## What is bioinformatics?

an interdisciplinary field at the interface of the computational and life sciences

"The ultimate goal of the field is to enable the discovery of new biological insights as well as to create a global perspective from which unifying principles in biology can be discerned. "

National Center for Biotechnology Information
http://www.ncbi.nlm.nih.gov/

## What is bioinformatics?

- the analysis and interpretation of nucleotide and protein sequences and structures

## What is bioinformatics?

- the analysis and interpretation of nucleotide and protein sequences and structures
- development of algorithms and software to support the acquisition … of biomolecular data

## What is bioinformatics?

- the analysis and interpretation of nucleotide and protein sequences and structures
- development of algorithms and software to support the acquisition … of biomolecular data
- the development of software that enables efficient access and management of biomolecular information

## Bioinformatics stems from parallel revolutions in biology and computing

At the beginning of World War II (1939-1944):

- The shared program computer had not yet been invented, and there were no programming languages, databases, or computer networks.

- The relationship between genes and proteins, the molecular basis of genes, the structure of DNA and the genetic code were all unknown.
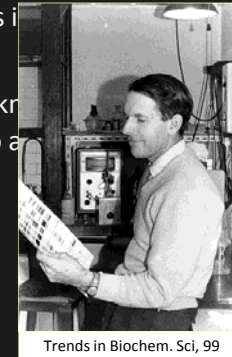
## The Origins of Computational Biology
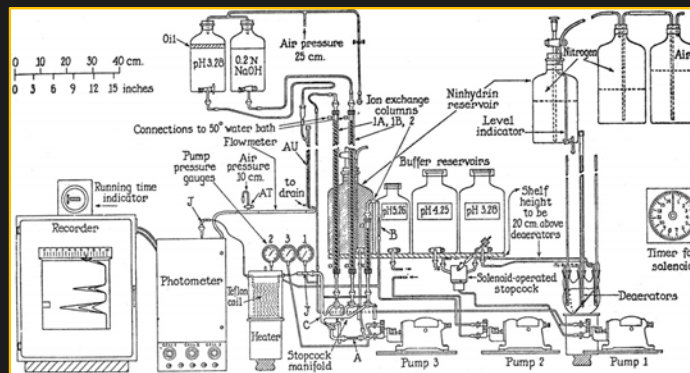
Amino acid sequencing

Sanger sequences i

Turing designs a stored
rogram computer

Stein, Moore, Spackr
automatic amino a
analyzer

ac: 1st stored
rogram computer



Trends in Biochem. Sci, 99

## Automatic recording apparatus used in the chromatographic analysis of mixtures of amino acids



Stein, Moore, Spackman, 1958

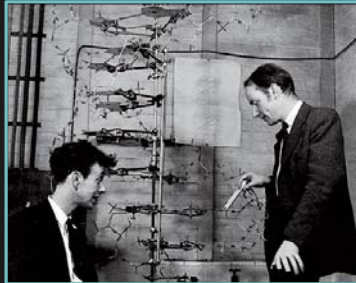## EDSAC: The first stored-program computer.

## The Origins of Computational Biology

Amino acid sequencing
   Sanger sequences insulin.

**1950**

Stein, Moore, Spackman:
   automatic amino acid
   analyzer

Turing designs a stored
   program computer

Edsac: 1st stored
   program computer

Grace Murray Hopper
   co-invents COBOL and
   finds the first bug

---



Grace Murray Hopper finds the first bug

*Rear Admiral Grace Hopper over-seeing her team of programmers*, Philadelphia Inquirer, 1957

---



Grace Murray Hopper

---

Discovery of DNA structure,
   Watson, Crick, Franklin

3

## The helical structure of DNA

**James Watson and Francis Crick**

**Rosalind Franklin**

---

## The Origins of Computational Biology

Determination of the
genetic code

Fortran, Basic, LISP

**1960**

**1970**

---

## The genetic code

The genetic code is a degenerate, non-overlapping,
triplet code. Crick, Barnett, Brenner, and Watts-Tobin, 1961

Determination of
the genetic code



Genetic code table:

| | | Second Base | | | |
| --- | --- | --- | --- | --- | --- |
| | U | C | A | G | |
| U | UUU / UUC Phe, UUA / UUG Leu | UCU / UCC / UCA / UCG Ser | UAU / UAC Tyr, UAA Stop / UAG Stop | UGU / UGC Cys, UGA Stop, UGG Trp | U C A G |
| C | CUU / CUC / CUA / CUG Leu | CCU / CCC / CCA / CCG Pro | CAU / CAC His, CAA / CAG Gln | CGU / CGC / CGA / CGG Arg | U C A G |
| A | AUU / AUC Ile, AUA, AUG Met / Start | ACU / ACC / ACA / ACG Thr | AAU / AAC Asn, AAA / AAG Lys | AGU / AGC Ser, AGA / AGG Arg | U C A G |
| G | GUU / GUC / GUA / GUG Val | GCU / GCC / GCA / GCG Ala | GAU / GAC Asp, GAA / GAG Glu | GGU / GGC / GGA / GGG Gly | U C A G |

---

## The Origins of Computational Biology

Determination of the
genetic code

Fortran, Basic, LISP

**1960**

On going protein sequen-
cing, Dayhoff publishes
the Protein Atlas

**1970**

4

## Atlas of Protein Sequence & Structure
## 1965 - 1978

Margaret Dayhoff
PhD in Chemistry, 47
Watson Computing Lab Fellow  47 - 48

---

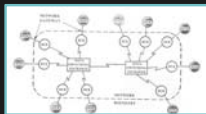## The Origins of Computational Biology

Determination of the
genetic code

Fortran,  Basic, LISP

**1960**

On going protein sequen-
cing, Dayhoff publishes
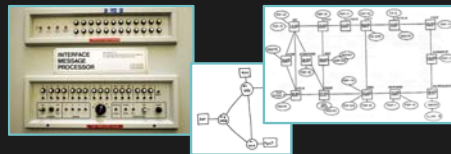the Protein Atlas

Arpanet

**1970**

---

## The ARPAnet is established,
## precursor to the Internet.

**Packet switching  is invented, which supports flexible, multi-node network**

**Interface Message Processor (IMP) network is constructed, linking 4, and then 15 nodes**

---

## The Origins of Computational Biology

Determination of the
genetic code

Fortran,  Basic, LISP

**1960**

On going protein sequen-
cing, Dayhoff publishes
the Protein Atlas
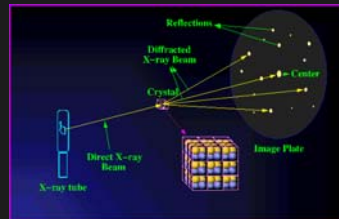
Arpanet

The Protein Data Bank is
established.

**1970**

## Protein Data Bank (PDB)

Growing collection of X-ray diffraction protein structure data

Development of molecular graphics display for 3D visualization
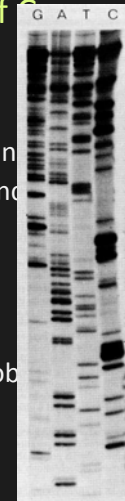
SEARCH program supports remote access



## The Origins of Computational Biology

Sanger-Coulson sequen...
Maxam-Gilbert sequenc...

TCP/IP
Internet

Gilbert, Sanger win Nob...
Prize

First royal email
USENET newsgroups



## The Origins of Computational Biology

**1970**

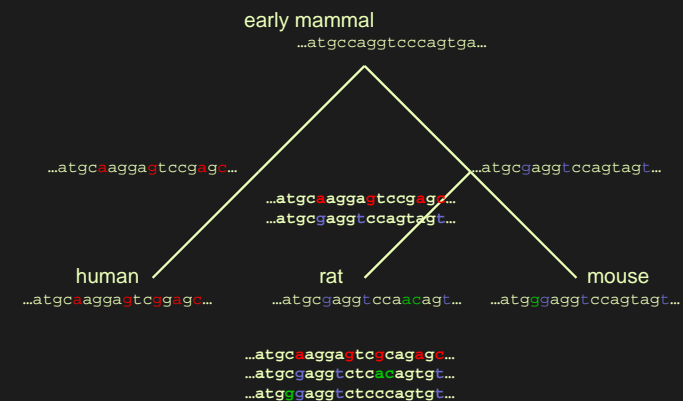Sanger-Coulson sequencing
Maxam-Gilbert sequencing

TCP/IP
Internet

Gilbert, Sanger win Nobel
Prize      **1980**

First royal email
USENET newsgroups

## Beginnings of molecular evolution

early mammal
...atgccaggtcccagtga...

...atgcaaggagtccgagc...

...atgcaaggagtccgagc...
...atgcgaggtccagtagt...

...atgcgaggtccagtagt...

human
...atgcaaggagtcgcagagc...

rat
...atgcgaggtccaacagt...

mouse
...atgggaggtccagtagt...

...atgcaaggagtcgcagagc...
...atgcgaggtctcacagtgt...
...atgggaggtctcccagtgt...

## Sequence similarity → structural similarity

**Structure?**

…v**k**ltpe**g**tr_wggh**p**ldekflske…

…v**h**ltpe**ttr**g**wgghmldek**ei**ske…



Estimate protein structure from
 a related protein with known
structure and similar sequence.

## Sequence similarity → structural similarity

**Structure?**

…v**k**ltpe**g**tr_wggh**p**ldekflske…
…v**h**ltpe**ttr**g**wgghmldek**ei**ske…



Estimate protein structure from a
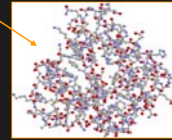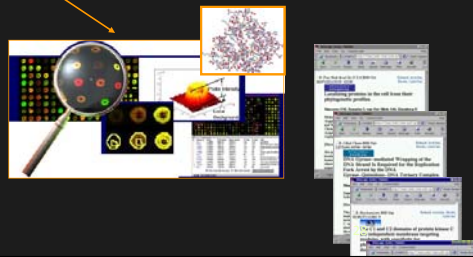related protein with known
structure and similar sequence.

## Sequence similarity → functional similarity

**?**

...atgc**g**agga**ttc**t**cg**aa**gacag**cg**a…
...atgc**a**agga**g**tc_cg**tt**gacag**ag**c…



## The Origins of Computational Biology
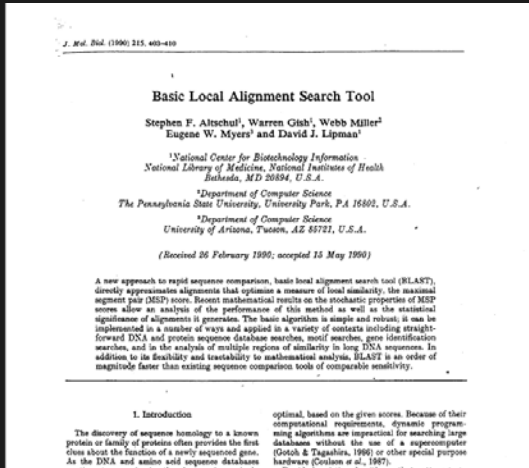
Congress establishes
  Genbank          **1990**

Human Genome Project
  begins (DOE)

Basic local alignment search
  tool (BLAST)

## BLAST



J. Mol. Biol. (1990) 215, 403-410

### Basic Local Alignment Search Tool

Stephen F. Altschul[1], Warren Gish[1], Webb Miller[2]
Eugene W. Myers[3] and David J. Lipman[1]

[1]National Center for Biotechnology Information
National Library of Medicine, National Institutes of Health
Bethesda, MD 20894, U.S.A.

[2]Department of Computer Science
The Pennsylvania State University, University Park, PA 16802, U.S.A.

[3]Department of Computer Science
University of Arizona, Tucson, AZ 85721, U.S.A.

(Received 26 February 1990; accepted 15 May 1990)

A new approach to rapid sequence comparison, basic local alignment search tool (BLAST), directly approximates alignments that optimize a measure of local similarity, the maximal segment pair (MSP) score. Recent mathematical results on the stochastic properties of MSP scores allow an analysis of the performance of this method as well as the statistical significance of alignments it generates. The basic algorithm is simple and robust; it can be implemented in a number of ways and applied in a variety of contexts including straightforward DNA and protein sequence database searches, motif searches, gene identification searches, and in the analysis of multiple regions of similarity in long DNA sequences. In addition to its flexibility and tractability to mathematical analysis, BLAST is an order of magnitude faster than existing sequence comparison tools of comparable sensitivity.

---

## The Origins of Computational Biology

Congress establishes
Genbank

**1990**

Human Genome Project
begins (DOE)

Information
superhighway

Basic local alignment search
tool (BLAST)

World Wide Web,

NCSA Mosaic

---

## The internet we know and love

**Al Gore Creates Bill to Fund "Information Superhighway"**





**Tim Berners-Lee proposes a design for information sharing that becomes the World Wide Web**



**The first web browser**

---

## The Origins of Computational Biology

Congress establishes
Genbank

**1990**

Human Genome Project
begins (DOE)

Information
superhighway

Basic local alignment search
tool (BLAST)

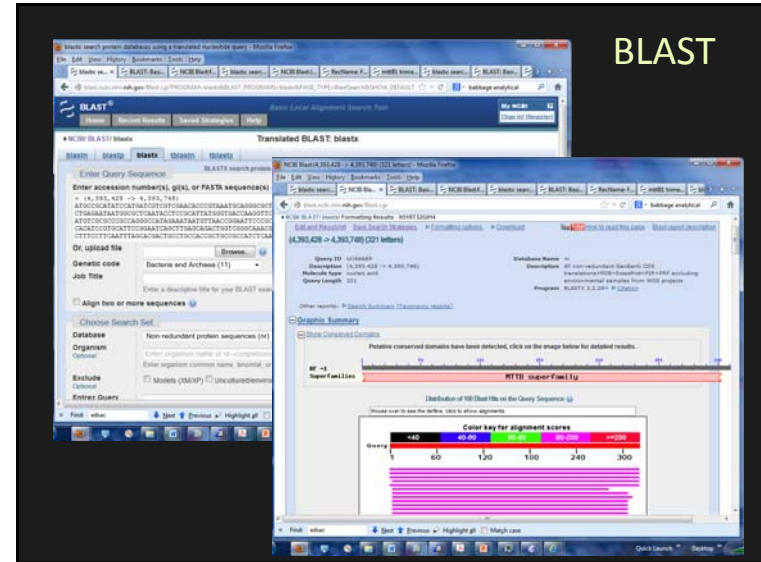World Wide Web,

GenBank goes online.
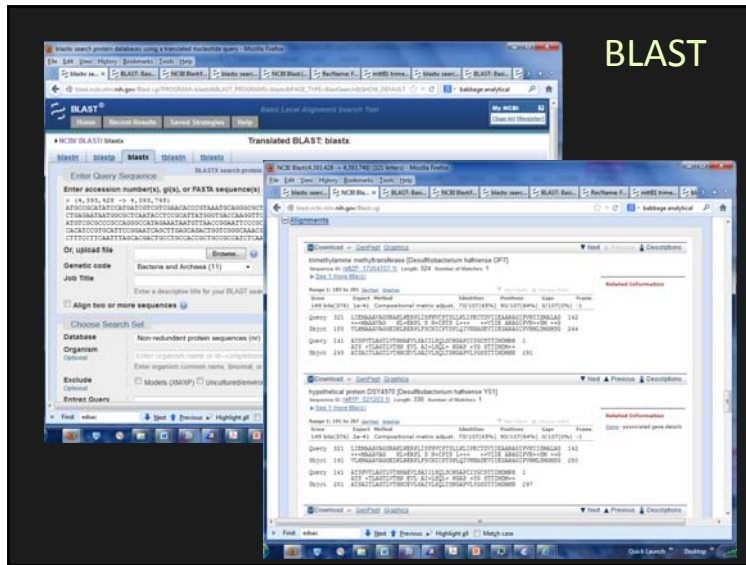
NCSA Mosaic

Pizza Hut goes on line

## Genbank

National Center for Biotechnology Information
http://www.ncbi.nlm.nih.gov/

## BLAST

## BLAST

## The Origins of Computational Biology

|  |  | Google |
| --- | --- | --- |
| *H. Influenzae: 1*st whole genome sequence | **1995** |  |
| Yeast – 1st eukaryote |  |  |
|  | **1998** |  |
| *C. elegans* - 1st animal |  | IBM's Blue Gene architecture for protein modeling |
| Fly genome | **1999** |  |
| *A.thaliana* – 1st plant |  |  |
|  | **2000** |  |

9

## Sequence Assembly

Limits of gel electrophoresis: ~ 500bp in one "read"

To sequence more than 500 bp:

- Sequence 500bp fragments separately
- Combine *computationally* using sequence comparison
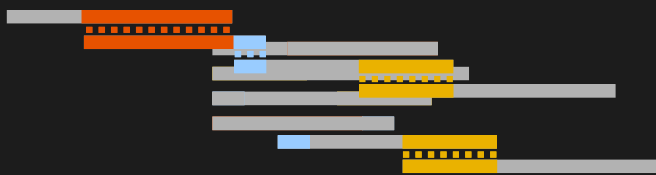


## Shotgun sequencing



DNA molecule

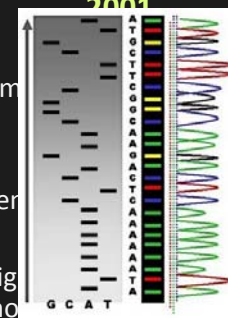shotgun reads

*Sequence assembly*

contigs

## Sequence Assembly



Alignments do not necessarily imply sequence identity

## The Origins of Computational Biology

**2001**

Draft human genom

ABI Capillary sequen

Chimp genome, alig
with human geno

pedia, Captcha

la Foundation:
en source, W3P
tent policy



10

## How to interpret whole genome sequence

- Where are the genes?
- When and where are those genes expressed?
- What proteins do they encode?
- What do they do?
  - Molecular function?
  - Biological pathway or process?

## How to interpret whole genome sequence

- ➢ Where are the genes?
- When and where are those genes expressed?
- What proteins do they encode?
- What do they do?
  - Molecular function?
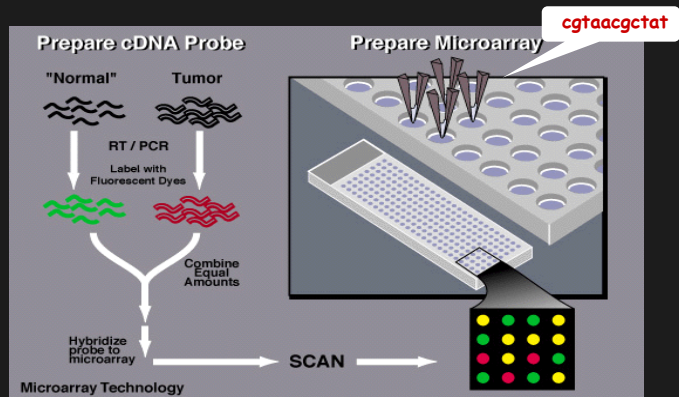  - Biological pathway or process?

```
atatactcacagcataactgtatatacacccaggggggcggaatgaaagcgttaacggcca
ggcaacaagaggtgtttgatctcatccgtgatcacatcagccagacaggtatgccgccga
cgcgtgcggaaatcgcgcagcgtttggggttccgttccccaaacgcggctgaagaacatc
tgaaggcgctggcacgcaaagcgttattgaaattgtttccggcgcatcacgcgggattc
gtctgttgcaggaagaggaagaaggggttgccgctggtaggtcgtgtggctgccggtgaac
cacttctggcgcaacagcatattgaaggtcattatcaggtcgatccttccttattcaagc
cgaatgctgatttcctgctgcgcgtcagcgggatgtcgatgaaagatatcggcattatgg
atggtgacttgctggcagtgcataaaactcaggatgtacgtaacggtcaggtcgttgtcg
cacgtattgatgacgaagttaccgttaagcgcctgaaaaaacagggcaataaagtcgaac
tgttgccagaaaatagcgagtttaaaccaattgtcgttgaccttcgtcagcagagcttca
cgcgtgcggaaatcgcgcagcgtttggggttccgttccccaaacgcggctgaagaacatc
tgaaggcgctggcacgcaaagcgttattgaaattgtttccggcgcatcacgcgggattc
gtctgttgcaggaagaggaagaaggggttgccgctggtaggtcgtgtggctgccggtgaac
ccattgaaaggctggcggttggggttattcgcaacggcgactggctgtaacatatctctg
gaattcgataaatctctggtttattgtgcagtttatggttccaaaatcgccttttgctgt
agaccgcgatgccgcctggcgtcgcggtttgtttttcatctctcttcatcaggcttgtct
gcatggcattcctcacttcatctgataaagcactctggcatctcgccttacccatgattt
cgaatgctgatttcctgctgcgcgtcagcgggatgtcgatgaaagatatcggcattatgg
atggtgacttgctggcagtgcataaaactcaggatgtacgtaacggtcaggtcgttgtcg
cacgtattgatgacgaagttaccgttaagcgcctgaaaaaacagggcaataaagtcgaac
tgttgccagaaaatagcgagtttaaaccaattgtcgttgaccttcgtcagcagagcttca
ccattgaaaggctggcggttggggttattcgcaacggcgactggctgtaacatatctctg
agaccgcgatgccgcctggcgtcgcggtttgtttttcatctctcttcatcaggcttgtct
gcatggcattcctcacttcatctgataaagcactctggcatctcgccttacccatgattt
tctccaatatcaccgttccgttgctgggactggtcgatacggcggtaattggtcatcttg
```

## DNA PATTERNS IN THE *E.coli* lexA GENE



```
                    Repressor binding site              Promotor sequences
1   gaattcgataaatcctggtttattgtgcagtttatggttccaaaatcgccttttgctgt
                                                      TTCCAA -35
61  atatactcacagcataactgtatatacacccagggggcggaatgaaagcgttaacggcca
-10  TATACT   mRNAstart+       +10GGGGG Ribosomal binding site
121 ggcaacaagaggtgtttgatctcatccgtgatcacatcagccagacaggtatgccgccga
181 cgcgtgcggaaatcgcgcagcgtttggggttccgttccccaaacgcggctgaagaacatc
241 tgaaggcgctggcacgcaaagcgttattgaaattgtttccggcgcatcacgcgggattc
301 gtctgttgcaggaagaggaagaaggggttgccgctggtaggtcgtgtggctgccggtgaac
361 cacttctggcgcaacagcatattgaaggtcattatcaggtcgatccttccttattcaagc
421 cgaatgctgatttcctgctgcgcgtcagcgggatgtcgatgaaagatatcggcattatgg        ATG..TAA
481 atggtgacttgctggcagtgcataaaactcaggatgtacgtaacggtcaggtcgttgtcg    open reading frame
541 cacgtattgatgacgaagttaccgttaagcgcctgaaaaaacagggcaataaagtcgaac
601 tgttgccagaaaatagcgagtttaaaccaattgtcgttgaccttcgtcagcagagcttca
661 ccattgaaaggctggcggttggggttattcgcaacggcgactggctgtaacatatctctg
721 agaccgcgatgccgcctggcgtcgcggtttgtttttcatctctcttcatcaggcttgtct
781 gcatggcattcctcacttcatctgataaagcactctggcatctcgccttacccatgattt
841 tctccaatatcaccgttccgttgctgggactggtcgatacggcggtaattggtcatcttg
901 atagcccggtttatttgggcggcgtggcggttggcgcaacggcggaccagct
```
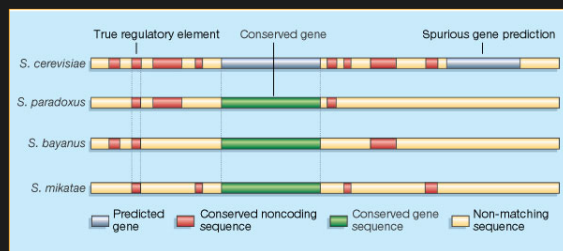
## DNA microarrays



cgtaacgctat

## The Origins of Computational Biology

5 more yeast genomes **2004**

**2005**

**2006**

## Align genomes to confirm gene predictions and identify regulatory regions



Salzberg, Nature, 2003

## Global alignment of upstream sequences to identify regulatory regions



Kellis *et al, Nature*, 03

## The Origins of Computational Biology

5 more yeast genomes        **2004**

454 pyrosequencer                    Rosetta@home

  Hi-thruput, short read               Facebook,  Twitter
  sequencing
                            **2005**
12 *Drosophila* genomes              US cyberworm attacks
                                       Iranian centrifuges

                            **2006**

## Next-generation, short read sequencing

<u>Sanger Sequencing</u>                <u>Illumina sequending</u>

• read lengths up to 1,000 bp        • read lengths up to 36 bp
• accuracy 99.999%                   • error rates 1-1.5%
• costs $500 per megabase            •cost $2 per megabase

        <u>454 sequencing</u>

        • read lengths 200-300 bp
        • accuracy problem with homopolymers
        • costs $60 per megabase

## Next-generation. short read sequencing

<u>Advantages</u>

    • High throughput
    •Does not require PCR amplification
    •Accurate measures of abundance
    •Cheaper

<u>Disadvantages</u>

    • Short reads are unlikely to be unique.
    • Difficult  to identify the origin of a given read
    • Particular challenge for genome assembly

13

## Some next generation sequencing applications

- Bacterial genomes

- Sample diversity in a bacterial population (e.g., your throat when you have strep)

- Transcription: more accurate and quantitative compared with microarrays

- Medical diagnostics: sequence short genomic regions to identify mutations associated with disease

## The Origins of Computational Biology

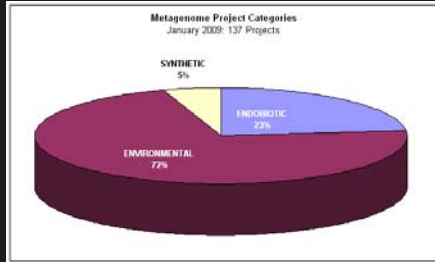| | | |
|---|---|---|
| 1000 Genomes project | **2007** | Apple iPhone |
| | | Estonia: First national elections via Internet |
| Human microbiome project. | | |
| First tumor/normal genome published | **2008** | Foldit: Crowd-sourced protein folding game |
| Draft Neanderthal genome | | |
| | **2009** | |



A crowd-sourcing game developed at CMU.

## Metagenomics

- Sample communities of microbial organisms directly from their natural environments, bypassing the need for isolation and lab cultivation of individual species.

- Result: a collection of DNA fragments that characterize the organismal and functional diversity of the envirment

## Metagenomics

- Production-scale plant fermenter
- Fungal communities from the Arctic
- Singapore indoor air filters
- Yellowstone Obsidian Hot Spring
- Fossil microbiome
- Human microbiome



Metagenome Project Categories
January 2009: 137 Projects

SYNTHETIC 5%
ENDOBIOTIC 23%
ENVIRONMENTAL 72%

---

# What makes us human?

- Human metabolic features- combo of human and microbial traits
- Microbiota- microrganisms that live inside and on humans
- Microbiome- the genomes of the microbial symbionts

---

# The Origins of Computational Biology

**2010**

Chocolate (*Theobroma cacao*) genome

Social networking topples regime in Egypt

**2011**

3rd Generation sequencing: Pac Bio, Ion Torrent

*Crystal structure …solved by protein folding game players,*
Nature Structural Biology

**2012**

---

# What is bioinformatics?

Development of algorithms and software to support the acquisition and interpretation of biomolecular data

- Acquisition:  Microarray design,  Sequence assembly
- Interpretation: Sequence comparison, clustering of microarrays, gene finding, phylogenetic profiling, …

15

## What is bioinformatics?

Development of <u>software</u> that enables <u>efficient access and management</u> of <u>biomolecular information</u>

- Genbank

- BLAST

- Protein visualization

## What is bioinformatics?

The <u>analysis and interpretation</u> of <u>nucleotide</u> and <u>protein sequences</u> and <u>structures</u>

Application of all of these methods to address specific questions about specific biological systems