# Reverse-engineering the human genome

**Mathieu Blanchette**
**McGill School of Computer Science**
**McGill Centre for Bioinformatics**

---

## Reverse engineering

**ie.exe**

010001011101010101010101
010101010101010101010101
010111110001010010101000
101001011101010010101010
110110101010001010110

?

**ie.cpp**

```
if (!strcmp(language,"sun.java"))
{
    printf("Unrecognized format");
    crashComputer();
}
```
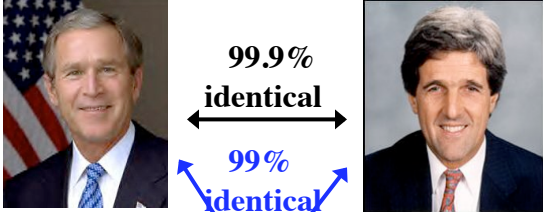
Goals: Understand...

• The function of each part of the code

• The interactions between different parts of the code

Motivations:

• Understand how a given problem is solved
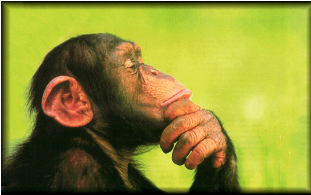
• Modify the code for our own purposes

# Genomes

THE STRUCTURE OF DNA

- Human genome: $\{A,C,G,T\}^{3 \times 10^9}$
- Each of your $10^{14}$ cells has two copies

**99.9% identical**

**99% identical**

one helical turn = 3.4 nm

Sugar-phosphate backbone

Base
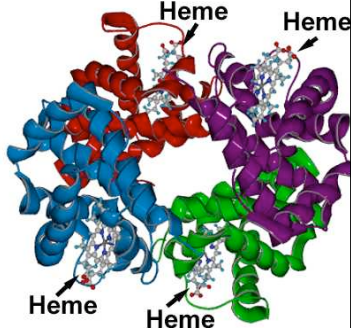
---

# Roles of the genome

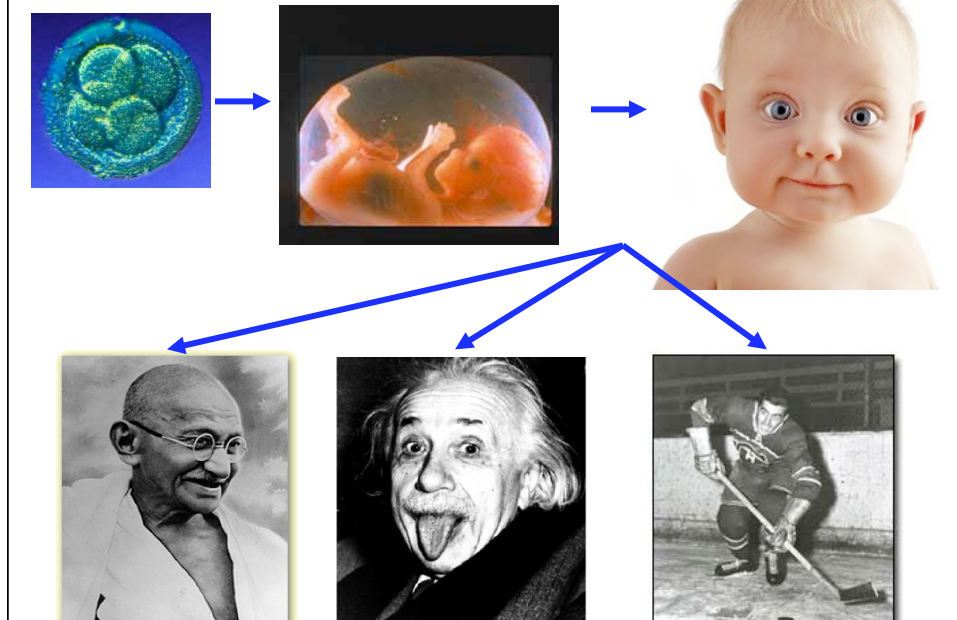Heme    Heme

- Genome is a blue print for a cell
- Describes *how* to build proteins
  - 25,000 genes --> 25,000 proteins (+variations)
  - Each protein has its biochemical function
- Describes *when* to build each protein
  - Under which situations should a gene be expressed?
- Proteins allow:
  - Cell administration and maintenance
  - Reaction to stimuli
  - Protocols for communication between cells
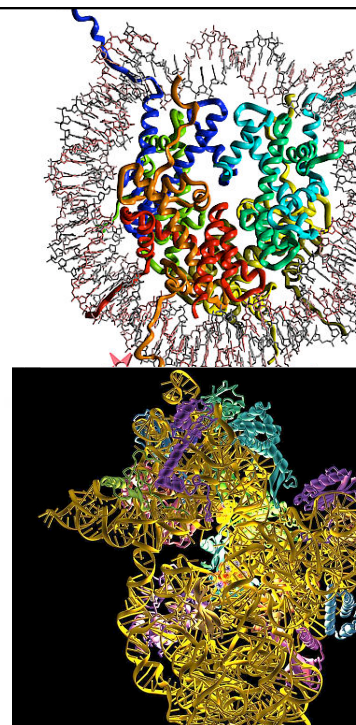
Heme    Heme

# Roles of the genome: development



# The hardware - Proteins

- Molecules only obey laws of physics and chemistry.
- Cell organization only relies on interactions between molecules
- Stochastic, dynamic, "chaotic" system
- High error rate in interactions
- Replicating, self-assembling, self-repairing system

Maybe software engineers have something to learn here...

# Content of the genome - Genes

• Gene: region of DNA that encodes one protein

DNA sequences                    Protein sequence

$\{A,C,G,T\}^L$

**START**    **STOP**

Second Letter of Codon

First letter of Codon (5' end)

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

$\{A,C,D,E,F,G,H,I,K,L,$
$M,N,P,Q,R,S,T,V,W,Y\}^K$

---

BLAST    PubMed    Nucleotide    Protein    Genome    Structure    PopSet    Taxonomy    Help

OVERVIEW

Views

Graphical View

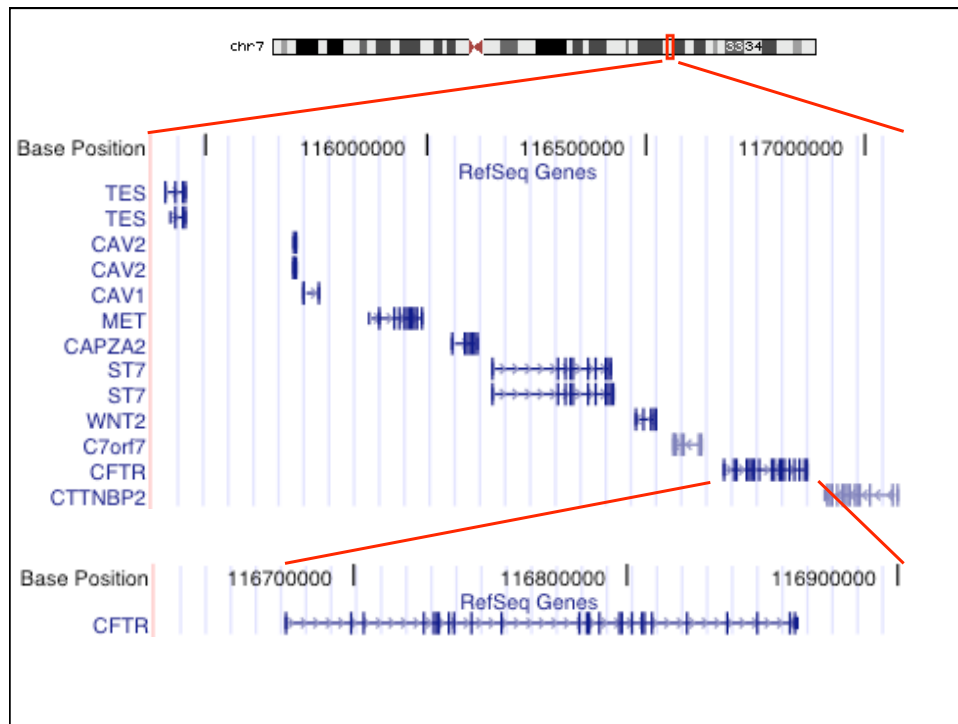**Haemophilus influenzae Rd KW20, complete genome - 75651..125650**

Start from : [ ] Go     Search for gene [ ] Find

43 protein coding genes    Find Open Reading Frames

Click on the rectangle to get BLAST neighbors for the gene of interest or click on the overview below to see a distant region

Haemophilus influenzae

Genome: 1.8 Mb

Number of genes: 1700

Translation, ribosomal structure and biogenesis
Transcription
DNA replication, recombination and repair
Cell division and chromosome partitioning
Posttranslational modification, protein turnover
Cell envelope biogenesis, outer membrane
Cell motility and secretion
Inorganic ion transport and metabolism
Signal transduction mechanisms
Energy production and conversion
Carbohydrate transport and metabolism
Amino acid transport and metabolism
Nucleotide transport and metabolism
Coenzyme metabolism
Lipid metabolism
Secondary metabolites biosynthesis, transport an
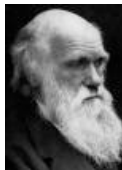General function prediction only
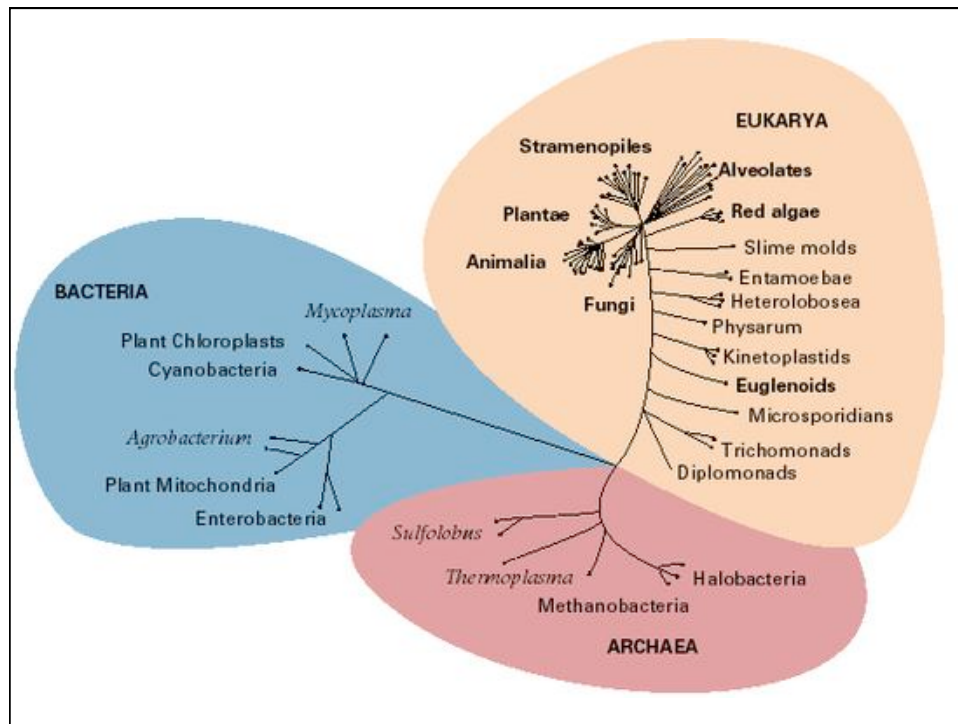Function unknown
No COG match

# The Programmer - Evolution

- Design principles:
  - Random modifications (variation)
  - Survival of the fittest (natural selection)
- 3 Billion years of evolution
- Today's species are the current solution of the fitness optimization problem

# Central dogma of comparative genomics

...ACTAGTCGCATACGATCAGCA...

Unavailable!

S1: GTTACGGTCACATACTGAAACA
S2: GTTATGGTCACATACTGAAACTGA
S3: ATTACTCGCATACGGTCTAACA
S4: ATTTACTCGCATACGGTCTAGCAC

# Mammalian evolution

- Rapid radiation ~75 Myrs ago

- Many nearly independent phyla

- Many "noisy" copies of ancestor

- **Accurate reconstruction of ancestors may be feasible**

Margulies et al., PNAS 2005

## Ancestral mammalian genome reconstruction

[Miller, Haussler, Blanchette]

Base-by-base reconstruction of complete ancestral genomes
• Including coding, non-coding, repetitive regions

Boreoeutherian ancestor

Expected reconstruction accuracy[*]:

• From ideal choice of extant mammals  99%

• From soon-to-be available genomes:  96%
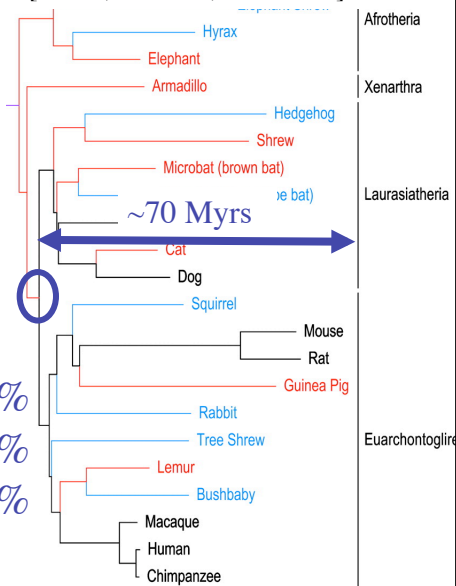
• With currently available sequences:  90%
  (full or 2X coverage)

[*] For >90% of euchromatic genome

~70 Myrs

Afrotheria — Hyrax, Elephant
Xenarthra — Armadillo
Laurasiatheria — Hedgehog, Shrew, Microbat (brown bat), (...e bat), Cat, Dog
Euarchontoglire — Squirrel, Mouse, Rat, Guinea Pig, Rabbit, Tree Shrew, Lemur, Bushbaby, Macaque, Human, Chimpanzee

Tree from Margulies et al., PNAS 200

---

# Why should we care?



• See ... st see its ...

• Bo ... an stuff

# One program, many functions

Neuron (nerve cell)

Muscle cell

**Same
genome!**

# Regulation of gene expression

Transcription factor binding sites:
- ☹ Short: 6 to 20 nucleotides
- ☹ No specific signature; each TF has different binding site
- ☹ Can be up to 1 million nucleotides upstream of gene regulated
- ☺ Often clustered with other binding sites, forming modules