**Boris Steipe**
Department of Biochemistry
Department of Molecular Genetics
University of Toronto

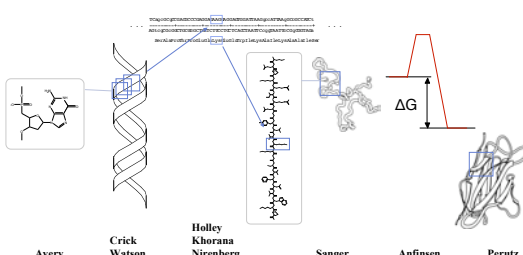B C H 4 4 1      B I O I N F O R M A T I C S

I N T R O D U C T I O N  -  2 0 1 2

---

# BCH441

## Why is Bioinformatics relevant?

---

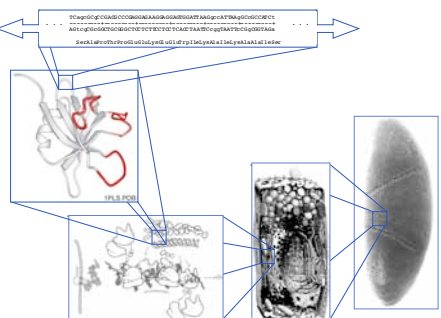# molecular biology

Inheritable information is a substance.



| Avery | Crick Watson | Holley Khorana Nirenberg | Sanger | Anfinsen | Perutz |

---

# genomic biology



---

# (post)genomic biology current practice

1. **Industrial scale** (Data intensive)
2. **Multiple genes** (Cross-sectional)
3. **Model Organisms** (Inference by analogy)
4. **Complete, exhaustive description** (Missing entities are important)
5. **Discovery Science** (Association, not Hypothesis)

**Too much data to process by hand,
too many entities to keep in your mind.**

---

# (post)genomic biology current practice

In any modern life-science paper you will likely find that conclusions derived from computational inference exceed conclusions derived from direct observation.
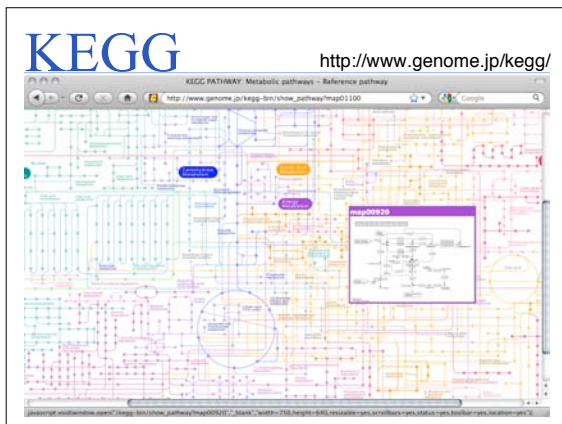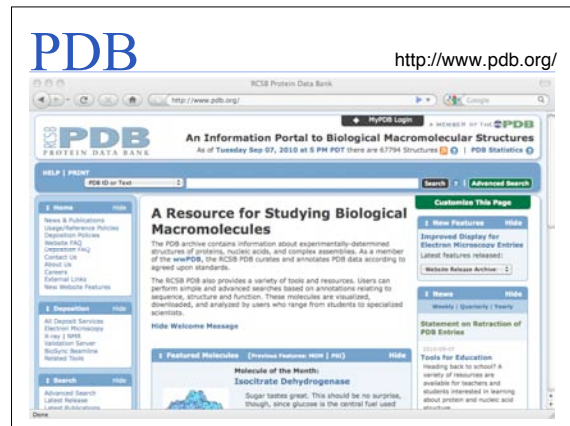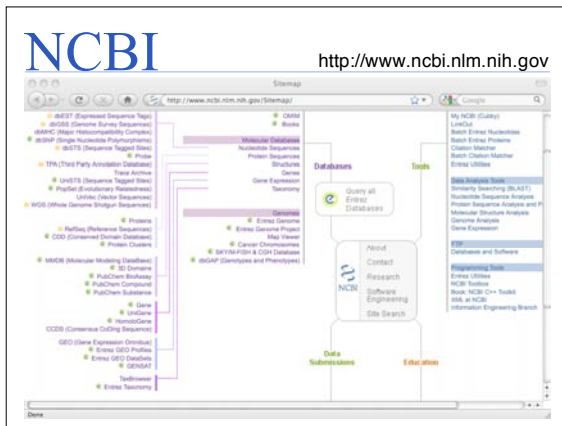
That is to be expected, given the importance of "context".

# BCH441

## What is Bioinformatics ?

# bioinformatics

Data management is the fundamental task of bioinformatics.

# NCBI
http://www.ncbi.nlm.nih.gov



# PDB
http://www.pdb.org/



# KEGG
http://www.genome.jp/kegg/



# challenges

1. Data overload (2009 NAR: 179 databases, 95 new - 1170 in the Molecular Biology Database Collection )
2. Service overload (2009 NAR: 112 Web services)
3. Poor integration
4. Peer review and expert opinions lacking
5. Cultural gap between life- and computer sciences

**How will bioinformatics contribute to our understanding of biology?**

The question becomes less: "What can you do?" but: "What should you do?" !

## bioinformatics

Modeling is the fundamental task of bioinformatics.

## bioinformatics

Problems of modeling:

Models can be right or wrong ...

## bioinformatics

Problems of modeling:

Models can be right or wrong ...
... but worse, they can also be irrelevant.

## *cargo cult* science



## *cargo cult* science

[...] In the South Seas there is a cargo cult of people. During the war they saw airplanes land with lots of good materials, and they want the same thing to happen now. So they've arranged to imitate things like runways, to put fires along the sides of the runways, to make a wooden hut for a man to sit in, with two wooden pieces on his head like headphones and bars of bamboo sticking out like antennas--he's the controller--and they wait for the airplanes to land. They're doing everything right. The form is perfect. It looks exactly the way it looked before. But it doesn't work. No airplanes land. So I call [some examples of pseudoscience] cargo cult science, because they follow all the apparent precepts and forms of scientific investigation, but they're missing something essential, because the planes don't land.

Now it behooves me, of course, to tell you what they're missing. But it would be just about as difficult to explain to the South Sea Islanders how they have to arrange things so that they get some wealth in their system. It is not something simple like telling them how to improve the shapes of the earphones. [...].

*Richard Feynman*

## *cargo cult* science



Avro Lancaster, WikiMedia Commons

# BCH441

What we will cover in this course:
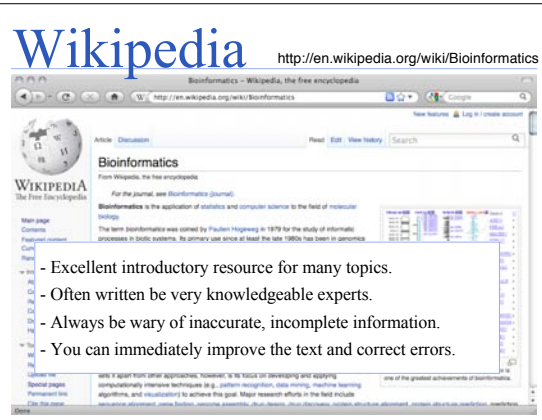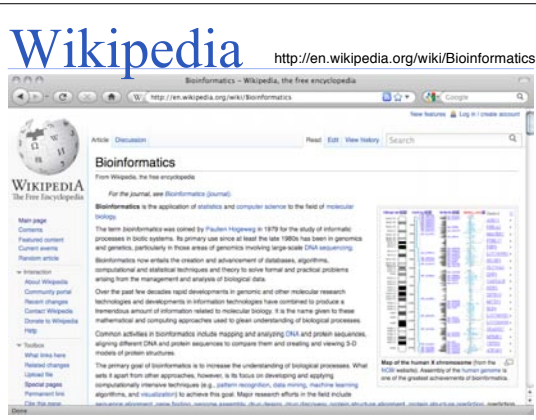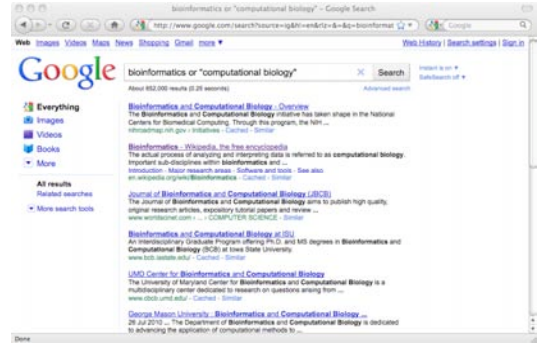
# learn for change

- **We'll discuss principles and examples of how the facts of biology can be expressed as computable abstractions.** As our knowledge of the facts changes, you should be able to think of novel models.

- **We'll use key databases that store publicly available molecular data.** As the databases grow and change, you should be able to work with new types of data, because you are familiar with the principles.

- **We'll use key procedures that analyse sequence, structure, function and phylogeny.** As new tools become available, you should be able to identify those that are useful to support **your own, changing objectives**.
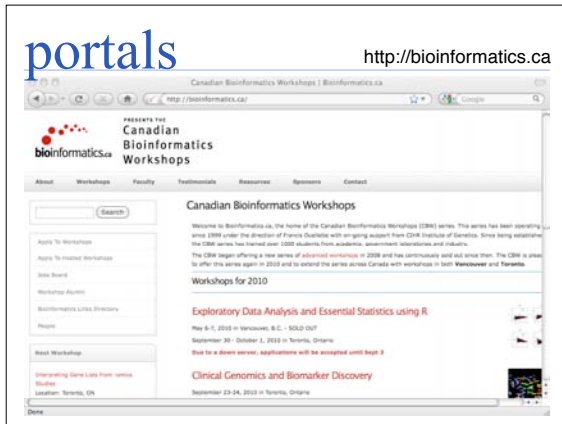
# Bioinformatics

Sources of information

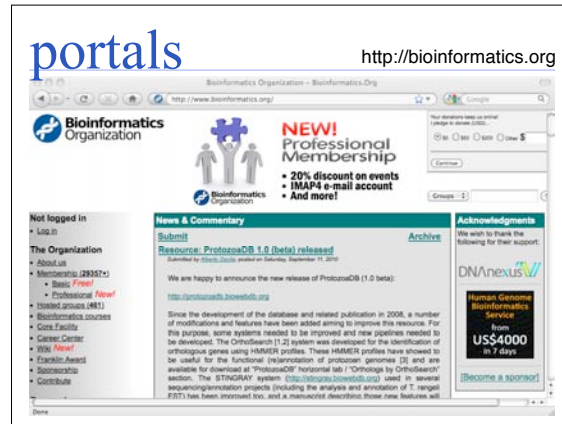Problem: outdated information has much inertia

# web

# Wikipedia    http://en.wikipedia.org/wiki/Bioinformatics

# Wikipedia    http://en.wikipedia.org/wiki/Bioinformatics

- Excellent introductory resource for many topics.
- Often written be very knowledgeable experts.
- Always be wary of inaccurate, incomplete information.
- You can immediately improve the text and correct errors.

# portals

http://bioinformatics.ca



# portals

http://bioinformatics.org



# portals

http://gchelpdesk.ualberta.ca



# societies

http://iscb.org



# journals

**Bioinformatics**
**NAR** (esp. Databases and WebServices issues)
**BMC Bioinformatics**
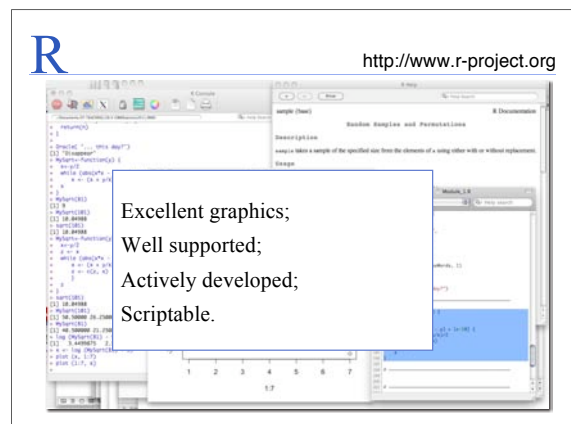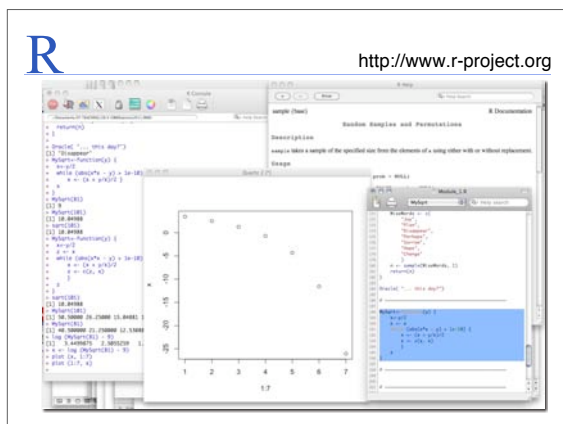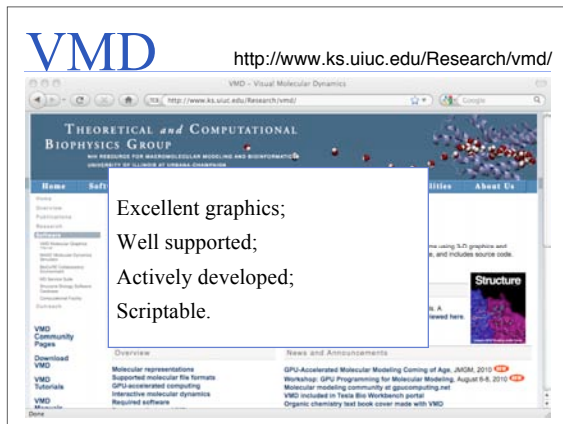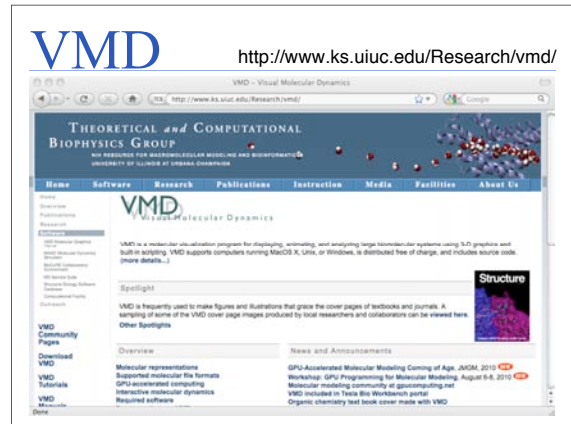**PLoS Computational Biology**
others ...

... all available electronically via U of T Library
... all have e-mail contents alert service or RSS.

# textbook

Zvelebil & Baum:
Understanding Bioinformatics

Garland Science, 2008

# Tools

- VMD
- R
- Perl
- Jalview
- Phylip
- UCSD genome browser
- [...]

---

# VMD

---

# VMD

Excellent graphics;

Well supported;

Actively developed;

Scriptable.



---

# VMD

Stereo vision ...



---

# R

---

# R

Excellent graphics;

Well supported;

Actively developed;

Scriptable.

# Perl

http://www.perl.org

**Perl** is a programming language.

**perl** is actually a program that runs commands in the Perl programming language. But from a user's perspective, that really doesn't make a difference.

**Perl** is

- free-format (whitespace is optional)
- compiled (everything is looked at before its executed)
- interpreted (works from code, step by step)

with

- automatic typing and memory management.

---

# Perl

http://www.perl.org

What is **Perl** good at?

Text processing
Rapid prototyping
Easy to learn for easy tasks
Powerful enough for difficult tasks
Programming for the Web

Use of large libraries of useful code modules

"Magic"
"There's more than one way to do it".

---

# Perl

http://www.perl.org

What is **Perl** poor at?

Complex, long-lived software projects with multiple authors
Need for complex datastructures
Performance-critical applications

"Magic"
"There's more than one way to do it".

---

http://biochemistry.utoronto.ca/
undergraduates/courses/
BCH441H/wiki

boris.steipe@utoronto.ca