analytical and theoretical methods, mathematical modeling and computational simulation techniques to the study of biological, behavioral and social systems.

after:

NIH Biomedical information Science & Technology Initiative Consortium

ISLS #1, FJV

13

Definitions of BioInformatics (3)

BioInformatics:

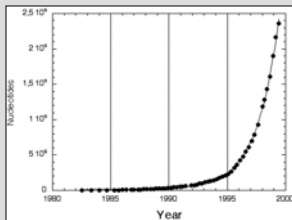
- Integration of mathematical, statistical and computer methods to analyze biological, biochemical and biophysical data.
- The science of developing computer **databases** and **algorithms** for the purpose of speeding up biological research (Human Genome Project)

ISLS #1, FJV

14

Why the Need? Massive Data Explosion

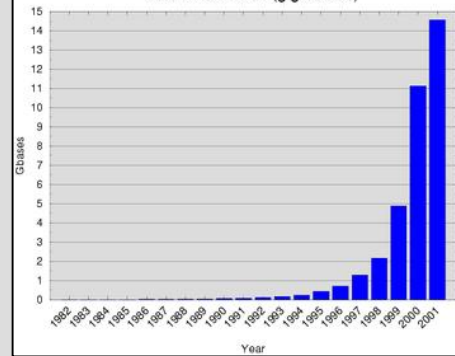
- Growth of EMBL nucleotide database
- And this is just the Nucleotide database.
- What is a nucleotide



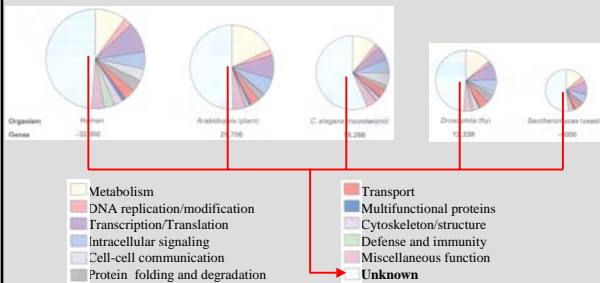
ISLS #1, FJV

15

EMBL Database Growth
total nucleotides (gigabases)



Necessity of BioInformatics



Adapted from human genome consortium: Nature 2001

ISLS #1, FJV

17

Data in Information Systems

ISLS #1, FJV

18

Molecular Basis of Living Systems

- A **gene** is a unit of information within the **chromosome** that can be inherited
- Expression of genomic information involves a complex sequence of steps
 - DNA to mRNA (transcription)
 - mRNA to proteins (translation)
- DNA contains only “alphabets”
 - Nucleotides
- The next significant step will come from deep understanding of protein expressions and interactions
 - *Functional genomics*

ISLS #1, FJV

19

101

- Information systems & Molecular Biology
 - Requires knowledge of databases
 - Data communication: Internet
 - Data type, i.e. Molecular biology
- Understanding basic Principles of
 - Molecular biology
 - BioChemistry
 - Molecular genetics
 - 101 will be provided (pdf)

ISLS #1, FJV

20

Key Components

- Cell
- Cell Nucleus
- Chromosome

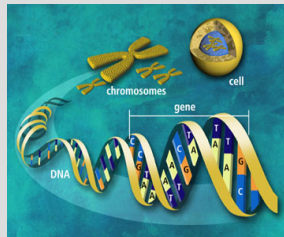


Photo from DCC/Biomed Center/Photic/CORBIS

- Gene:
 - smallest physical unit of heredity coding: information carrier of a feature
- Let us further decompose a gene ...

ISLS #1, FJV

21

Building Blocks: Nucleotides

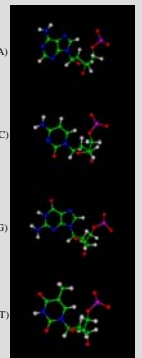
- A **nucleotide** is the building block of DNA and RNA
 - Nucleotide = bases + sugars + phosphate
- Bases:
 - Adenine
 - Cytosine
 - Guanine
 - Thymine, replaced by Uracil in RNA
- Complementary pairs
 - A complements with T
 - C complements with G

Adenine (A)

Cytosine (C)

Guanine (G)

Thymine (T)



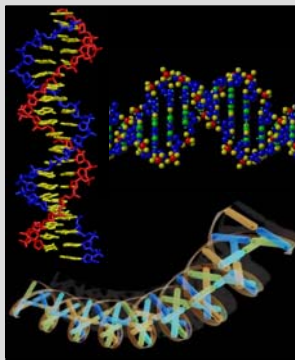
Adenine from NYU; Cytosine, B; Guanine, R; Thymine, W; Hydrogen, P; Phosphate, R

ISLS #1, FJV

22

Nucleic Acids

- DNA and RNA
- In eukaryotes, DNA most commonly occurs as a double helix
 - Sugar-phosphate backbone on outside
 - Base paired by hydrogen bonds stacked on the inside
- DNAs are highly stable
 - Dipole-dipole interactions
 - Hydrophobic
- Complementary chains
 - The sequence of bases in one chain determines the sequence in the other chain

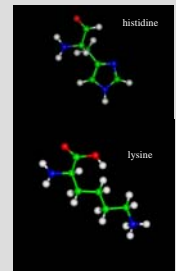


ISLS #1, FJV

23

Amino Acids

- Translation
 - Bases of mRNA are translated in groups of 3 (codons)
 - A codon translates into amino acid
 - CAU, CAC -> His
 - AAA, AAG -> Lys
 - Etc.
 - Specific codons for each of the amino-acids
- Proteins are chains of amino acids
 - 20 amino acids
 - 20 letter alphabet

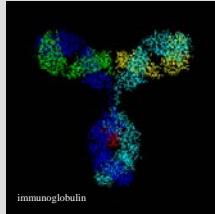


ISLS #1, FJV

24

Proteins

- Proteins are chains of amino acids
 - 20 amino acids,
 - 20 letter alphabet
- Variety of functions
 - Enzymes
 - Membrane receptors
 - Transport (e.g., hemoglobin)
 - Structure (e.g., collagen)
 - Nutrition (e.g., ovalbumin)
 - Immunity (e.g., antibodies)
 - Regulatory



ISLS #1, FJV

25

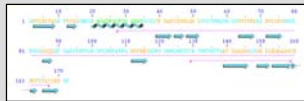
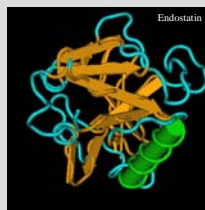
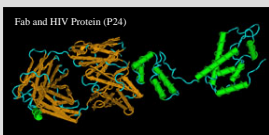
Protein Structure

- Complex structures
 - Primary – sequence
 - Secondary – α -helix and β -sheets
 - Tertiary – folding structures
 - Quaternary – multi-chain (multimeric) arrangements
- Protein structure determines function
 - Where is the active site(s)
 - What is the catalytic strength
 - Interaction in complexes

ISLS #1, FJV

26

Examples Protein Structure



ISLS #1, FJV

27

Key Principles Molecular Biology

- DNA acts as a template to replicate itself
- DNA is transcribed into RNA
- RNA is translated into a protein
- Sequence and structural homology (similarity) between molecules can be used to infer structural and functional similarity

ISLS #1, FJV

28

Fields of Application

- Genomic sequencing
- Comparative genomics
 - Comparisons to find similarities/differences
- Expression quantification
 - Relative abundance of expression during development or disease stages
- Functional genomics
 - Large-scale mapping of gene functions and associations
- Proteomics
 - Catalog of activities that characterize interactions among gene products
- Structural genomics
 - Protein structure mapping and predictions
- Research informatics and data management
 - Experimental data management

ISLS #1, FJV

29

Applications of Bioinformatics

- Sequence analysis, Alignment
- Comparison BLAST, FASTA
- Molecular modeling, Prediction modeling
- Databases for EST's, Sequences (HGP), Linkage Maps (Synteny), Physical Maps, Probes, Gene Array data.
- Databases of Gene Expression

ISLS #1, FJV

30

Bio Molecular Databases

- Used for 3 major tasks
 - Lookup
 - Is there a gene known for my protein?
 - Is there mutations known causing this disease?
 - Compare
 - Are there sequences available resembling my cloned protein?
 - Are these two sequences similar (to what extent)?
 - Predict
 - Can the active site residues of the this enzyme be predicted?
 - Can a 3D model of this protein be made?
- Answers are not necessarily found in “1” database
- Combined search, Integrate results
 - How can this be realized
 - Interoperable databases

ISLS #1, FJV

31

Core Databases

- Sequence search, BLAST:
 - Basic Local Alignment Search Tool
 - <http://www.ncbi.nlm.nih.gov>
- Protein structure: PHD
 - <http://www.embl-heidelberg.de/predictprotein/predictprotein>
- Molecular modeling and imaging: RasMol
 - <http://www.umass.edu/microbio/rasmol/>

ISLS #1, FJV

32

Core Databases

- Data repositories
 - GenBank: NCBI Nucleotide database
 - Protein DataBank (PDB):
<http://www.rcsb.org/pdb>
 - Repository for processing and distribution of 3D biological macromolecular structure data
- KEGG
 - Kyoto Encyclopedia of Genes and Genomes
 - From Genes to BioChemical Pathways

ISLS #1, FJV

33

Searching in Databases

- Key issue is the different information that is available in the different databases
- Added value is obtained if these databases are transparently accessible
- Information is shared!
- Learning the specific contents of a database
- Ontology's

ISLS #1, FJV

34

Searching in Sequences

- Complications
 - Sequence DBs contain enormous amounts of nucleotides
 - Query sequence is not exact
 - It is important to find non-exact matches (homologues)
- Techniques
 - Sequence alignments
 - Multiple sequence alignments
 - Sequences of common structure or function

ISLS #1, FJV

35

Sequence Alignment

Drosophila “eyeless” (S) gene vs human aniridia (Q)

```
pir||A1644 homeotic protein aniridia - human
Length = 447

Score = 256 bits (647), Expect = 5e-67
Identities = 128/146 (87%), Positives = 134/146 (91%), Gaps = 1/146 (0%)

Query: 24 IERLPSLEEMRNGHSGVWQLGGVFNQGRPLPDSTRQKIVELAHSGARPCDISRILQVSN 83
I R P+ M = HSGVWQLGGVFN GRPLPDSTRQKIVELAHSGARPCDISRILQVSN
Sbjct: 17 IFRPPARASMQNS-HSGVWQLGGVFNQGRPLPDSTRQKIVELAHSGARPCDISRILQVSN 75

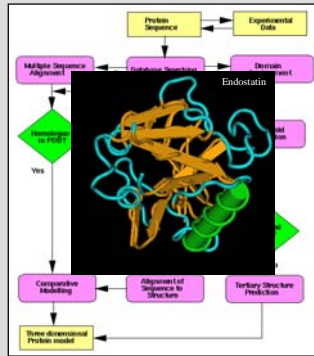
Query: 84 GCVSKILGRVYETGSIKPRALGSGKPRVATAEVSKISQYKRCPCPSIFAMEIRDRLQEN 143
GCVSKILGRVYETGSIKPRALGSGKPRVATAEVSKISQYKRCPCPSIFAMEIRDRL E
Sbjct: 76 GCVSKILGRVYETGSIKPRALGSGKPRVATAEVSKISQYKRCPCPSIFAMEIRDRLSEG 135

Query: 144 VCTNINI PSVSSINRVLRLNLAQKQ 169
VCTNINI PSVSSINRVLRLNLA++K+Q
Sbjct: 136 VCTNINI PSVSSINRVLRLNLA SEQ 161
```

ISLS #1, FJV

36

Structure Prediction, Drug design

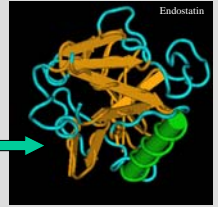
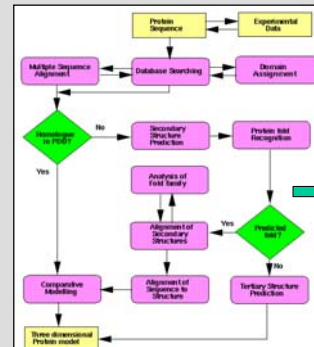


Adapted from EMBL

ISLS #1, FJV

37

Structure Prediction, Drug design



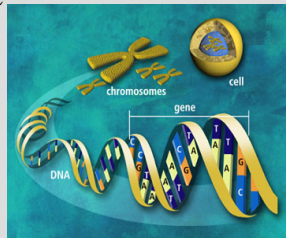
Adapted from EMBL

ISLS #1, FJV

38

Key definitions & relations

- Cell % Cell Nucleus
- Cell Nucleus % Chromosome
- Chromosome % Gene
- Gene % DNA
- % = has part
- Or reverse: is part of
- ontology



ISLS #1, FJV

39

Facts on Ontology

- Share common understanding of a domain
- Make domain knowledge explicit
- There is no explicit method of writing an ontology
 - Depending on the application in mind
 - Obtained through iteration
- Concepts
 - Objects & Relationships in domain of interest
 - Nouns & Verbs in domain to be described
- Biology (Life Sciences)
 - GO-BO initiative coordinated by EBI

ISLS #1, FJV

40

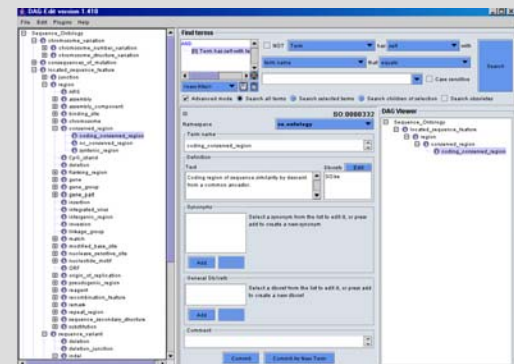
Open Biological Ontologies (OBO)

- Produced with DAG edit
 - Directed Acyclic Graph
- Ontologies stored in MySQL database
- Regular updates of the ontologies submitted to the database
 - Sequence ontology
 - Microarray Gene Expression Data (MGED)
 - Generic Model Organism databases

ISLS #1, FJV

41

DAG Edit



ISLS #1, FJV

42

DNA-chip

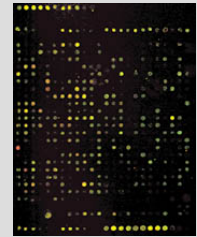
- Glass chip consisting of array of spots,
 - each spot 20-100 μm diameter
- Each spot contains a RNA of interest – “probe”
- Fluorescently tagged mRNA samples flow over probes
- Two-color fluorescence can be used to identify:
 - normal from abnormal,
 - over-expression from under-expression, etc.

ISLS #1, FJV

43

MicroArray

- DNA chip a.k.a.
 - = MicroArray
 - = DNA Array
- Miniaturization
- Temporal-Information
- Little Space Information
- Lots of genes tested at the same time
- Renders a pattern of gene expression



ISLS #1, FJV

44

Applications of MicroArrays

- Genomics
 - Fundamental research
 - Systems biology
- Toxicogenomics
 - Field of functional genomics focusing on environmental health
 - Samples are taken from
 - environment or
 - from a food production process
- Food genomics

ISLS #1, FJV

45

Spatio-Temporal Frameworks

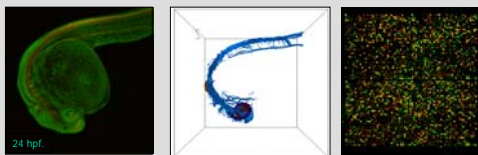
- Lots of data are generated which are not stored in one single repository
- Link repositories, create conditions to make that feasible
- Model system specific
 - Fast (short generation time)
 - Slow, mammalian, close to human genome (rat, mouse)
- Gene-expression can be applied on
 - Micro-arrays
 - In situ (in vivo), whole mount
 - Different model systems
- Relate gene expression to a model system
- Relate model systems

ISLS #1, FJV

46

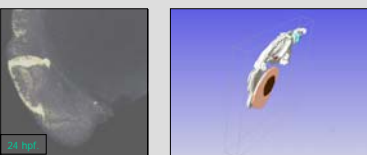
Patterns of Gene Expression

translation



24 hpf
topro anti-tubulin

transcription

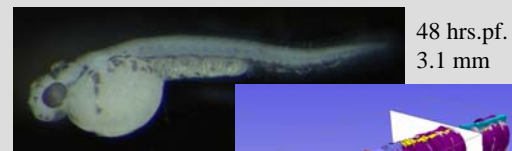


24 hpf
wnt1

ISLS #1, FJV

47

3D Atlas Zebrafish Development



48 hrs.pf.
3.1 mm

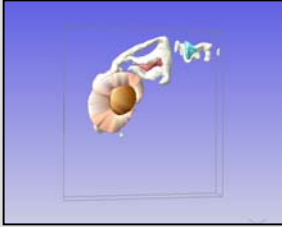
200 Mb image data,
< 2 Mb annotations

ISLS #1, FJV

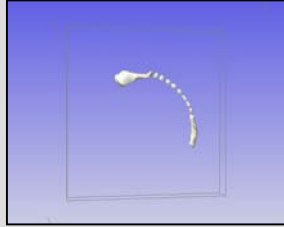
48

"FISH" Gene Expression Patterns

- CLSM images
- Imaging protocol for processing



wnt1 (24 hrs pf.)



myoD (24 hrs pf.)

ISLS #1, FJV

49

Integration of Information



ISLS #1, FJV

50

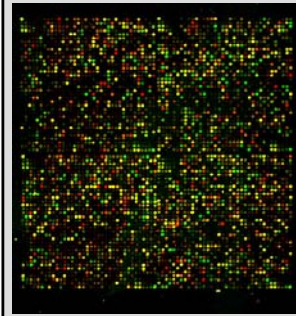
Finding new genes

- Genetic code is an alphabet (4 letters)
- Genome is a string "...ATTGCGTA ..." (very long)
- Looking for meaning in that string
 - Codon: group of 3 coding for an amino acid
- Scanning for Open Reading Frames (ORF)
 - Start codon
 - Stop codon
 - Complication: Intron (non-coding) and Exon (coding)
- Algorithmic approaches, all available data i.e.:
 - Molecular Biology,
 - Model Fitting to intron-exon boundaries,
 - Similarity to other organisms

ISLS #1, FJV

51

When is a gene expressed



- MicroArray (DNA Chip)
- Large collection of genes in 20-100 μm spots
- Two samples
- Different conditions
- Lots of data
 - Co-expression
 - Level of expression
 - Relations between genes
- Cluster analysis
 - Algorithmic approach
- Database
 - Storage, Retrieval
 - Integration with other data

ISLS #1, FJV

52

Where is a gene expressed

- MicroArrays (and other tools): temporal expression
- Microscopy of whole organism with *in situ* hybridisation: spatial expression (3D)
- Microscopy of whole organism with *immuno-histochemistry*: spatial expression (3D)
- Combine different images
 - Built a framework to look for expression
 - Learn about the genetic networks underpinning a function

ISLS #1, FJV

53

Systems Biology

- Molecular components of a system
- Understand interactions at system level
- Integration of different data
- Integration of different disciplines



ISLS #1, FJV

54

Data Deluge

ISLS #1, FJV

55

Solutions to Data Deluge

- GRID
 - Bringing data to computer power, Sharing data
 - High energy physics (CERN)
- E.g.: SETI
 - Searching for Extra-Terrestrial Intelligence
 - Distributing data to grid of computers
 - Berkeley University CA, USA
- E.g. compound screening (chemistry)
 - Looking for right molecule as cancer drug
 - Lifesaver project, Oxford university, UK
- E.g. VL-e
 - Using the GRID for a Virtual Lab (environment)

ISLS #1, FJV

56

Sharing data & Computing power



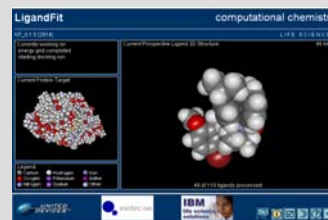
- Search for Extraterrestrial Intelligence
- program that downloads and analyzes radio telescope data

ISLS #1, FJV

57

Screensaver

- Oxford sends molecules and protein targets to UD.com
- UD.com Global MetaProcessor Grid sends tasks to members
- Screensaver processes the tasks and returns the results to central UD.com server
- Computers are typically idle 90% of the time.



ISLS #1, FJV

58

LifeSaver: Computational Chemistry

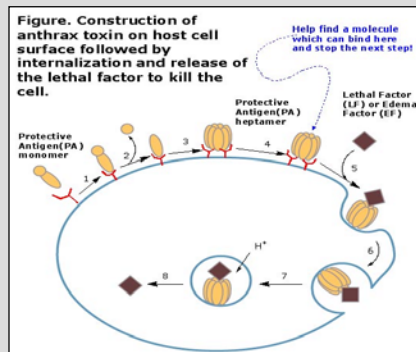
Two approaches used in the project

- THINK – Keith Davis
 - Cancer and Anthrax targets
- LigandFit – Accelrys Inc.
 - Cancer and Smallpox targets

ISLS #1, FJV

59

LifeSaver & Anthrax

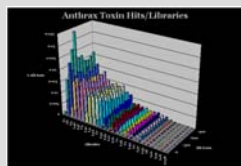
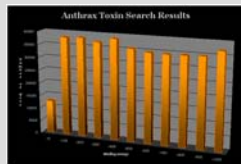


Taken from K. Harrison Oxford University
ISLS #1, FJV

60

Results THINK: Anthrax

- Four Weeks
- 376,064 molecules as hits from the 3.5 billion screened.



ISiLS #1, FJV

61

The Questions in ISiLS

ISiLS #1, FJV

62

The Questions

- What are good information systems
- Why is it a good information system (HCI)
 - Success (measurable)
 - tools offered
- How to improve the information system
- What extra tools & techniques are
 - Required
 - To be developed

ISiLS #1, FJV

63

Summary

- Typical issues in BioTech Information systems
- BioInformatics
- Molecular Biology primer
- Databases & Frameworks
- Examples

- Case studies are worked out in this course

ISiLS #1, FJV

64