

[www.bio.ind.in](http://www.bio.ind.in)



## **INDIAN INSTITUTE OF BIOINFORMATICS**

**New Delhi**

This website titled [www.bio.ind.in](http://www.bio.ind.in) is the official website of the Indian Institute of Bioinformatics (IIB) inaugurated on the auspicious occasion of the World Environment Day on 5<sup>th</sup> June 1991 by the Prime Minister of India Hon'ble Shri Chandra Shekhar. The activities of the IIB include training, research, publications, conference organizations and consultancy in the areas of bioprogramming, genetic engineering, molecular biology, biochemistry, biotechnology, microbiology, structural bioinformatics, biocomputing languages, bioinformatics software, algorithmic bioinformatics etc.

During the last 28 years of its existence, the IIB has organised many seminars, symposia, conventions, congresses and summits on different subjects relating to bioinformatics. The IIB has designed the following Certificate Courses for the benefit of all those interested in acquiring expert knowledge of bioinformatics and related subjects :

**Certificate in Bioinformatics**

**Certificate in Bioprogramming Languages**

**Certificate in Genetics**

**Certificate in Statistical Methods and DBMS**

**Certificate in Molecular Biology and Bio Chemistry**

**Certificate in Biotechnology and Industrial Microbiology**

**Certificate in Structural Bioinformatics**

**Certificate in Biocomputing Languages**

**Certificate in Bioinformatics Software**

## **Certificate in Algorithmic Bioinformatics**

## **Certificate in Computational Biology**

**Duration :** Three Months

**Eligibility :** No Minimum Educational Qualification has been led down. All those interested in acquiring expert knowledge of Bioinformatics are eligible to apply.

**Fee :** Rs. 3500 or US\$ 85 only to be paid on account of admission, registration and evaluation fee.

This amount is to be transferred to our Bank Account having the following details :

**Name of the Account : Indian Institute of Bioinformatics**

**Bank Name : Indian Bank, Saket Branch, New Delhi, India**

**Account Number : 6755880232**

**IFS Code : IDIB000S097**

Step by step method of learning at the Indian Institute of Bioinformatics :

1. Get the Admission Form downloaded and complete the same
2. Email the filled-up Admission Form
3. Pay the Admission Fee by Cheque / Draft / Electronic Transfer
4. Receive the Roll Number and Study Materials
5. Go through the e-book carefully
6. Complete the assignments and send the same to the Institute by Email / Post
7. Submit the Project Report based on your experience and knowledge acquired regarding any topic relevant to the admitted student.
8. Wait for the announcement of results.
9. Receive the Certificate (Online)

In case of any clarification, contact the Facilitation Officer, Indian Institute of Bioinformatics, A 14-15-16, Paryavaran Complex, New Delhi – 110030, India by post or by Email : [bioinformatics@ecology.edu](mailto:bioinformatics@ecology.edu)

**For any clarification, contact may be made through telephone by calling on 011-29533801, 011-29533830, 011-29535053.**

**24-Hours Helpline : 9999833886**

Roll Number Allotted

Stamp Size Photo

.....



# INDIAN INSTITUTE OF BIOINFORMATICS

A 14-15-16, Paryavaran Complex, South of Saket, New Delhi-110030

Email : bioinformatics@ecology.edu Tel. : 011-29533801, 011-29533830

## ADMISSION FORM

**NAME OF THE COURSE SELECTED .....**

Name of the Candidate .....

Father's Name .....

Mother's Name .....

Date of Birth ..... Nationality.....

Address .....

.....

.....

Email .....Website (if any).....

Mobile.....Telephone.....

Educational Qualification .....

.....

.....

Mention how will this course help you ?

.....

.....

Details of Fee paid (Cheque / Draft / Electronic Transfer)

.....

.....

Date

Signature



*The Prime Minister of India Hon'ble Shri Chandra Shekhar inaugurating the Indian Institute of Bioinformatics (IIB) on the occasion of the World Environment Day on 5<sup>th</sup> June 1991.*

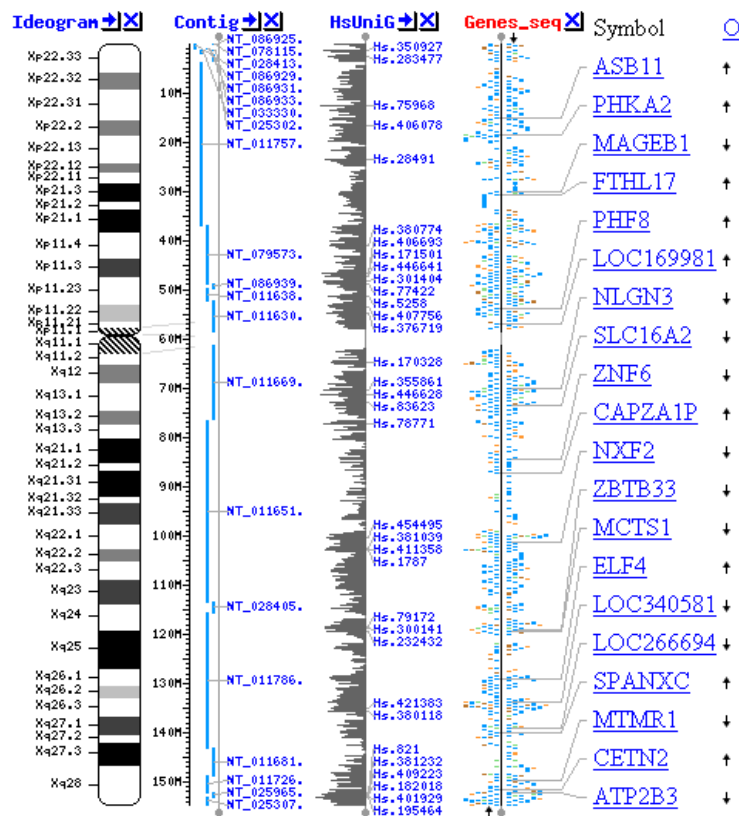
## **OTHER PUBLICATION**

The Indian Institute of Bioinformatics (IIB) has also come out with new publications for the benefit of scientists, researchers, industrial houses, pharmaceutical industries, bioinformatics organizations and companies, hospitals and medical laboratories. The following is the illustrative but not exhaustive list of topics :

1. Introduction to Bioinformatics
2. Sequence Analysis
3. Genome Annotation
4. Computational Evolutionary Biology
5. Comparative Genomics
6. Genetics of Diseases
7. Analysis of Mutations in Cancer
8. Gene and Protein Expression
9. Structural Bioinformatics
10. Prediction of Protein Structure
11. Network and Systems Biology
12. Molecular Interaction Networks

13. High-throughput Image Analysis
14. High-throughput Single Cell Data Analysis
15. Biodiversity Informatics
16. Web Services in Bioinformatics
17. Bioinformatics Workflow Management System
18. Flow Cytometry Bioinformatics
19. Computational Genomics
20. Health Informatics
21. Computational Bio-Modelling
22. Functional Genomics
23. Phylogenetics
24. Proteomics
25. Intelligent Bioinformatics

## BIOINFORMATICS



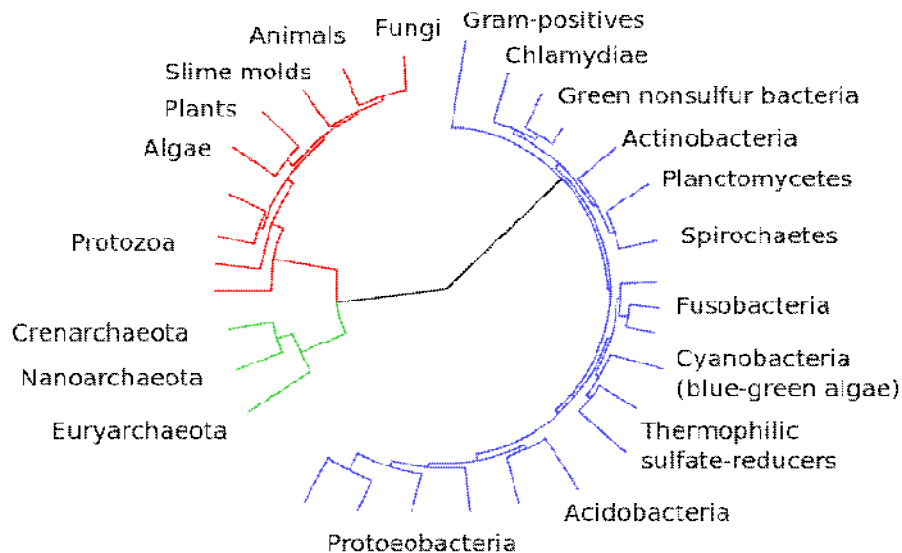
Map of the human X chromosome (from the NCBI website).

Assembly of the human genome is one of the greatest achievements of bioinformatics.

## EVOLUTIONARY BIOLOGY

Bioinformatics is an interdisciplinary scientific field that develops methods for storing, retrieving, organizing and analyzing biological data. A major activity in bioinformatics is to develop software tools to generate useful biological knowledge. Bioinformatics is a distinct science from biological computation, the latter being a computer science and computer engineering subfield using bioengineering and biology to build biological

computers, whereas bioinformatics simply uses computers to better understand biology. Bioinformatics is similar to computational biology and has similar aims to it but differs on scale: whereas bioinformatics works with basic biological data (e.g. DNA bases), i.e. it works on the small scale paying attention to details, computational biology is a subfield of computer science which builds large-scale general theoretical models of biological systems seeking to expand our understanding of them from an abstract point of view, just as mathematical biology does with mathematical models.



Bioinformatics uses many areas of computer science, statistics, mathematics and engineering to process biological data. Complex machines are used to read in biological data at a much faster rate than before and used in decoding the code of life. Databases and information systems are used to store and organize biological data. Analyzing biological data may involve algorithms in artificial intelligence, soft computing, data mining, image processing, and simulation. The algorithms in turn depend on theoretical foundations such as discrete mathematics, control theory, system theory, information theory, and statistics. Commonly used software tools and technologies in the field include Java, C#, XML, Perl, C, C++, Python, R, SQL, CUDA, MATLAB, and spreadsheet applications.

## INTRODUCTION

### HISTORY

Paulien Hogeweg coined the term "Bioinformatics" in 1970 to refer to the study of information processes in biotic systems. This definition placed bioinformatics as a field parallel to biophysics (the study of physical processes in biological systems) or biochemistry (the study of chemical processes in biological systems).

**Sequences** : Computers became essential in molecular biology when protein sequences became available after Frederick Sanger determined the sequence of insulin in the early 1950s. Comparing multiple sequences manually turned out to be impractical. A pioneer in the field was Margaret Oakley Dayhoff, who has been hailed by David

Lipman, director of the National Center for Biotechnology Information, as the "mother and father of bioinformatics." Dayhoff compiled one of the first protein sequence databases, initially published as books and pioneered methods of sequence alignment and molecular evolution. Another early contributor to bioinformatics was Elvin A. Kabat, who pioneered biological sequence analysis in 1970 with his comprehensive volumes of antibody sequences released with Tai Te Wu between 1980 and 1991.

**Genomes** : As whole genome sequences became available, again with the pioneering work of Frederick Sanger, the term bioinformatics was re-discovered to refer to the creation of databases such as GenBank in 1982. With the public availability of data tools for their analysis were quickly developed and described in journals such as Nucleic Acids Research which published specialized issues on bioinformatics tools as early as 1982.

**Goals** : In order to study how normal cellular activities are altered in different disease states, the biological data must be combined to form a comprehensive picture of these activities. Therefore, the field of bioinformatics has evolved such that the most pressing task now involves the analysis and interpretation of various types of data. This includes nucleotide and amino acid sequences, protein domains, and protein structures. The actual process of analyzing and interpreting data is referred to as computational biology. Important sub-disciplines within bioinformatics and computational biology include:

- the development and implementation of computer programs that enable efficient access to, use and management of, various types of information.
- the development of new algorithms (mathematical formulas) and statistical measures with which to assess relationships among members of large data sets. For example, there are methods to locate a gene within a sequence, to predict protein structure and/or function, and to cluster protein sequences into families of related sequences.

The primary goal of bioinformatics is to increase the understanding of biological processes. What sets it apart from other approaches, however, is its focus on developing and applying computationally intensive techniques to achieve this goal. Examples include: pattern recognition, data mining, machine learning algorithms, and visualization. Major research efforts in the field include sequence alignment, gene finding, genome assembly, drug design, drug discovery, protein structure alignment, protein structure prediction, prediction of gene expression and protein-protein interactions, genome-wide association studies, and the modeling of evolution.

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

Over the past few decades rapid developments in genomic and other molecular research technologies and developments in information technologies have combined to produce a tremendous amount of information related to molecular biology. Bioinformatics is the name given to these mathematical and computing approaches used to glean understanding of biological processes.

## **APPROACHES**

Common activities in bioinformatics include mapping and analyzing DNA and protein sequences, aligning DNA and protein sequences to compare them, and creating and viewing 3-D models of protein structures.

There are two fundamental ways of modelling a Biological system (e.g., living cell) both coming under Bioinformatic approaches.

- Static
  - Sequences – Proteins, Nucleic acids and Peptides
  - Interaction data among the above entities including microarray data and Networks of proteins, metabolites
- Dynamic
  - Structures – Proteins, Nucleic acids, Ligands (including metabolites and drugs) and Peptides (structures studied with bioinformatics tools are not considered static anymore and their dynamics is often the core of the structural studies)
  - Systems Biology comes under this category including reaction fluxes and variable concentrations of metabolites
  - Multi-Agent Based modelling approaches capturing cellular events such as signalling, transcription and reaction dynamics

A broad sub-category under bioinformatics is structural bioinformatics.

## **MAJOR RESEARCH AREAS**

Bioinformatics has become an important part of many areas of biology. In experimental molecular biology, bioinformatics techniques such as image and signal processing allow extraction of useful results from large amounts of raw data. In the field of genetics and genomics, it aids in sequencing and annotating genomes and their observed mutations. It plays a role in the textual mining of biological literature and the development of biological and gene ontologies to organize and query biological data. It plays a role in the analysis of gene and protein expression and regulation. Bioinformatics tools aid in the comparison of genetic and genomic data and more generally in the understanding of evolutionary aspects of molecular biology. At a more integrative level, it helps analyze and catalogue the biological pathways and networks that are an important part of systems biology. In structural biology, it aids in the simulation and modeling of DNA, RNA, and protein structures as well as molecular interactions.

## **SEQUENCE ANALYSIS**

Since the Phage  $\Phi$ -X174 was sequenced in 1977, the DNA sequences of thousands of organisms have been decoded and stored in databases. This sequence information is analyzed to determine genes that encode polypeptides (proteins), RNA genes, regulatory sequences, structural motifs, and repetitive sequences. A comparison of genes within a species or between different species can show similarities between protein functions, or relations between species (the use of molecular systematics to construct phylogenetic trees). With the growing amount of data, it long ago became impractical to analyze DNA sequences manually. Today, computer programs such as



BLAST are used daily to search sequences from more than 260 000 organisms, containing over 190 billion nucleotides. These programs can compensate for mutations (exchanged, deleted or inserted bases) in the DNA sequence, to identify sequences that are related, but not identical. A variant of this sequence alignment is used in the sequencing process itself. The so-called shotgun sequencing technique (which was used, for example, by The Institute for Genomic Research to sequence the first bacterial genome, *Haemophilus influenzae*) does not produce entire chromosomes. Instead it generates the sequences of many thousands of small DNA fragments (ranging from 35 to 900 nucleotides long, depending on the sequencing technology). The ends of these fragments overlap and, when aligned properly by a genome assembly program, can be used to reconstruct the complete genome. Shotgun sequencing yields sequence data quickly, but the task of assembling the fragments can be quite complicated for larger genomes. For a genome as large as the human genome, it may take many days of CPU time on large-memory, multiprocessor computers to assemble the fragments, and the resulting assembly will usually contain numerous gaps that have to be filled in later. Shotgun sequencing is the method of choice for virtually all genomes sequenced today, and genome assembly algorithms are a critical area of bioinformatics research.

Another aspect of bioinformatics in sequence analysis is annotation. This involves computational gene finding to search for protein-coding genes, RNA genes, and other functional sequences within a genome. Not all of the nucleotides within a genome are part of genes. Within the genomes of higher organisms, large parts of the DNA do not serve any obvious purpose. This so-called junk DNA may, however, contain unrecognized functional elements. Bioinformatics helps to bridge the gap between genome and proteome projects — for example, in the use of DNA sequences for protein identification.

## **GENOME ANNOTATION**

In the context of genomics, annotation is the process of marking the genes and other biological features in a DNA sequence. The first genome annotation software system was designed in 1995 by Owen White, who was part of the team at The Institute for Genomic Research that sequenced and analyzed the first genome of a free-living organism to be decoded, the bacterium *Haemophilus influenzae*. White built a software system to find the genes (fragments of genomic sequence that encode proteins), the transfer RNAs, and to make initial assignments of function to those genes. Most current genome annotation systems work similarly, but the programs available for analysis of genomic DNA, such as the GeneMark program trained and used to find protein-coding genes in *Haemophilus influenzae*, are constantly changing and improving.

## **COMPUTATIONAL EVOLUTIONARY BIOLOGY**

Evolutionary biology is the study of the origin and descent of species, as well as their change over time. Informatics has assisted evolutionary biologists by enabling researchers to:

- trace the evolution of a large number of organisms by measuring changes in their DNA, rather than through physical taxonomy or physiological observations alone,

- more recently, compare entire genomes, which permits the study of more complex evolutionary events, such as gene duplication, horizontal gene transfer, and the prediction of factors important in bacterial speciation,
- build complex computational models of populations to predict the outcome of the system over time
- track and share information on an increasingly large number of species and organisms

Future work endeavours to reconstruct the now more complex tree of life.

The area of research within computer science that uses genetic algorithms is sometimes confused with computational evolutionary biology, but the two areas are not necessarily related.

## **COMPARATIVE GENOMICS**

The core of comparative genome analysis is the establishment of the correspondence between genes (orthology analysis) or other genomic features in different organisms. It is these intergenomic maps that make it possible to trace the evolutionary processes responsible for the divergence of two genomes. A multitude of evolutionary events acting at various organizational levels shape genome evolution. At the lowest level, point mutations affect individual nucleotides. At a higher level, large chromosomal segments undergo duplication, lateral transfer, inversion, transposition, deletion and insertion. Ultimately, whole genomes are involved in processes of hybridization, polyploidization and endosymbiosis, often leading to rapid speciation. The complexity of genome evolution poses many exciting challenges to developers of mathematical models and algorithms, who have recourse to a spectra of algorithmic, statistical and mathematical techniques, ranging from exact, heuristics, fixed parameter and approximation algorithms for problems based on parsimony models to Markov Chain Monte Carlo algorithms for Bayesian analysis of problems based on probabilistic models.

Many of these studies are based on the homology detection and protein families computation.

## **GENETICS OF DISEASE**

With the advent of next-generation sequencing we are obtaining enough sequence data to map the genes of complex diseases such as infertility, breast cancer or Alzheimer's Disease. Genome-wide association studies are essential to pinpoint the mutations for such complex diseases.

## **ANALYSIS OF MUTATIONS IN CANCER**

In cancer, the genomes of affected cells are rearranged in complex or even unpredictable ways. Massive sequencing efforts are used to identify previously unknown point mutations in a variety of genes in cancer. Bioinformaticians continue to produce specialized automated systems to manage the sheer volume of sequence data produced, and they create new algorithms and software to compare the sequencing results to the growing collection of human genome sequences and germline

polymorphisms. New physical detection technologies are employed, such as oligonucleotide microarrays to identify chromosomal gains and losses (called comparative genomic hybridization), and single-nucleotide polymorphism arrays to detect known *point mutations*. These detection methods simultaneously measure several hundred thousand sites throughout the genome, and when used in high-throughput to measure thousands of samples, generate terabytes of data per experiment. Again the massive amounts and new types of data generate new opportunities for bioinformaticians. The data is often found to contain considerable variability, or noise, and thus Hidden Markov model and change-point analysis methods are being developed to infer real copy number changes.

Another type of data that requires novel informatics development is the analysis of lesions found to be recurrent among many tumors.

## **GENE AND PROTEIN EXPRESSION**

### **ANALYSIS OF GENE EXPRESSION**

The expression of many genes can be determined by measuring mRNA levels with multiple techniques including microarrays, expressed cDNA sequence tag (EST) sequencing, serial analysis of gene expression (SAGE) tag sequencing, massively parallel signature sequencing (MPSS), RNA-Seq, also known as "Whole Transcriptome Shotgun Sequencing" (WTSS), or various applications of multiplexed in-situ hybridization. All of these techniques are extremely noise-prone and/or subject to bias in the biological measurement, and a major research area in computational biology involves developing statistical tools to separate signal from noise in high-throughput gene expression studies. Such studies are often used to determine the genes implicated in a disorder: one might compare microarray data from cancerous epithelial cells to data from non-cancerous cells to determine the transcripts that are up-regulated and down-regulated in a particular population of cancer cells.

### **ANALYSIS OF PROTEIN EXPRESSION**

Protein microarrays and high throughput (HT) mass spectrometry (MS) can provide a snapshot of the proteins present in a biological sample. Bioinformatics is very much involved in making sense of protein microarray and HT MS data; the former approach faces similar problems as with microarrays targeted at mRNA, the latter involves the problem of matching large amounts of mass data against predicted masses from protein sequence databases, and the complicated statistical analysis of samples where multiple, but incomplete peptides from each protein are detected.

### **ANALYSIS OF REGULATION**

Regulation is the complex orchestration of events starting with an extracellular signal such as a hormone and leading to an increase or decrease in the activity of one or more proteins. Bioinformatics techniques have been applied to explore various steps in this process. For example, promoter analysis involves the identification and study of sequence motifs in the DNA surrounding the coding region of a gene. These motifs influence the extent to which that region is transcribed into mRNA. Expression data can be used to infer gene regulation: one might compare microarray data from a wide

variety of states of an organism to form hypotheses about the genes involved in each state. In a single-cell organism, one might compare stages of the cell cycle, along with various stress conditions (heat shock, starvation, etc.). One can then apply clustering algorithms to that expression data to determine which genes are co-expressed. For example, the upstream regions (promoters) of co-expressed genes can be searched for over-represented regulatory elements. Examples of clustering algorithms applied in gene clustering are k-means clustering, self-organizing maps (SOMs), hierarchical clustering, and consensus clustering methods such as the Bi-CoPaM. The later, namely Bi-CoPaM, has been actually proposed to address various issues specific to gene discovery problems such as consistent co-expression of genes over multiple microarray datasets.

## **STRUCTURAL BIOINFORMATICS**

### **PREDICTION OF PROTEIN STRUCTURE**

Protein structure prediction is another important application of bioinformatics. The amino acid sequence of a protein, the so-called primary structure, can be easily determined from the sequence on the gene that codes for it. In the vast majority of cases, this primary structure uniquely determines a structure in its native environment. (Of course, there are exceptions, such as the bovine spongiform encephalopathy – a.k.a. Mad Cow Disease – prion.) Knowledge of this structure is vital in understanding the function of the protein. For lack of better terms, structural information is usually classified as one of *secondary*, *tertiary* and *quaternary* structure. A viable general solution to such predictions remains an open problem. Most efforts have so far been directed towards heuristics that work most of the time.

One of the key ideas in bioinformatics is the notion of homology. In the genomic branch of bioinformatics, homology is used to predict the function of a gene: if the sequence of gene *A*, whose function is known, is homologous to the sequence of gene *B*, whose function is unknown, one could infer that *B* may share *A*'s function. In the structural branch of bioinformatics, homology is used to determine which parts of a protein are important in structure formation and interaction with other proteins. In a technique called homology modeling, this information is used to predict the structure of a protein once the structure of a homologous protein is known. This currently remains the only way to predict protein structures reliably.

One example of this is the similar protein homology between hemoglobin in humans and the hemoglobin in legumes (leghemoglobin). Both serve the same purpose of transporting oxygen in the organism. Though both of these proteins have completely different amino acid sequences, their protein structures are virtually identical, which reflects their near identical purposes.

Other techniques for predicting protein structure include protein threading and *de novo* (from scratch) physics-based modeling.

## **NETWORK AND SYSTEMS BIOLOGY**

Main articles: Computational systems biology, Biological network and Interactome

Network analysis seeks to understand the relationships within biological networks such as metabolic or protein-protein interaction networks. Although biological networks can be constructed from a single type of molecule or entity (such as genes), network biology often attempts to integrate many different data types, such as proteins, small molecules, gene expression data, and others, which are all connected physically and/or functionally.

Systems biology involves the use of computer simulations of cellular subsystems (such as the networks of metabolites and enzymes which comprise metabolism, signal transduction pathways and gene regulatory networks) to both analyze and visualize the complex connections of these cellular processes. Artificial life or virtual evolution attempts to understand evolutionary processes via the computer simulation of simple (artificial) life forms.

## **MOLECULAR INTERACTION NETWORKS**

Tens of thousands of three-dimensional protein structures have been determined by X-ray crystallography and protein nuclear magnetic resonance spectroscopy (protein NMR) and a central question in structural bioinformatics is whether it is practical to predict possible protein-protein interactions only based on these 3D shapes, without performing protein-protein interaction experiments. A variety of methods have been developed to tackle the protein-protein docking problem, though it seems that there is still much work to be done in this field.

Other interactions encountered in the field include Protein-ligand (including drug) and protein-peptide. Molecular dynamic simulation of movement of atoms about rotatable bonds is the fundamental principle behind computational algorithms, termed docking algorithms, for studying molecular interactions.

## **OTHERS**

### **LITERATURE ANALYSIS**

The growth in the number of published literature makes it virtually impossible to read every paper, resulting in disjointed sub-fields of research. Literature analysis aims to employ computational and statistical linguistics to mine this growing library of text resources. For example:

- abbreviation recognition – identify the long-form and abbreviation of biological terms,
- named entity recognition – recognizing biological terms such as gene names
- protein-protein interaction – identify which proteins interact with which proteins from text

The area of research draws from statistics and computational linguistics.

### **HIGH-THROUGHPUT IMAGE ANALYSIS**

Computational technologies are used to accelerate or fully automate the processing, quantification and analysis of large amounts of high-information-content biomedical

imagery. Modern image analysis systems augment an observer's ability to make measurements from a large or complex set of images, by improving accuracy, objectivity, or speed. A fully developed analysis system may completely replace the observer. Although these systems are not unique to biomedical imagery, biomedical imaging is becoming more important for both diagnostics and research. Some examples are:

- high-throughput and high-fidelity quantification and sub-cellular localization (high-content screening, cytohistopathology, Bioimage informatics)
- morphometrics
- clinical image analysis and visualization
- determining the real-time air-flow patterns in breathing lungs of living animals
- quantifying occlusion size in real-time imagery from the development of and recovery during arterial injury
- making behavioral observations from extended video recordings of laboratory animals
- infrared measurements for metabolic activity determination
- inferring clone overlaps in DNA mapping, e.g. the Sulston score

## **HIGH-THROUGHPUT SINGLE CELL DATA ANALYSIS**

Computational techniques are used to analyse high-throughput, low-measurement single cell data, such as that obtained from flow cytometry. These methods typically involve finding populations of cells that are relevant to a particular disease state or experimental condition.

## **BIODIVERSITY INFORMATICS**

Biodiversity informatics deals with the collection and analysis of biodiversity data, such as taxonomic databases, or microbiome data. Examples of such analyses include phylogenetics, niche modelling, species richness mapping, or species identification tools.

## **DATABASES**

Databases are essential for bioinformatics research and applications. There is a huge number of available databases covering almost everything from DNA and protein sequences, molecular structures, to phenotypes and biodiversity. Databases generally fall into one of three types. Some contain data resulting directly from empirical methods such as gene knockouts. Others consist of predicted data, and most contain data from both sources. There are meta-databases that incorporate data compiled from multiple other databases. Some others are specialized, such as those specific to an organism. These databases vary in their format, way of accession and whether they are public or not. Some of the most commonly used databases are listed below. For a more comprehensive list, please check the link at the beginning of the subsection.

- Used in Motif Finding: GenomeNet MOTIF Search
- Used in Gene Ontology: DAVID, FuncAssociate, GATHER
- Used in Gene Finding: Hidden Markov Model
- Used in finding Protein Structures/Family: PFAM

- Used for Next Generation Sequencing: (Not database but data format), FASTQ Format
- Used in Gene Expression Analysis: GEO
- Used in Network Analysis: Interaction Analysis Databases(BioGRID, MINT, HPRD), Functional Networks (STRING, KEGG)

Please keep in mind that this is a quick sampling and generally most computation data is supported by wet lab data as well.

## **SOFTWARE AND TOOLS**

Software tools for bioinformatics range from simple command-line tools, to more complex graphical programs and standalone web-services available from various bioinformatics companies or public institutions.

## **OPEN-SOURCE BIOINFORMATICS SOFTWARE**

Many free and open-source software tools have existed and continued to grow since the 1980s. The combination of a continued need for new algorithms for the analysis of emerging types of biological readouts, the potential for innovative *in silico* experiments, and freely available open code bases have helped to create opportunities for all research groups to contribute to both bioinformatics and the range of open-source software available, regardless of their funding arrangements. The open source tools often act as incubators of ideas, or community-supported plug-ins in commercial applications. They may also provide *de facto* standards and shared object models for assisting with the challenge of bioinformation integration.

The range of open-source software packages includes titles such as Bioconductor, BioPerl, Biopython, BioJava, BioRuby, Bioclipse, EMBOSS, .NET Bio, Taverna workbench, and UGENE. In order to maintain this tradition and create further opportunities, the non-profit Open Bioinformatics Foundation have supported the annual Bioinformatics Open Source Conference (BOSC) since 2000.

## **WEB SERVICES IN BIOINFORMATICS**

SOAP- and REST-based interfaces have been developed for a wide variety of bioinformatics applications allowing an application running on one computer in one part of the world to use algorithms, data and computing resources on servers in other parts of the world. The main advantages derive from the fact that end users do not have to deal with software and database maintenance overheads.

Basic bioinformatics services are classified by the EBI into three categories: SSS (Sequence Search Services), MSA (Multiple Sequence Alignment), and BSA (Biological Sequence Analysis).

The availability of these service-oriented bioinformatics resources demonstrate the applicability of web-based bioinformatics solutions, and range from a collection of standalone tools with a common data format under a single, standalone or web-based interface, to integrative, distributed and extensible bioinformatics workflow management systems.

## **BIOINFORMATICS WORKFLOW MANAGEMENT SYSTEMS**

A Bioinformatics workflow management system is a specialized form of a workflow management system designed specifically to compose and execute a series of computational or data manipulation steps, or a workflow, in a Bioinformatics application. Such systems are designed to

- provide an easy-to-use environment for individual application scientists themselves to create their own workflows
- provide interactive tools for the scientists enabling them to execute their workflows and view their results in real-time
- simplify the process of sharing and reusing workflows between the scientists.
- enable scientists to track the provenance of the workflow execution results and the workflow creation steps.

Currently, there are at least four platforms giving this service: Galaxy, Kepler, Taverna Anduril and Anvaya.

## **EDUCATION PLATFORMS**

Software platforms designed to teach bioinformatics concepts and methods include Rosalind and online courses offered through the Swiss Institute of Bioinformatics Training Portal.

## **CONFERENCES**

There are several large conferences that are concerned with bioinformatics. Some of the most notable examples are Intelligent Systems for Molecular Biology (ISMB), European Conference on Computational Biology (ECCB), Research in Computational Molecular Biology (RECOMB) and American Society of Mass Spectrometry (ASMS).