

**Informatics in Biological Research**

*Agenda for Informatics in the Biological Sciences*

University of Chicago 13 Jan 2003

Peter Cooper  
National Center for Biotechnology Information

NCFBI

---

---

---

---

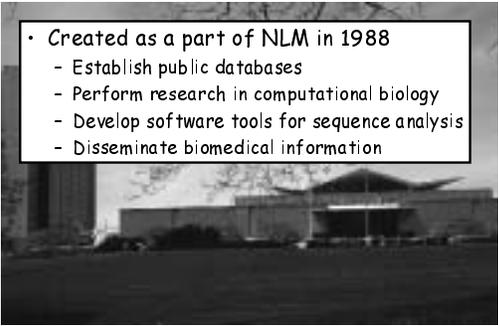
---

---

---

**The National Center for Biotechnology Information**

- Created as a part of NLM in 1988
  - Establish public databases
  - Perform research in computational biology
  - Develop software tools for sequence analysis
  - Disseminate biomedical information



NCFBI

---

---

---

---

---

---

---

**The Revolutions**

- **Biology**
  - DNA sequencing
  - High throughput methodologies
  - Improvements in cloning technology
  - PCR
  - Microarray techniques
- **Computer Science**
  - CPU speed
  - Storage devices
  - The Internet and the WWW

NCFBI

---

---

---

---

---

---

---

## The Selection Pressure

"It's sink or swim as a tidal wave of data approaches"

*Nature 399:517 10 June 1999*

IB2N

---

---

---

---

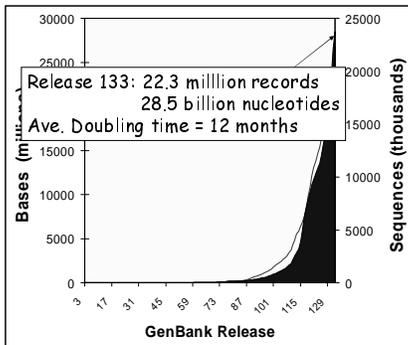
---

---

---

---

## The Growth of GenBank



IB2N

---

---

---

---

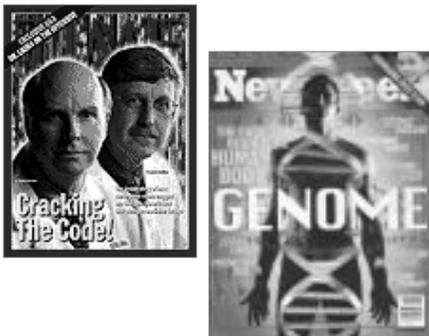
---

---

---

---

## The Genome Sequencing Era



IB2N

---

---

---

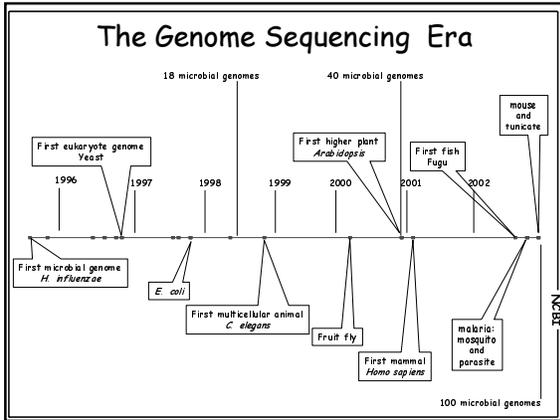
---

---

---

---

---




---

---

---

---

---

---

---

---

---

---

- ### Coming soon ...
- Nearly done
    - rat
    - purple sea urchin
    - zebrafish
  - NHGRI's Priority Organisms
    - Chicken
    - Cow
    - Dog
    - Chimpanzee
    - Honeybee
    - *Tetrahymena*
    - *Oxytricha*
    - Several fungi
  - Over 100 bacterial genomes

---

---

---

---

---

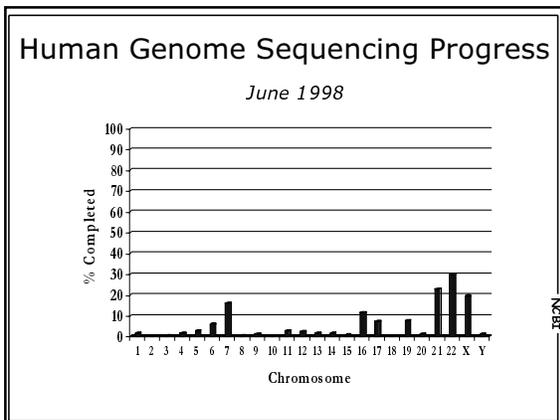
---

---

---

---

---




---

---

---

---

---

---

---

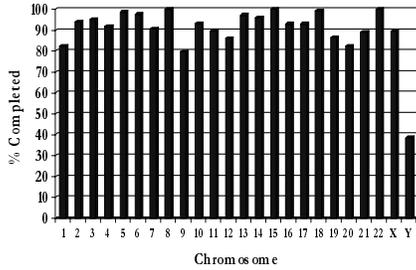
---

---

---

## Human Genome Sequencing Progress

Dec. 2002



IBCN

---

---

---

---

---

---

---

---

---

---

## The Rapid Evolution of Resources

IBCN

---

---

---

---

---

---

---

---

---

---

### Web Access

- Text
  - Entrez
- Sequence
  - BLAST
- Structure
  - VAST

IBCN

---

---

---

---

---

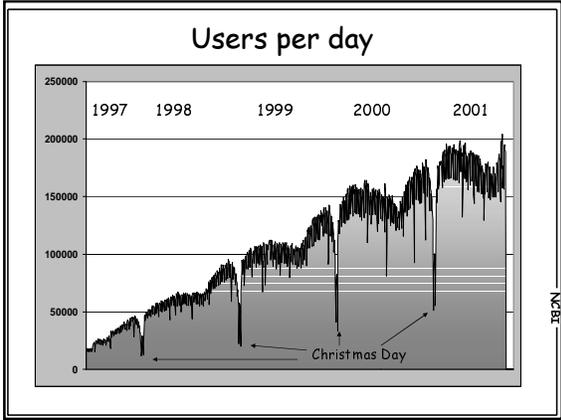
---

---

---

---

---




---

---

---

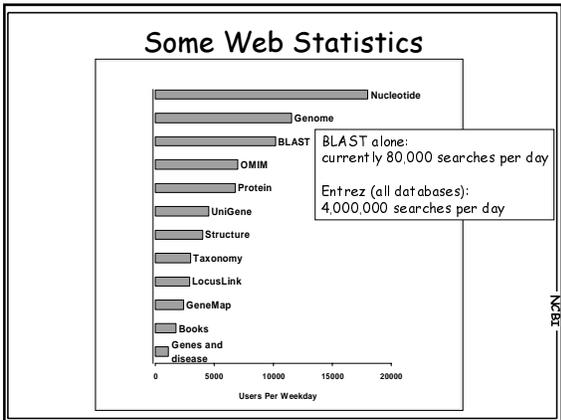
---

---

---

---

---




---

---

---

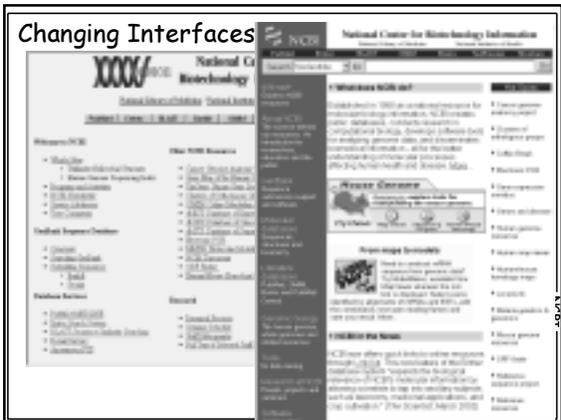
---

---

---

---

---




---

---

---

---

---

---

---

---

### Changing Interfaces




---

---

---

---

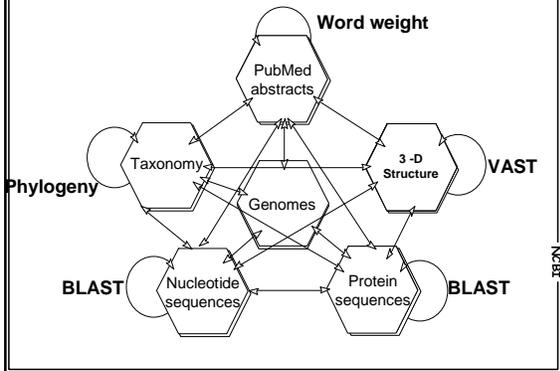
---

---

---

---

### Entrez: Database Integration




---

---

---

---

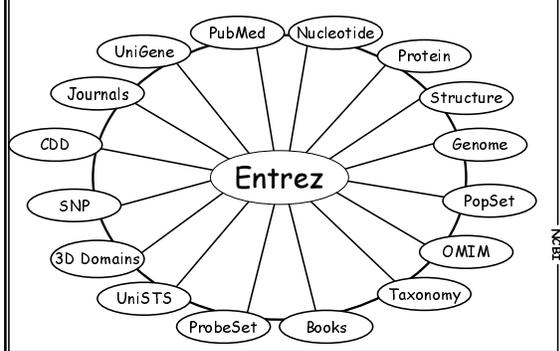
---

---

---

---

### The (ever) Expanding Entrez System




---

---

---

---

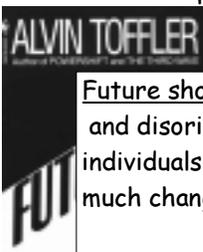
---

---

---

---

**Revolution**



**Future shock:** the shattering stress and disorientation that we induce in individuals by subjecting them to too much change in too short a time.

-Alvin Toffler

MORE THAN 5 MILLION COPIES IN PRINT

---

---

---

---

---

---

---

---

**The Promise of Computational Biology**

*Gene*



> DNA sequence  
AATTGATGAAATGTTATGCTGCTGCTGACCGGCAACAC  
TGGAAATGGGAGACTTATGCTTAAAGGTATGATGAA  
TCTGTAAAGAGCTCAACACCACTCACTGCTGACGTTA  
ACATGAAAGACTGCTGACGAGAGATATCTGATCTGGG  
TCTCTGCTGCTGGGCTGAACTTCTGAGGAGGAA  
TTTGAACGTTTCACTGAGAGATCTCTACCAAAATCTCTG  
GTAGAGAGTTGCTGCTGCTGCTGCTGCTGCTGCTGCTG  
CGTAACTGATGCTGCTGCTGCTGCTGCTGCTGCTGCTG  
TACGCTGCTGCTGCTGCTGCTGCTGCTGCTGCTGCTG  
ACGAGCTGAGAGCTGAGAGACTGCTGCTGCTGCTGCTG  
TAGAGATGCTGAGAGACTGCTGCTGCTGCTGCTGCTG

➔

*Structure & Function*



> Protein sequence  
MELVWIGTQTRKMAELAKGIIESKEDVNTINYSVRI  
SHELAREHLLGQKANDVLEKSEFFYFEELSTKISE  
EVALPGEVWQDKHMDPEERHNTGCVVYKPLLVQNE  
PDEARQDCIEFPGKILANI

*The power of computing  
on the data*

---

---

---

---

---

---

---

---

**The Cross-over to  
"Functional Genomics"**

"In the past we have had functions in search of sequences. In the future, pathology and physiology will become 'functionators' for the sequences."

Daniel C. Tosteson, Dean  
Harvard Medical School  
March 26, 1997

***The future is NOW.***

---

---

---

---

---

---

---

---

## Comparative Analysis of Genomes

NCBI



"What is true for *E. coli* is also true for elephant."  
Jacques Monod, c. 1961



"What is true for yeast is also true for human."  
David Botstein, 1988

IB2N

---

---

---

---

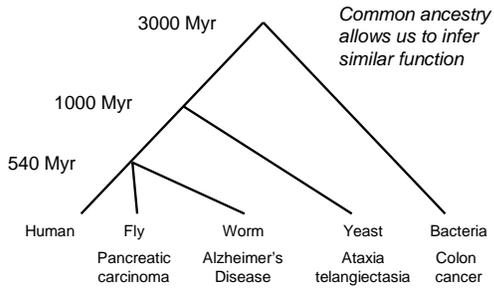
---

---

---

---

## Molecular Evolution



IB2N

---

---

---

---

---

---

---

---

## "Homology..."



... is the central concept for *all* of biology. Whenever we say that a mammalian hormone is the 'same' hormone as a fish hormone, that a human gene sequence is the 'same' as a sequence in a chimp or a mouse, that a HOX gene is the 'same' in a mouse, a fruit fly, a frog, and a human -- even when we argue that discoveries about a worm, a fruit fly, a frog, a mouse, or a chimp have relevance to the human condition -- we have made a bold and direct statement about homology. The aggressive confidence of modern biomedical science implies that we know what we are talking about."

David B. Wake



IB2N

---

---

---

---

---

---

---

---

**Application:**  
**MLH1 Homolog in *Anopheles gambiae***

---

---

---

---

---

---

---

---

**Annotation on the fly: mosquito genome**

Predicted protein

---

---

---

---

---

---

---

---

**Annotation on the fly**

What could it do?

---

---

---

---

---

---

---

---









### Polymorphisms

Chr	Position	Protein	Function	SNP	Protein	Codon	Amino acid
11	100000	100000	100000	100000	100000	100000	100000
11	100000	100000	100000	100000	100000	100000	100000
11	100000	100000	100000	100000	100000	100000	100000

Would this be functionally significant?

---

---

---

---

---

---

---

---

---

---

### BLink

Finding structural models

---

---

---

---

---

---

---

---

---

---

### BLink: Finding Modeling Template

ID	Name	Description	Accession
1	1	1	1
2	2	2	2
3	3	3	3

---

---

---

---

---

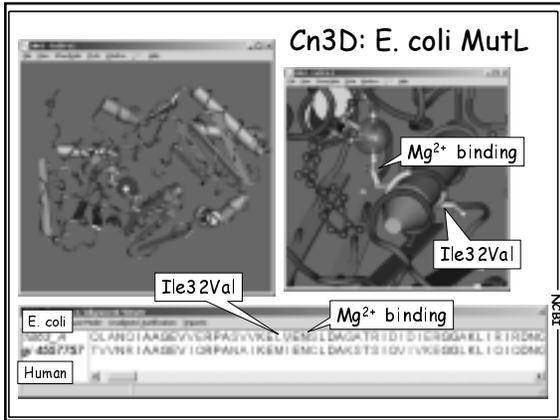
---

---

---

---

---




---

---

---

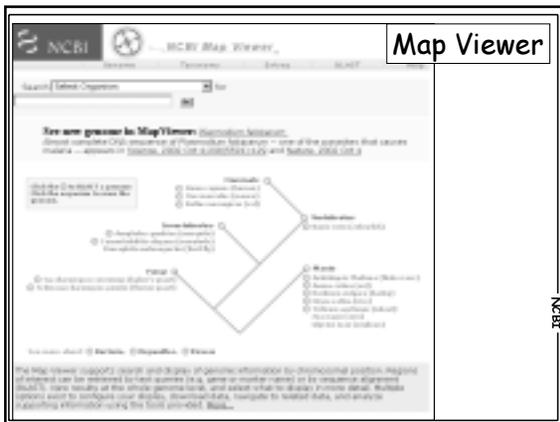
---

---

---

---

---




---

---

---

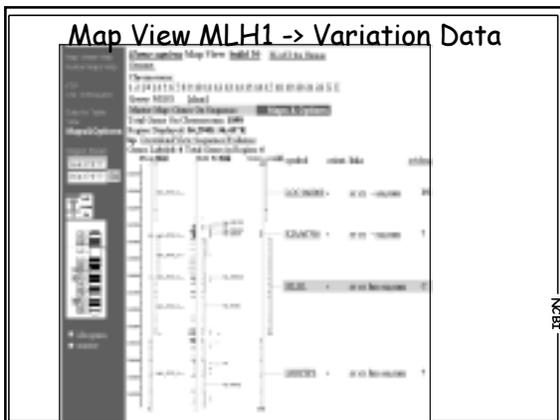
---

---

---

---

---




---

---

---

---

---

---

---

---





