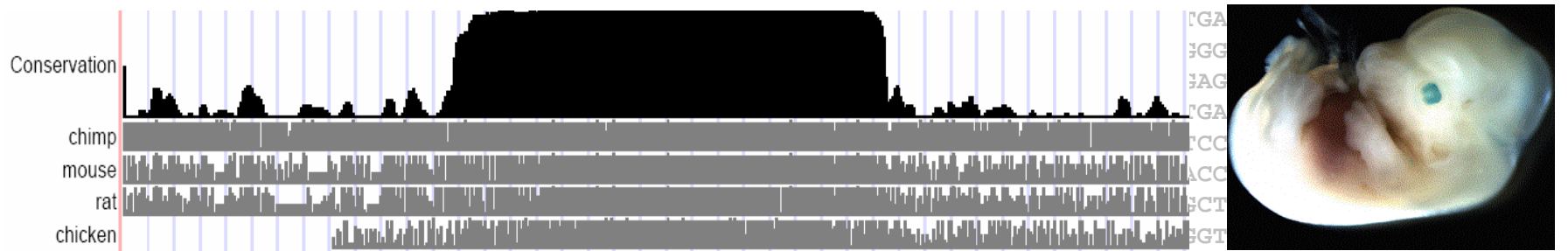


GGTGCCAGGGAAAGGGCAGGAGGTGAGTGCTGGGAGGCAGCTGAGGTCAACTTCTTTTGAACCTCCACGTGGTATTTACTCAGAGCAATTGGTGCCAGAG  
GCTCAGGGCCCTGGAGTATAAAGCAGAATGTCTGCTCTCTGTGCCAGACGTGAGCAGGTGAGCAGCTGGGGCGAAAAGACCTGTTGGAGGCTATGAATGC  
AATCAAGGTGACG...  
CGTGTGTATGAG...  
TTCCTCATGGGGCAAATCTCAC...  
CCTTGGCTG...  
AGGGG...  
ACAAACAGGTTCTTCTCTGTGGTGGGCCAGCCAGCAGGTGAGTGGGAAGGTTAAAGGTGATGGGGTTGGGAGAACTGGGTGAGGAGTTCAGCCCCATC  
CCCCGTAAGCTCCTGGGAAGCACTTCTCTACTGGGGCAGCCCCTGATACCAGGGCACTCATTAAACCCTCTGGGTGCCAGGGAAAGGGCAGGAGGTGAGT  
GCTGGGAGGCAGCTGAGGTCAACTTCTTTTGAACCTCCACGTGGTATTTACTCAGAGCAATTGGTGCCAGAGGCTCAGGGCCCTGGAGTATAAAGCAGAA

# Deciphering the Human Genome: Computational Insights & Opportunities



GAAACTGGGTGAGGAGTTCAGCCCCATCCCCGTAAGCTCCTGGGAAGCACTTCTCTACTGGGGCAGCCCCTGATACCAGGGCACTCATTAAACCCTCTG  
GGTGCCAGGGAAAGGGCAGGAGGTGAGTGCTGGGAGGCAGCTGAGGTCAACTTCTTTTGAACCTCCACGTGGTATTTACTCAGAGCAATTGGTGCCAGAG  
GCTCAGGGCCCTGGAGTATAAAGCAGAATGTCTGCTCTCTGTGCCAGACGTGAGCAGGTGAGCAGCTGGGGCGAAAAGACCTGTTGGAGGCTATGAATGC  
AATCAAGGTGACG...  
CGTGTGTATGAG...  
TTCCTCATGGGGCAAATCTCAC...  
CCTTGGCTG...  
AGGGG...  
ACAAACAGGTTCTTCTCTGTGGTGGGCCAGCCAGCAGGTGAGTGGGAAGGTTAAAGGTGATGGGGTTGGGAGAACTGGGTGAGGAGTTCAGCCCCATC  
CCCCGTAAGCTCCTGGGAAGCACTTCTCTACTGGGGCAGCCCCTGATACCAGGGCACTCATTAAACCCTCTGGGTGCCAGGGAAAGGGCAGGAGGTGAGT  
GCTGGGAGGCAGCTGAGGTCAACTTCTTTTGAACCTCCACGTGGTATTTACTCAGAGCAATTGGTGCCAGAGGCTCAGGGCCCTGGAGTATAAAGCAGAA  
TGTCTGCTCTCTGTGCCAGACGTGAGCAGGTGAGCAGCTGGGGCTGTCTGCTCTCTGTGCCAGACGTGAGCAGGTGAGCAGCTGGGGCTGTCTGCTCT  
CTGTGCCAGACGTGAGCAGGTGAGCAGCTGGGGCTGTCTGCTCTCTGTGCCAGACGTGAGCAGGTGAGCAGCTGGGGCTGTCTGCTCTCTGTGCCAG

**Gill Bejerano**

Assistant Professor

Dept. of Developmental Biology

& Dept. of Computer Science

Stanford University

2006

2007

Postdoc w/David Haussler  
School of Engineering  
UC Santa Cruz



# This is “the Century of Biology”

---

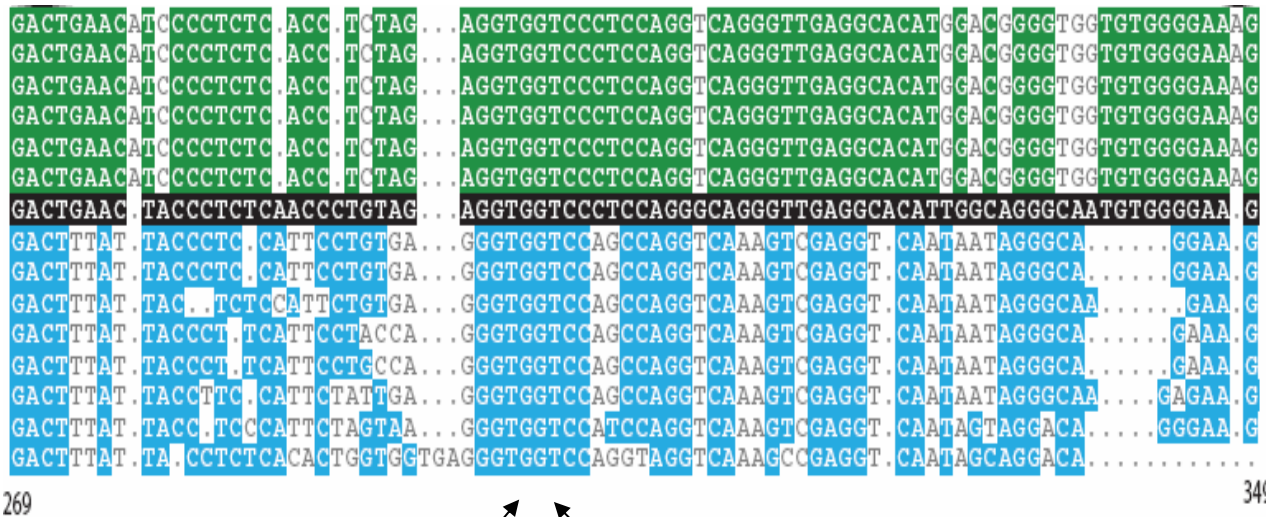
## EDITORIAL

### Unification in the Century of Biology

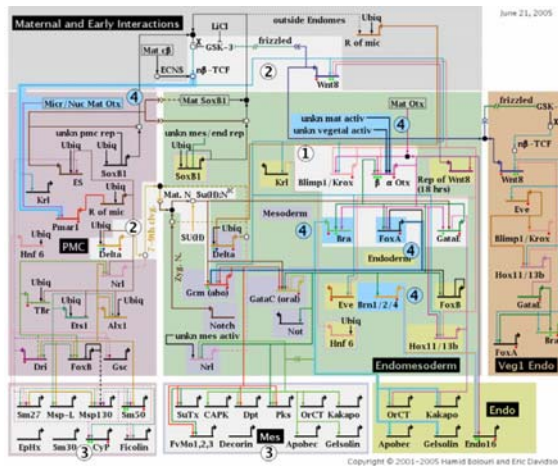
**S**cientific progress is based ultimately on unification rather than fragmentation of knowledge. At the threshold of what is widely regarded as the century of biology, the life sciences are undergoing a profound transformation. They have long existed as a collection of narrow, even parochial, disciplines with well-defined territories. Now they are undergoing consolidation, forming two major domains: one extending from the molecule to the organism, the other bringing together population biology, biodiversity studies, and ecology. Kept separate, these domains, no matter how fruitful, cannot hope to deliver on the full

# We can now cast Biology in “our” terms

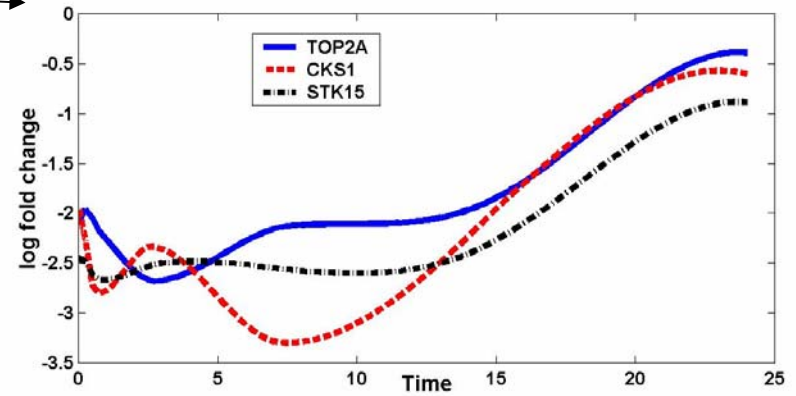
strings



circuits

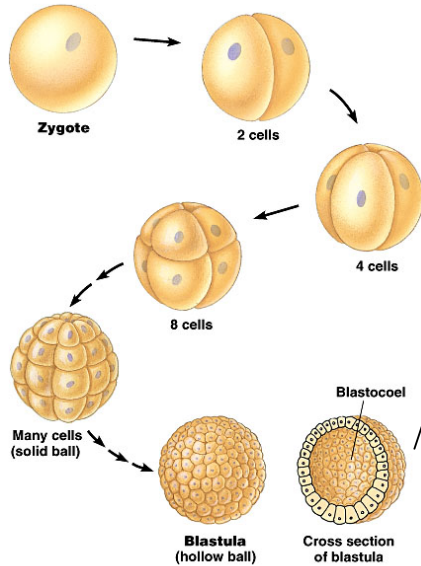


time series

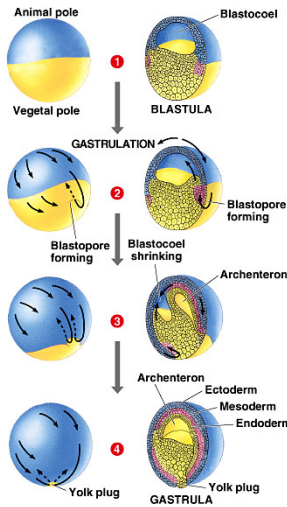


# Grand Challenge: Understanding Embryonic Development

one  
cell

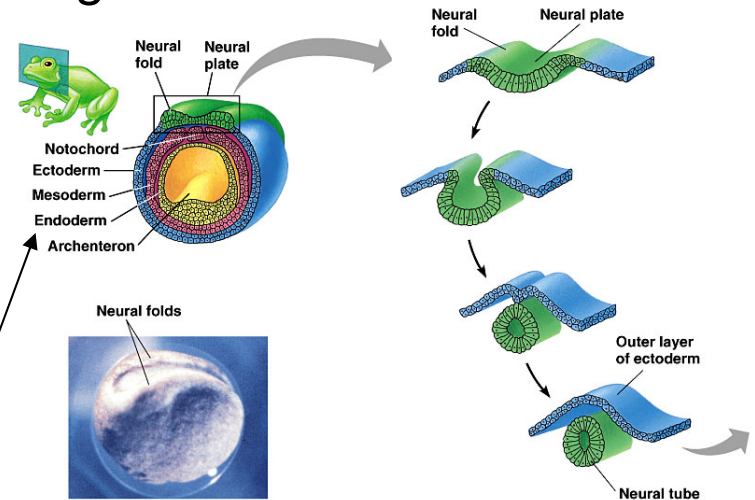


Copyright © 2003 Pearson Education, Inc., publishing as Benjamin Cummings.



Copyright © 2003 Pearson Education, Inc., publishing as Benjamin Cummings.

organism



Copyright © 2003 Pearson Education, Inc., publishing as Benjamin Cummings.

Enter DNA (“Merely the Secret of Life”) ...

# DNA: Functional and Non-Functional

---

DNA = linear molecule that carries instructions for making living organisms ~ long string(s) over a small alphabet

Alphabet of four {A,C,G,T}      Strings of length  $10^4$ - $10^{11}$

...ACGTACGACT**TGACTAGCATCGACTAC**GACTAGCAC...



“junk” DNA



**genetic  
instructions:**

**how to...**

**when to...**

**where to...**



“junk” DNA

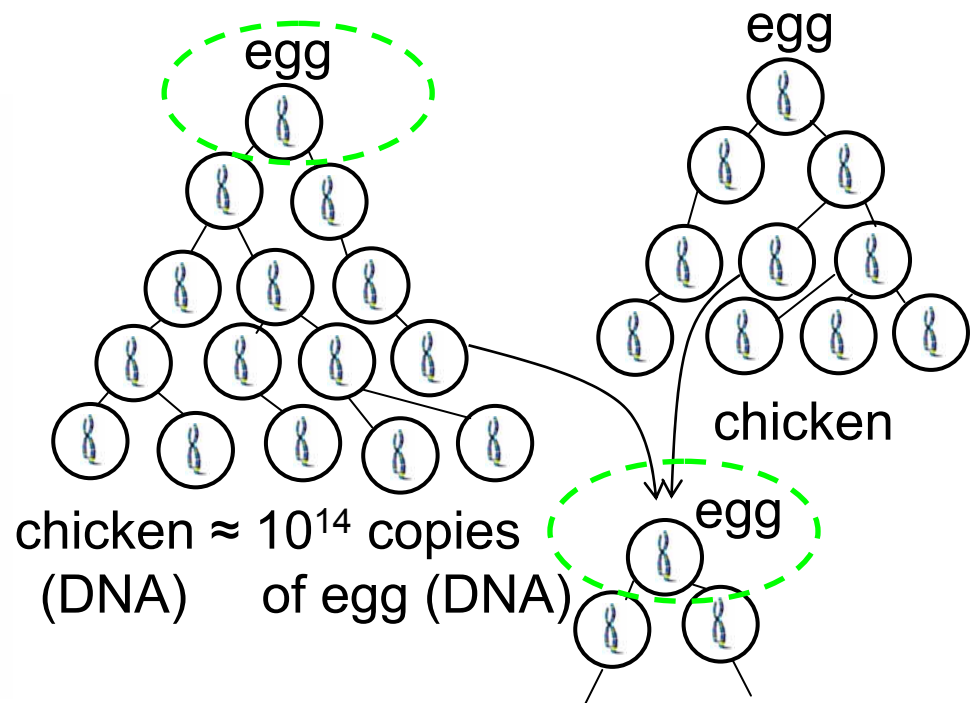
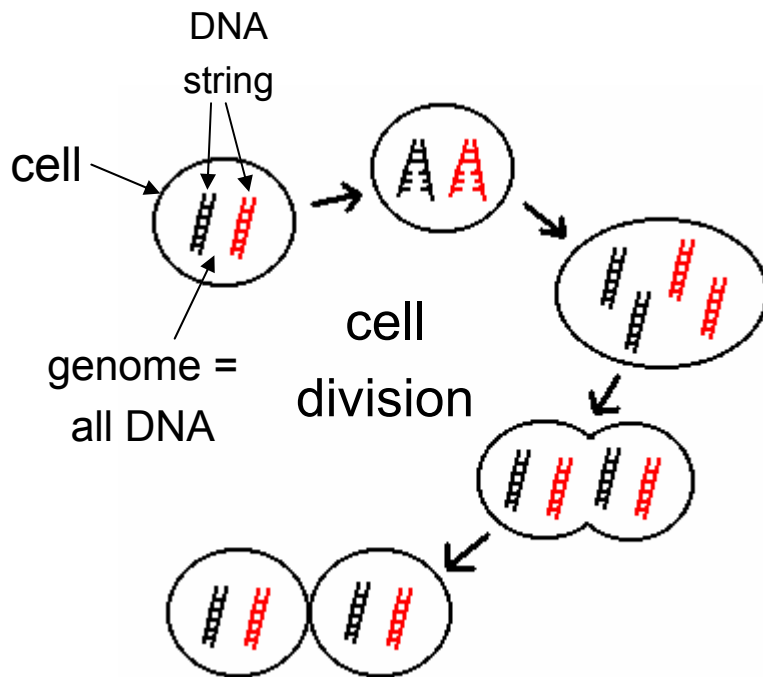
# One Cell, One Genome, One Replication

Every cell holds a single copy of all its DNA = its genome.

The genome is replicated every cell division.

The human body is made of  $\sim 10^{14}$  cells.

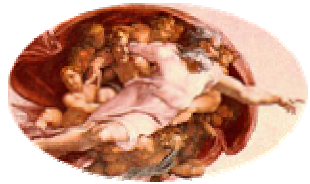
All originate from a *single* cell through cell division.



# Comparative Genomics

“Nothing in Biology Makes Sense  
Except in the Light of Evolution”

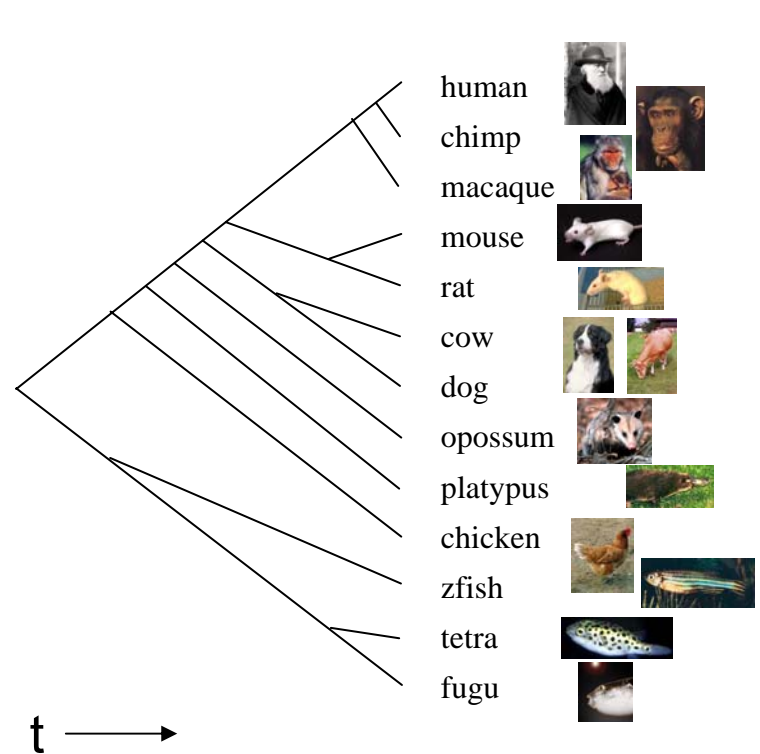
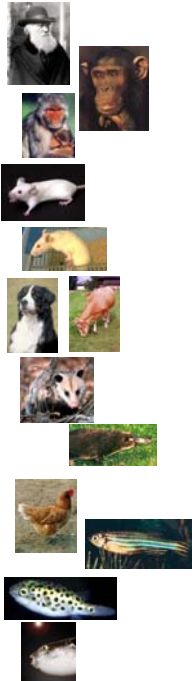
Theodosius Dobzhansky



Intelligent Designer

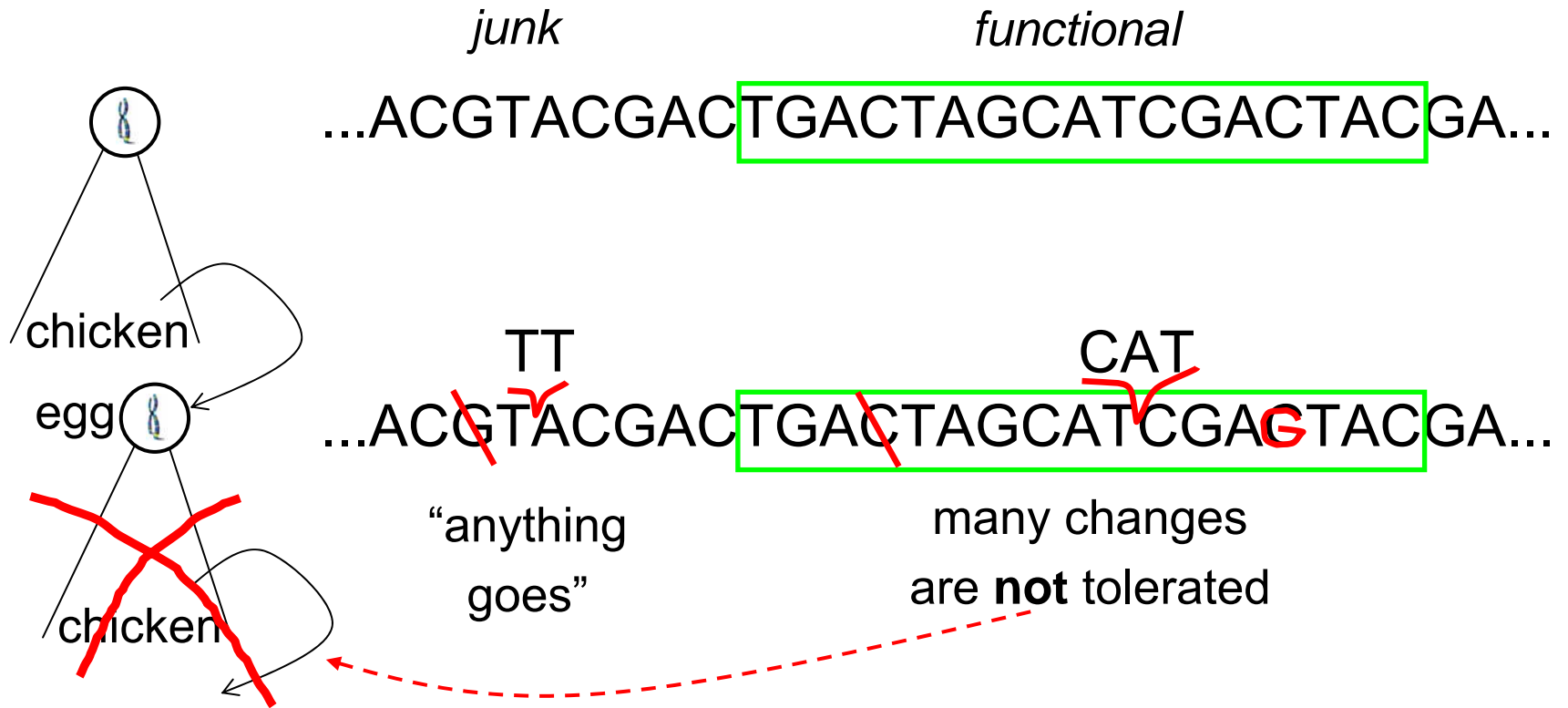


human  
chimp  
macaque  
mouse  
rat  
cow  
dog  
opossum  
platypus  
chicken  
zfish  
tetra  
fugu



# DNA Replication is Imperfect

Small Scale: single letters are substituted, erased, added



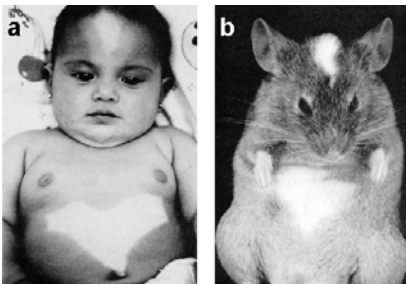
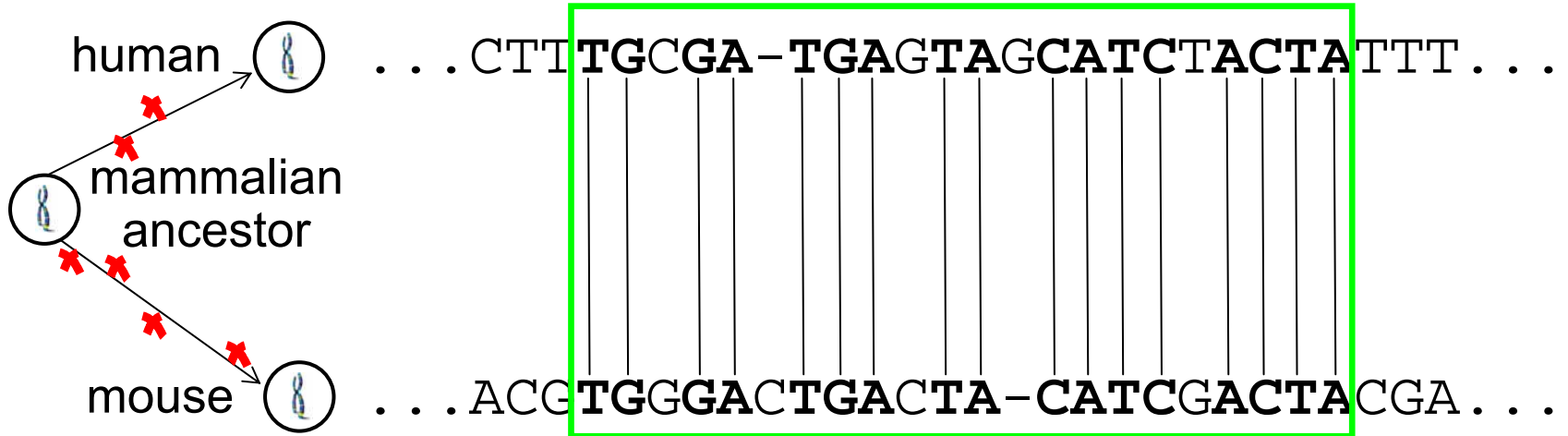
thus, sequence conservation over generations → function!



# Sequence Conservation implies Function

Comparative Genomics of Distantly related species:

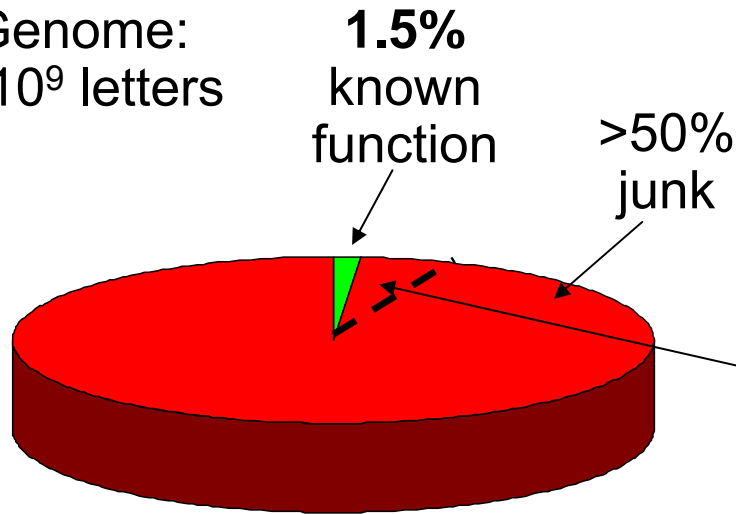
functional region!



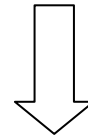
(but which function/s?...)

# The Human Genome is Full of Mysteries

Human  
Genome:  
 $3 \times 10^9$  letters



compare to other species



>5% human genome functional

**3x more functional DNA than known!**

**But what do these  $\sim 10^7$  substrings do??**



[*Science* 2004 Breakthrough of the Year, 5<sup>th</sup> runner up]

# Genes = How to make Proteins

 gene

DNA

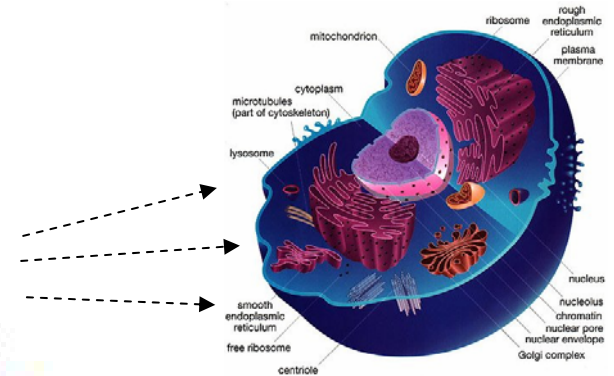


Translation

protein



“the workhorses of every living cell”



cell

# DNA Replication is Imperfect (contd)

Medium Scale: substrings are duplicated, deleted, inverted

Large Scale: whole DNA strings are duplicated, deleted

*junk*

*functional*

...ACGTACGACT **TGACTAGCATCGACTACGA**...

**substring  
duplication**

*functional*

*functional*

...ACGTACGACT **TGACTAGCATCGACTACGA**.....TCT**TGACTAGCATCGACTACGA**...

**functional  
divergence**

*functional'*

*functional''*

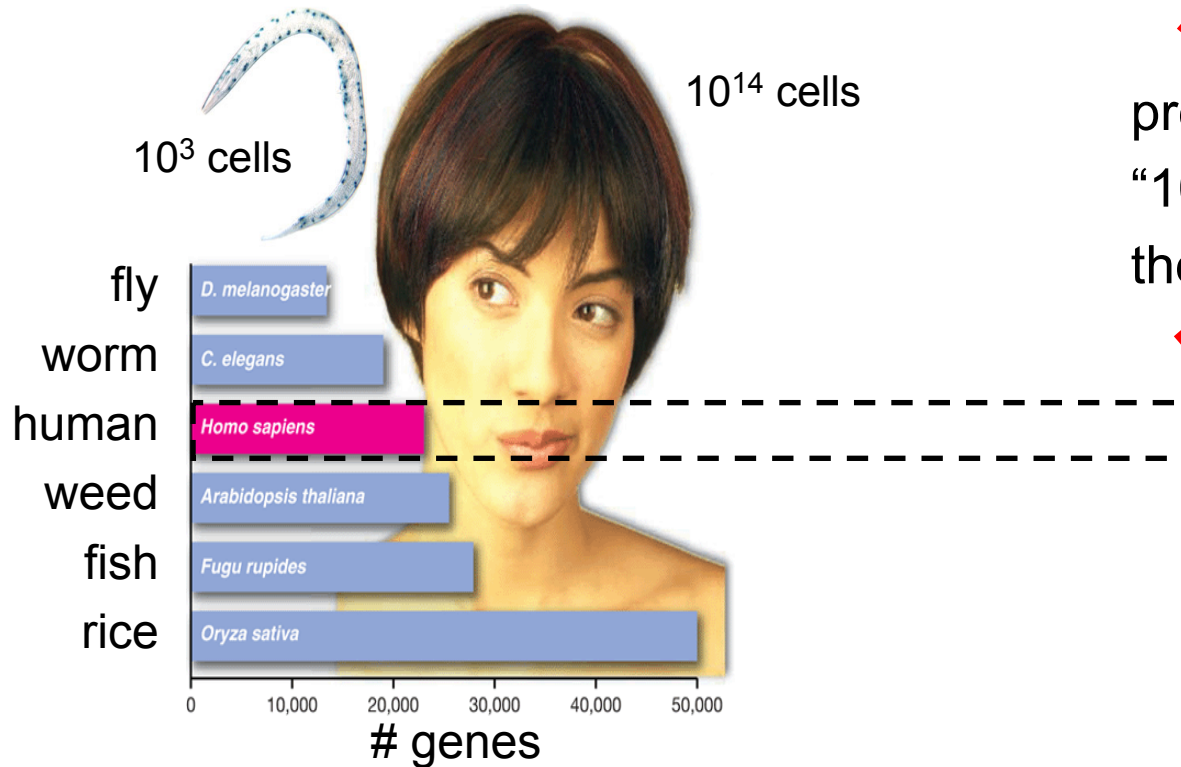
...ACGTACGACT **TGACTAGCATCGACTACGA**.....TCT**TGACTAGCATCGACTACGA**...

So...More Genes...More Complexity!...Right?

# Genes & Complexity

Gene numbers do **not** correlate with organism complexity.

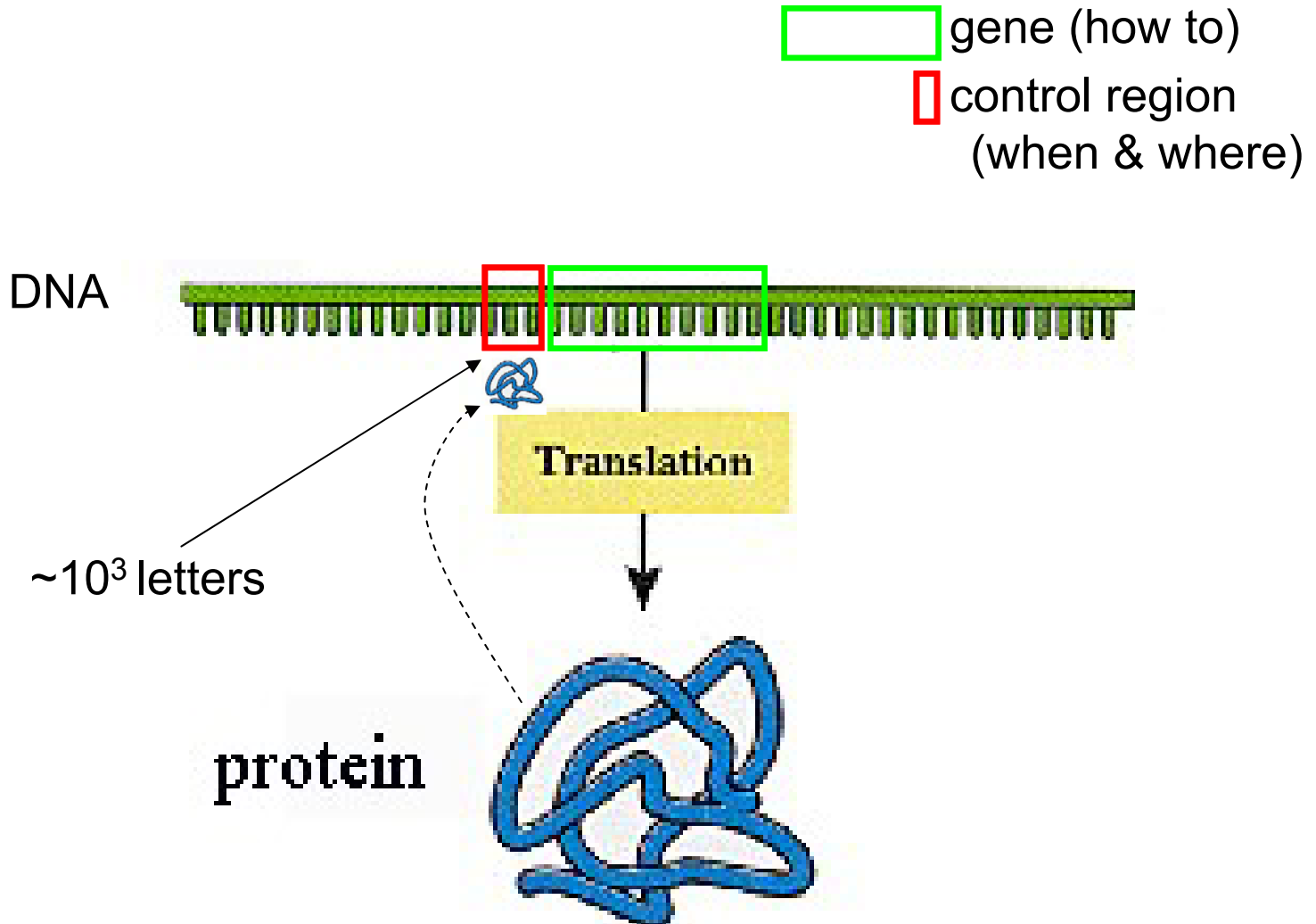
Many gene families are surprisingly old.



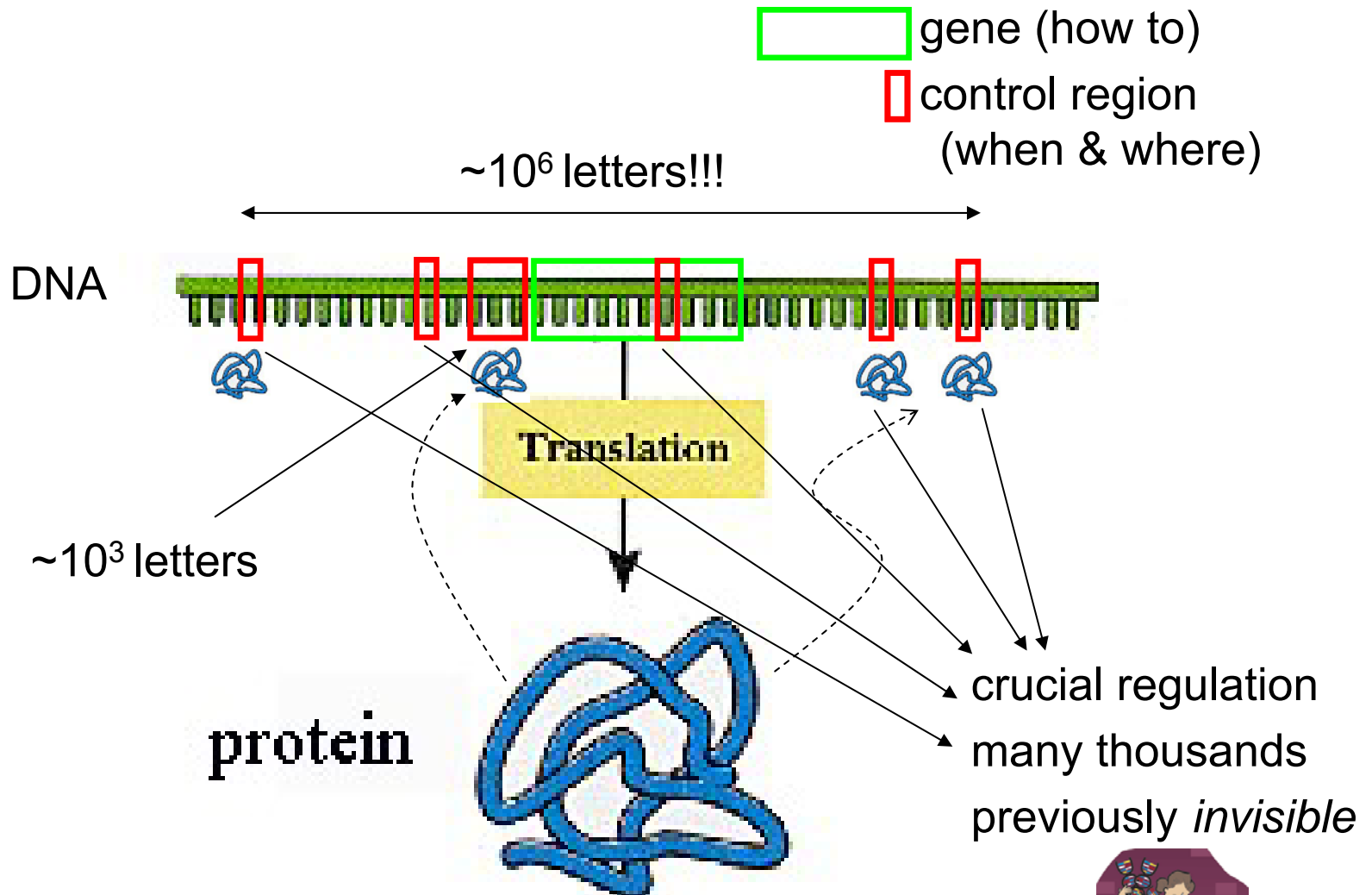
~~pre-genomic era:  
“100,000 genes to  
the human genome”~~

# Gene regulation = when/where to make protein

---



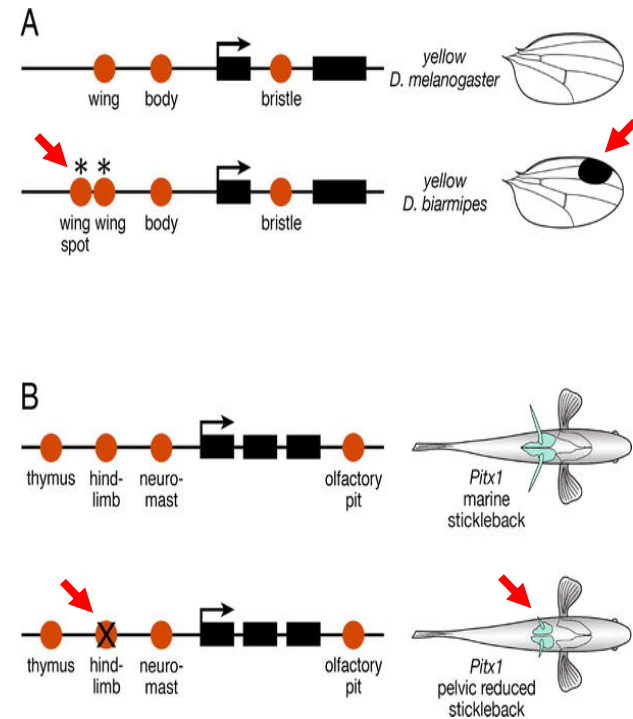
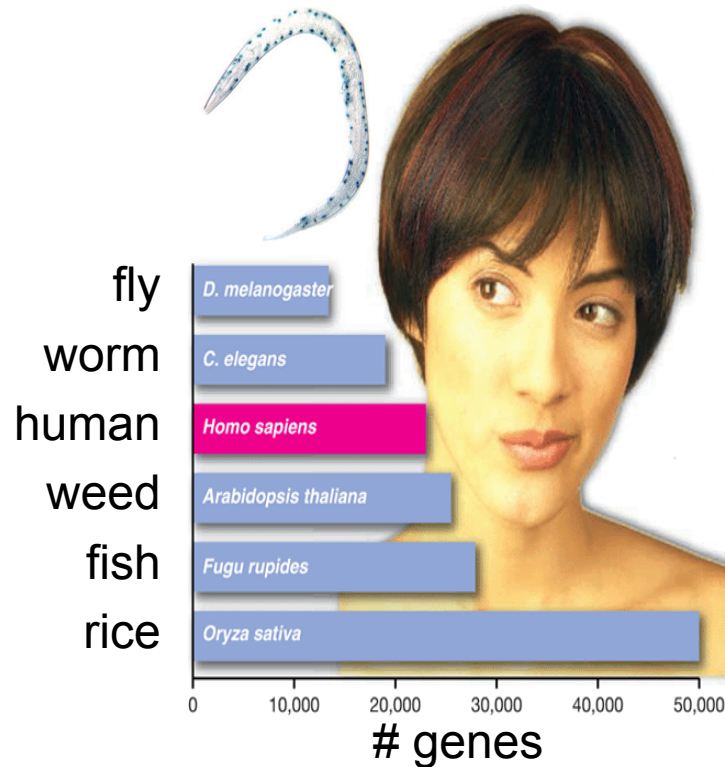
# Vertebrate Gene Regulation



# Regulatory regions drive morphological diversity

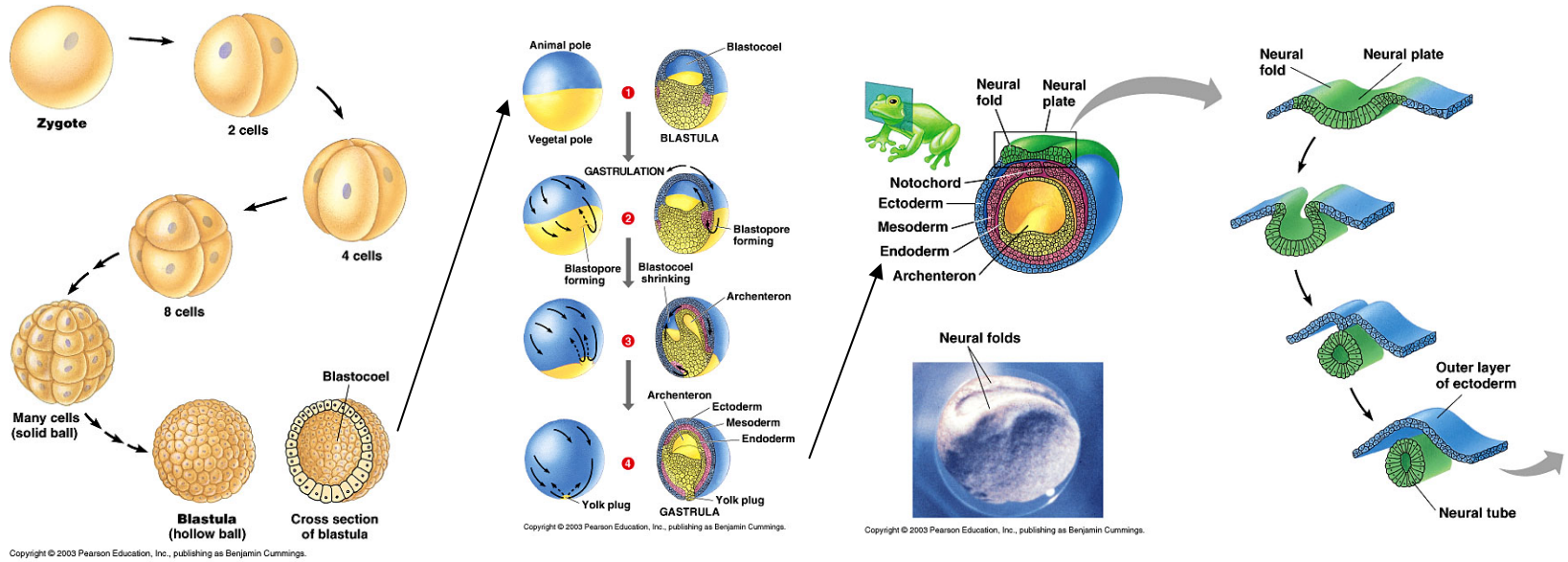
Gene numbers do **not** correlate with organism complexity.  
 Many gene families are surprisingly old.

“Regulatory sequence evolution must be the major contribution to the evolution of form.” [Carroll, Wilson memorial lecture, *PLoS Biol*, 2005]





# Embryonic Development and the Genome



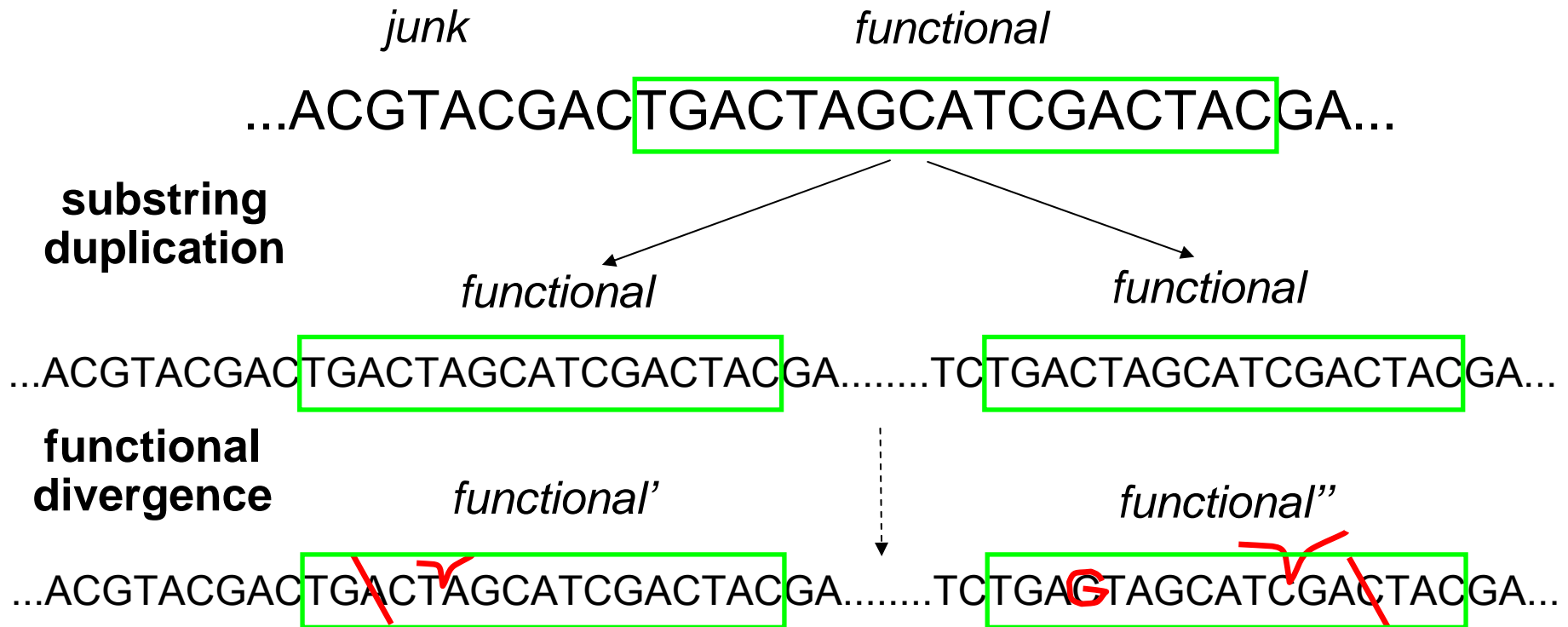
Many thousands of human conserved elements  
 congregate en-masse near developmental genes.  
 [Eg, Dog Genome Paper, Nature, 2005; Bejerano et al., Nature Methods, 2005]



# DNA Replication is Imperfect (reminder)

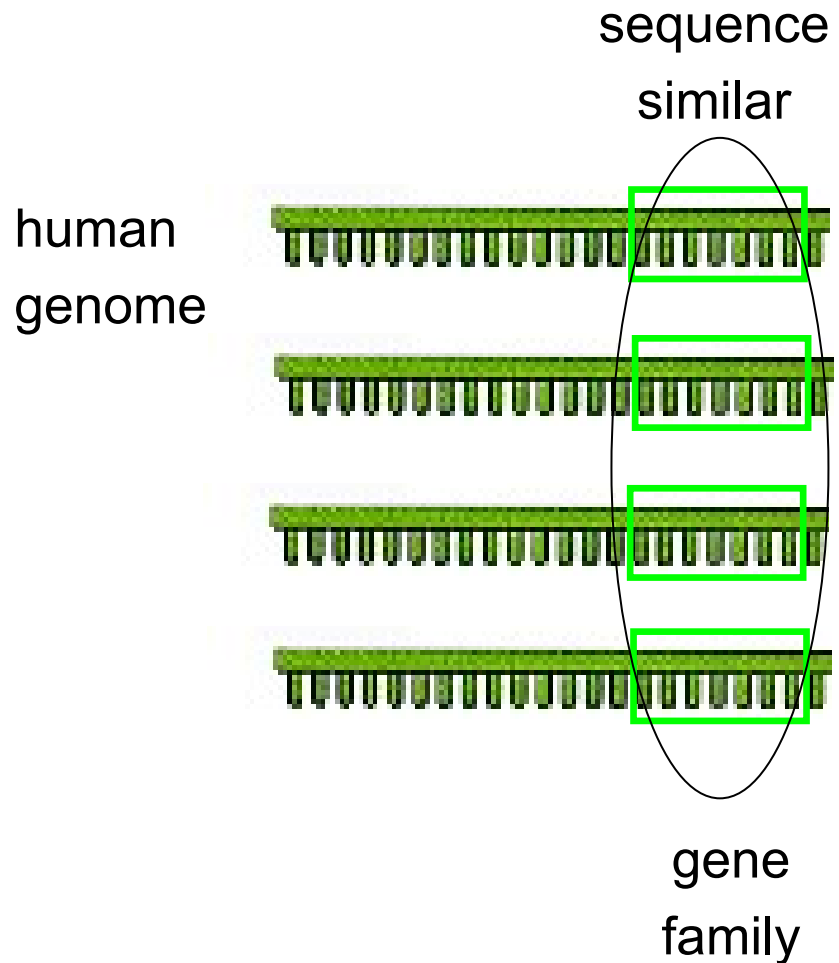
Medium Scale: substrings are duplicated, deleted, inverted

Large Scale: whole DNA strings are duplicated, deleted



# The Power of Computation

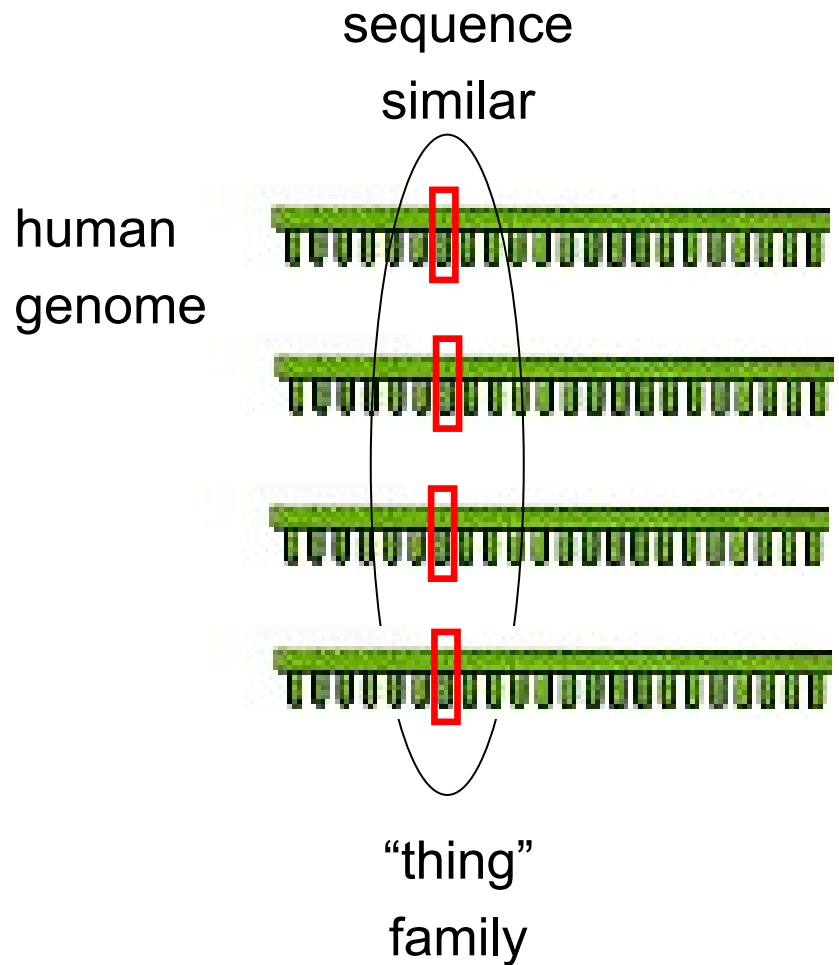
---



- similar function
- tale telling differences
- figure out one and you have a working hypothesis for all

# The Power of Computation

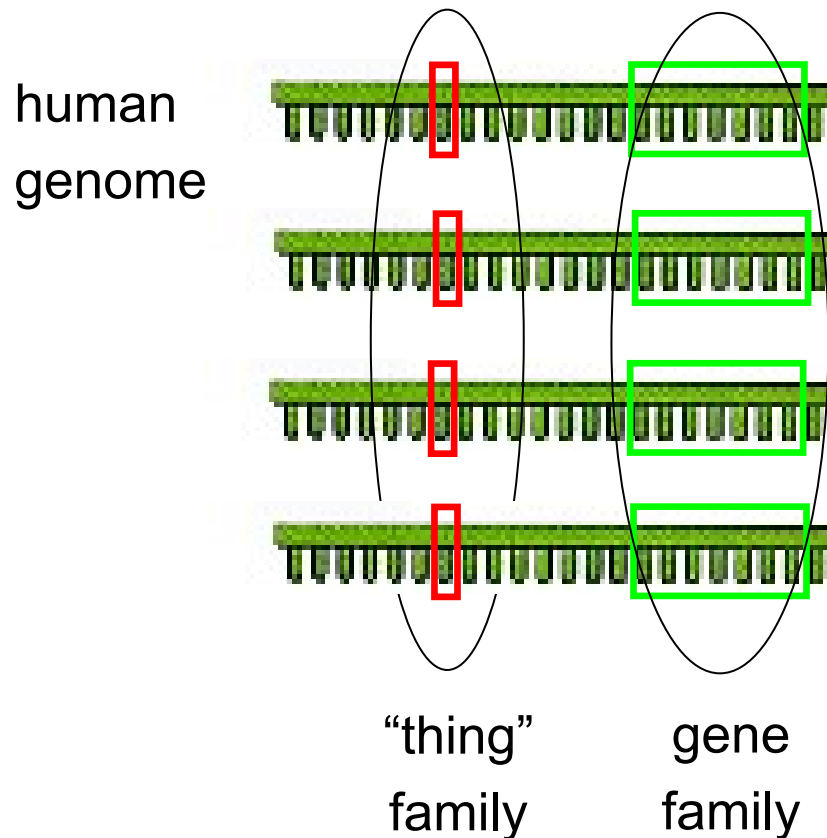
---



- similar function
- tale telling differences
- figure out one and you have a working hypothesis for all

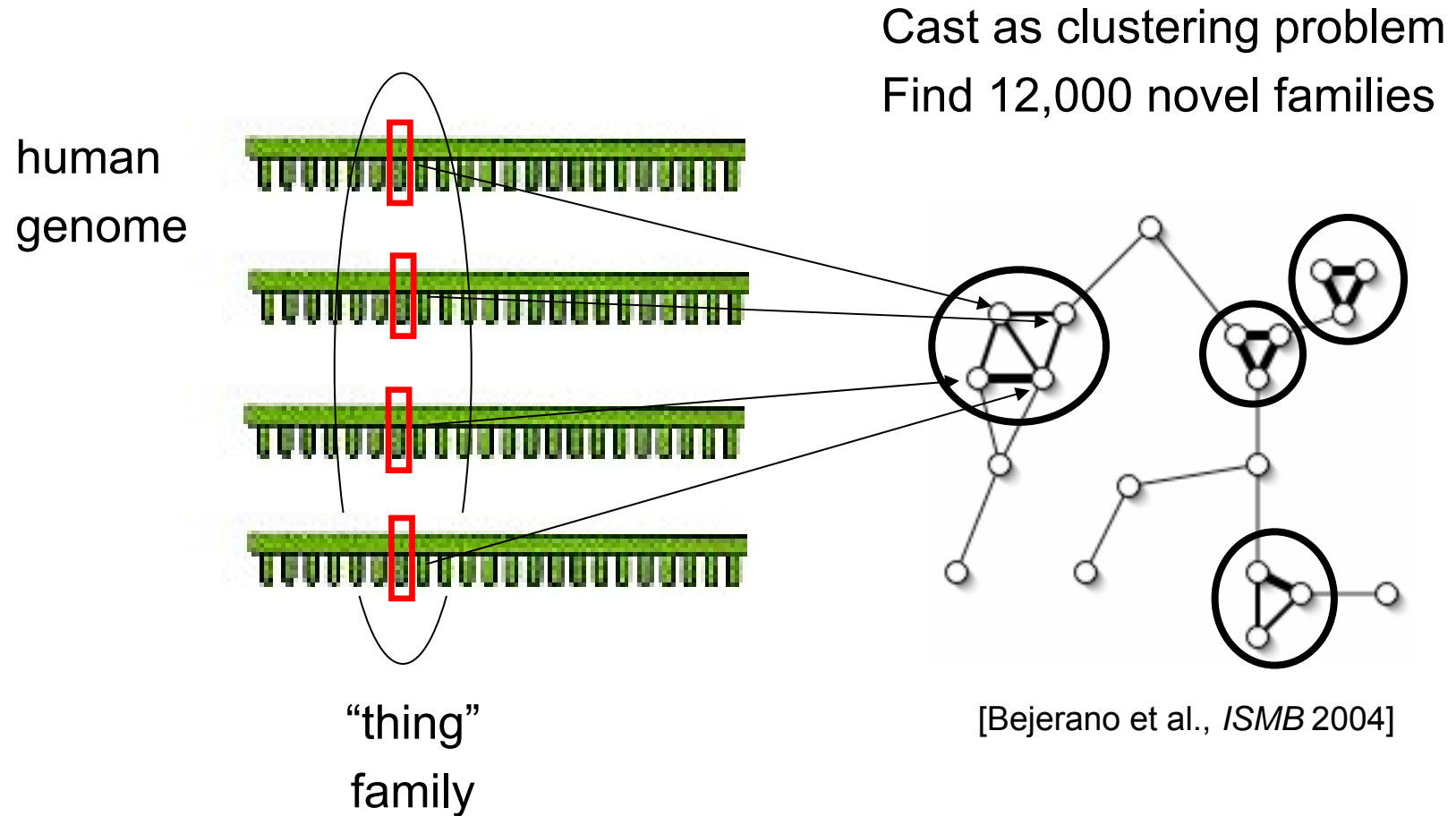
# The Power of Computation

---

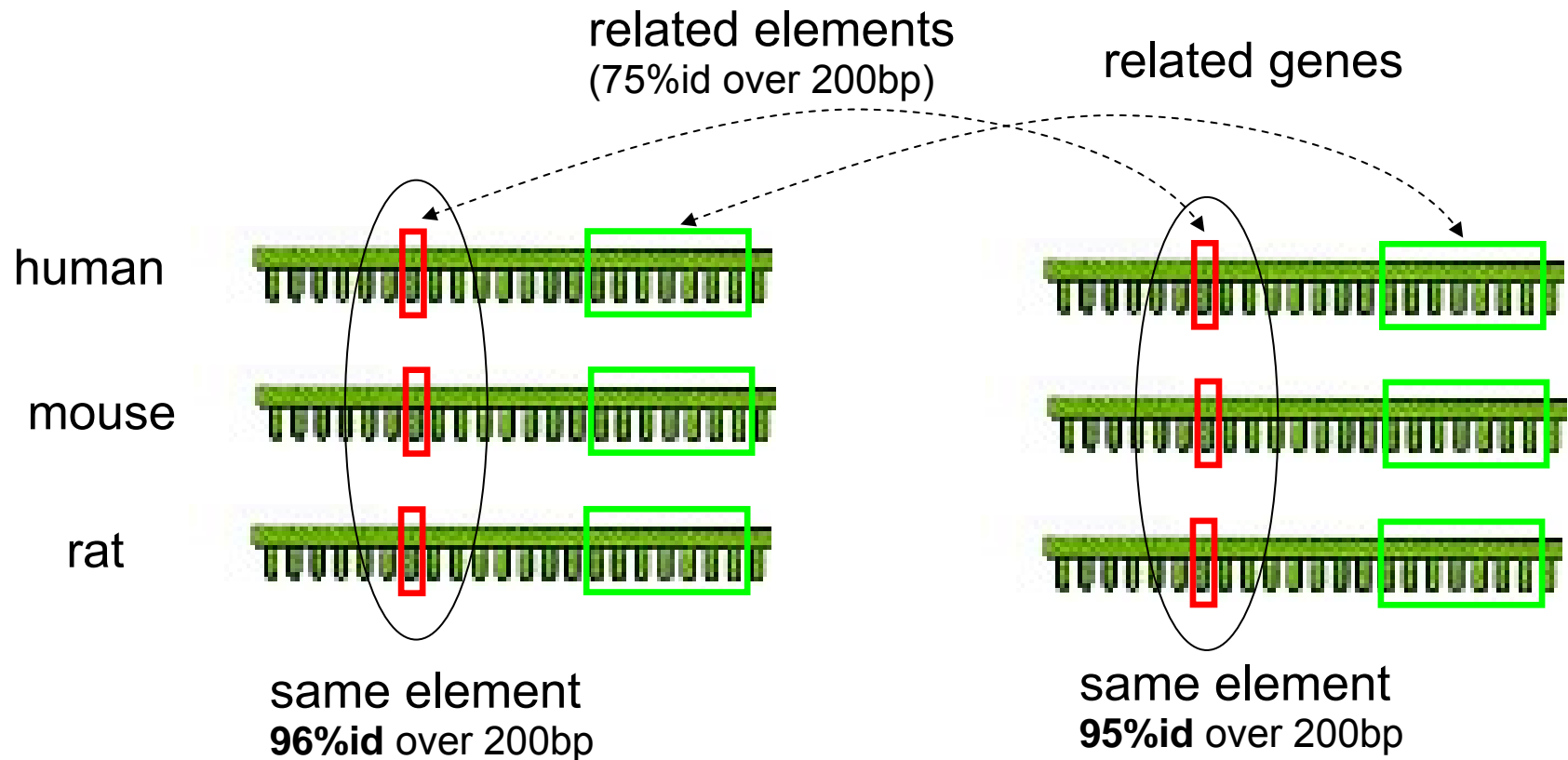


- similar function
- tale telling differences
- figure out one and you have a working hypothesis for all
- *“guilt by association” anchoring to genes (annotated landmarks)*

# Families of Conserved Non Coding Elements



# A Computational Question

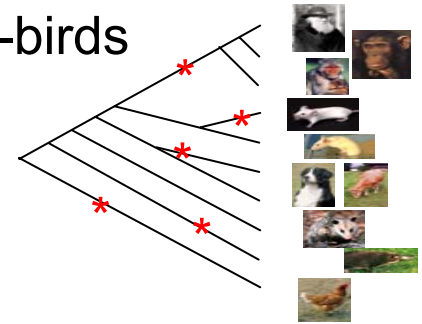


Classical Biological approach: experiment to understand *these* regions  
Computational approach: how many regions *like* this are there?

# Ultraconserved Elements

Hundreds of long substrings *identical* between human-birds  
 → they must have **rejected many different changes**.

But... *all* functions we understand in our genome are encoded using **redundant codes**.

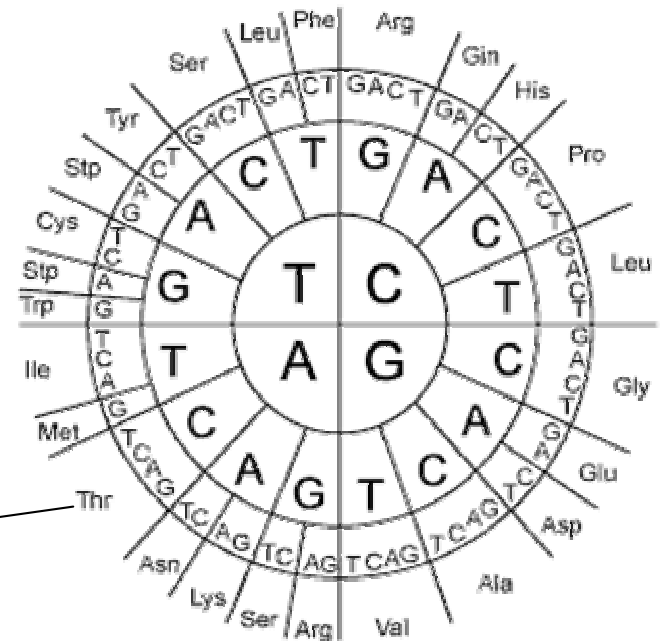
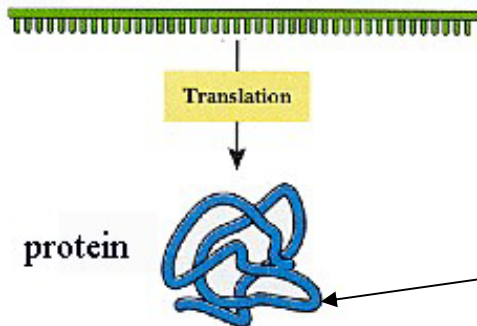


E.g. Protein Coding Genes:

DNA –  $10^8$  letters  
 over alphabet of 4.

Protein –  $10^2$  letters  
 over alphabet of 20.

Coding: 3 DNA letters → 1 Protein letter.

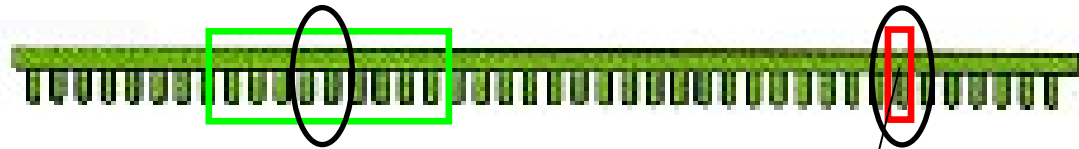




# Computational Hypotheses

Based on public domain genome wide data:

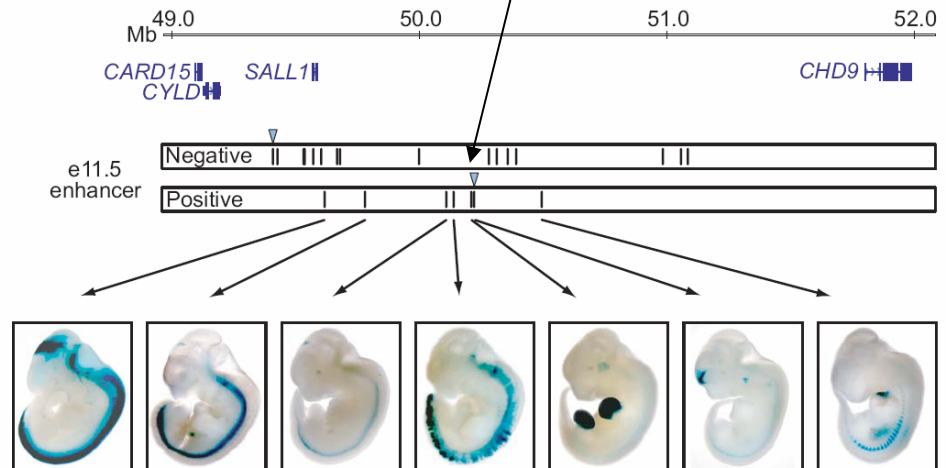
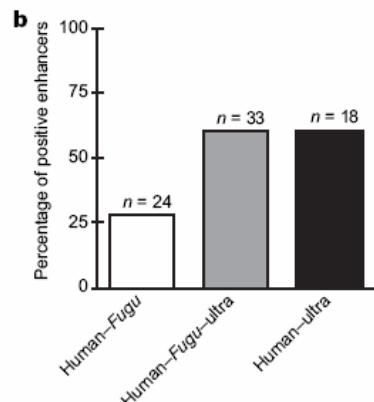
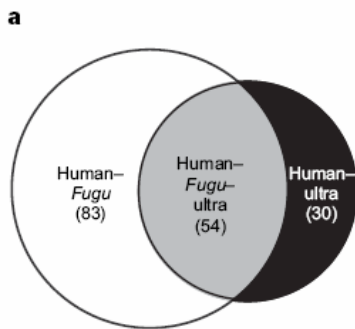
ultraconserved elements



one subset codes protein

larger subset does not

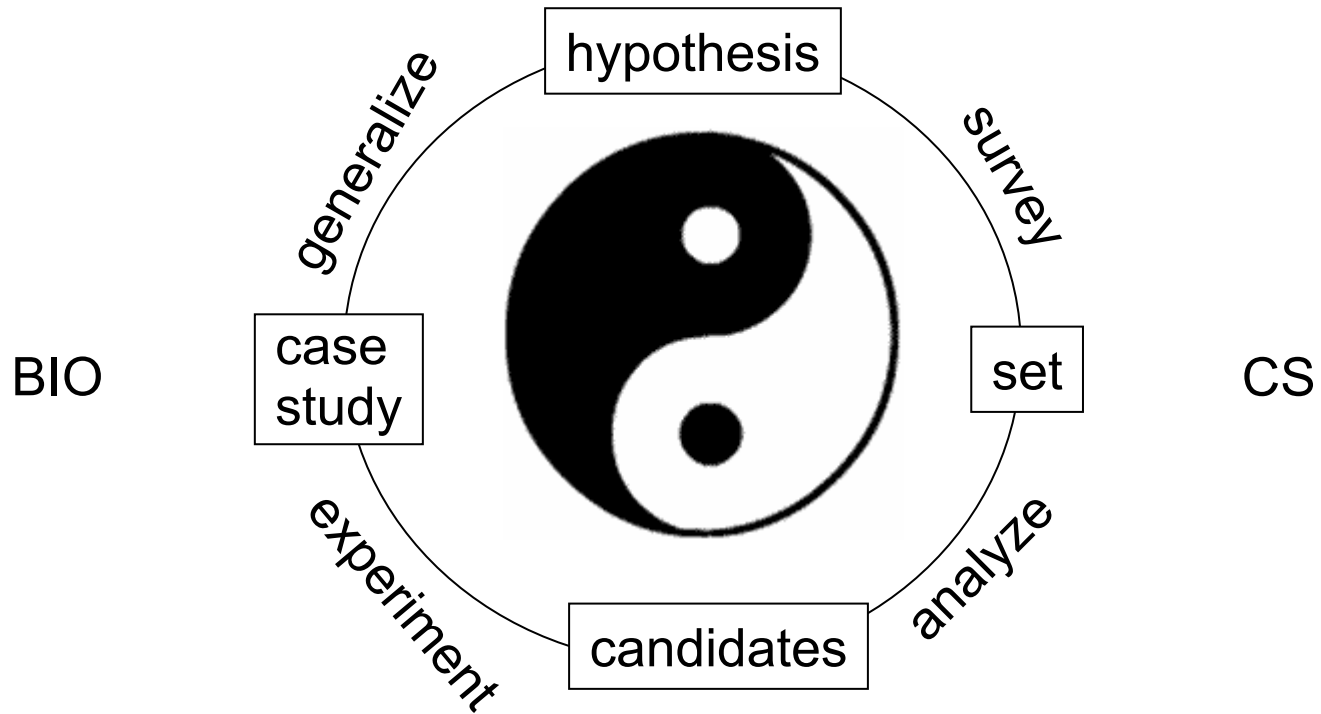
generate testable hypotheses for function from existing knowledge (2004)



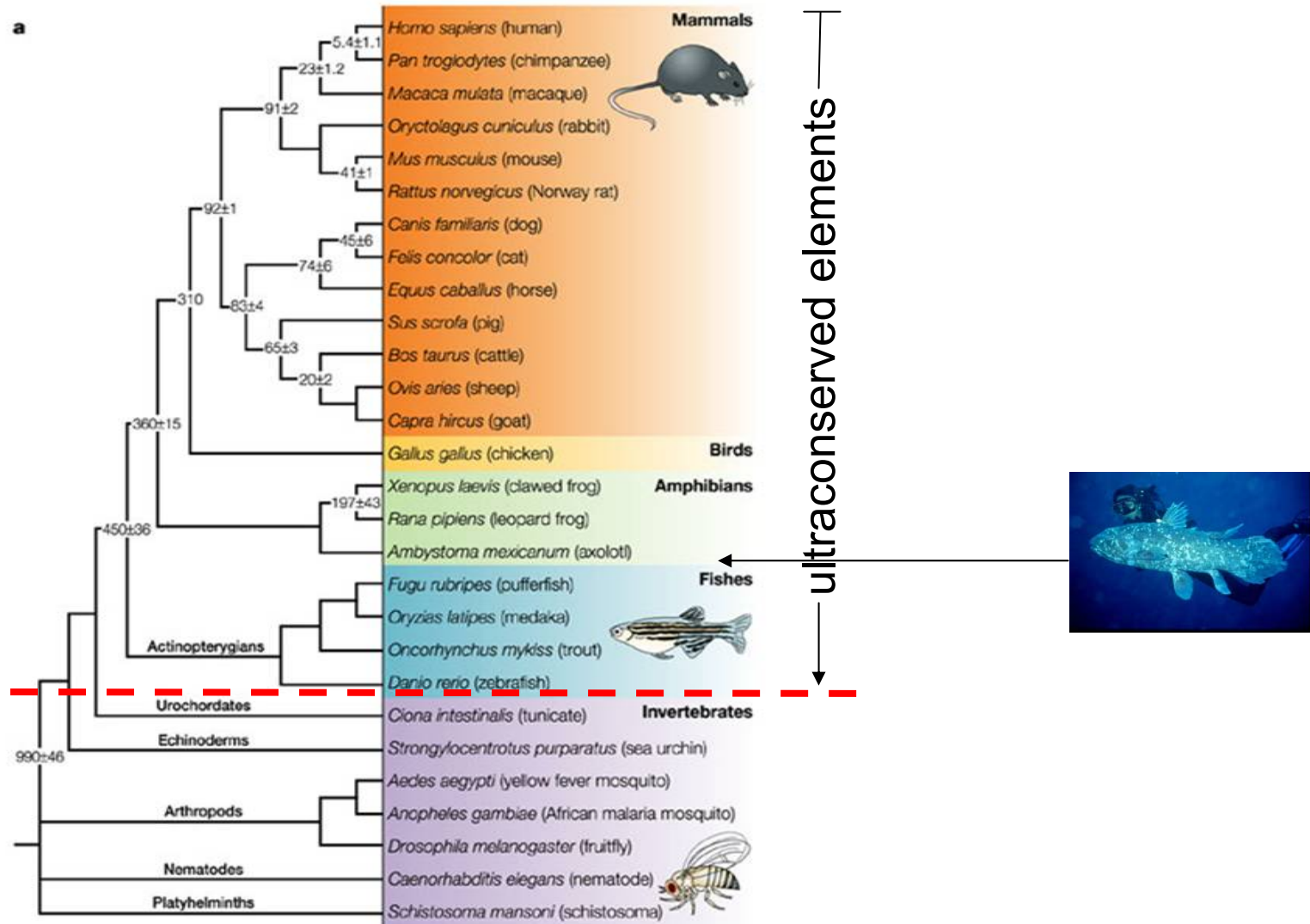
[Pennacchio et al., *Nature*, 2006]

# Computationally Driven Biology Simplified

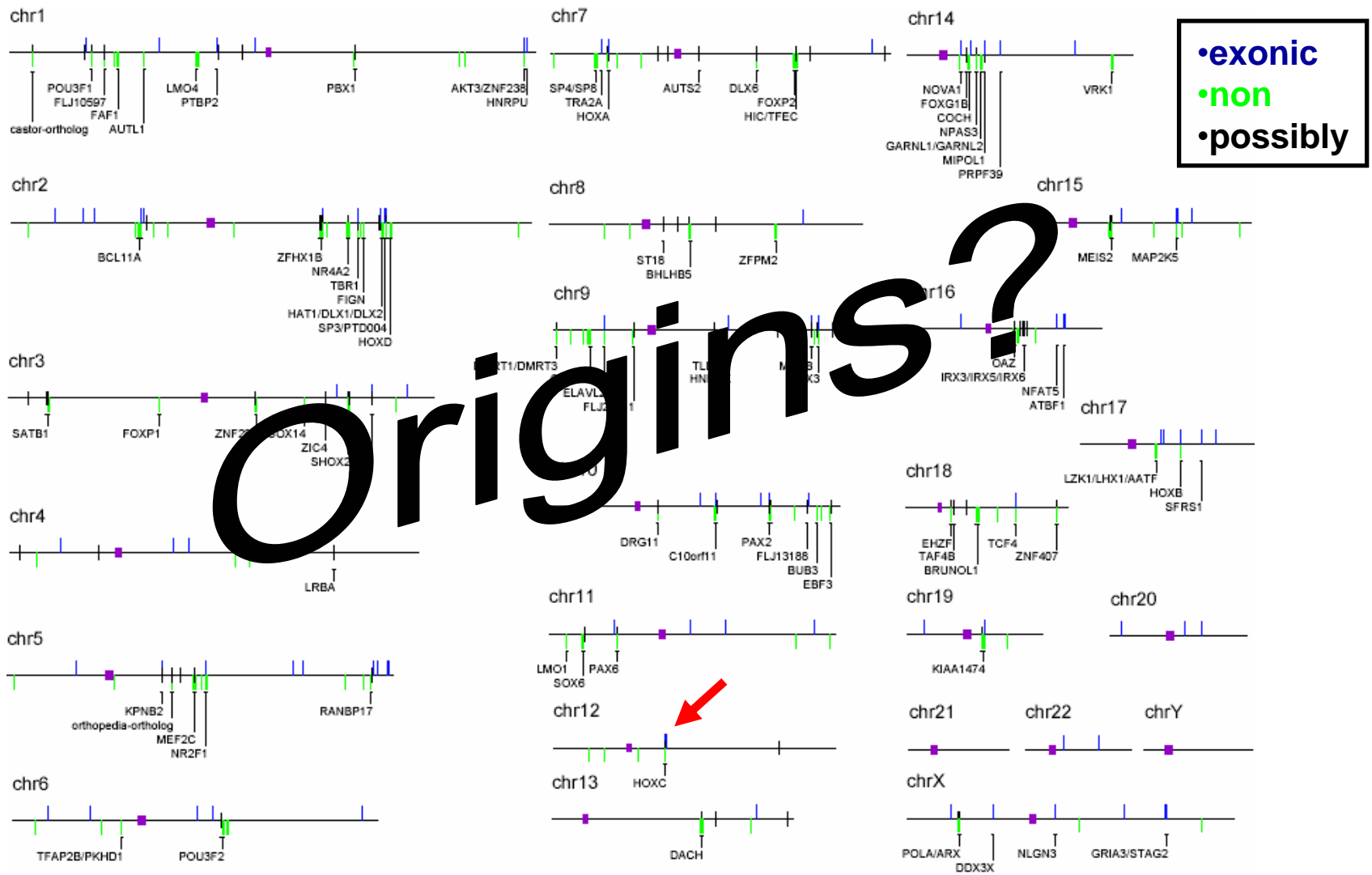
---



# Origins of Ultraconserved Elements?

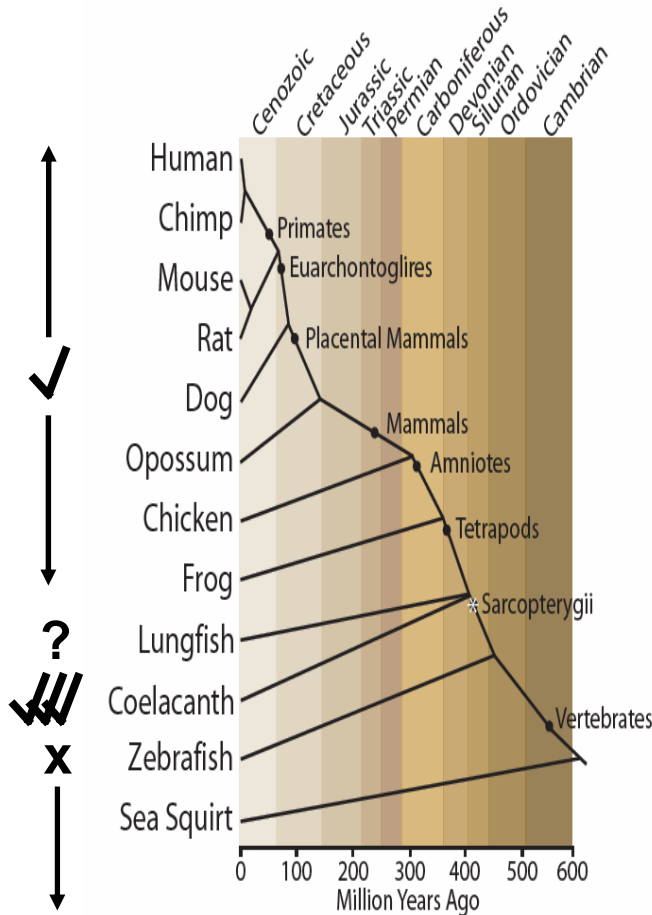


# Genomic Distribution of Ultraconserved Elements



# Uniquely Abundant in Coelacanth

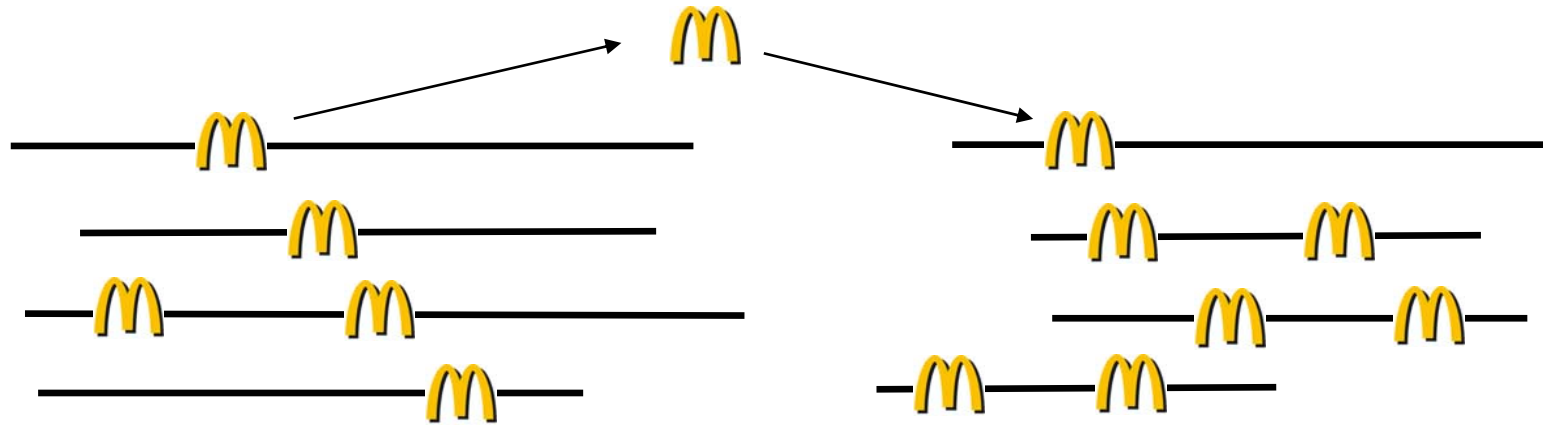
Upto 80%id between Coelacanth instances and some human instances, inc uc.338.



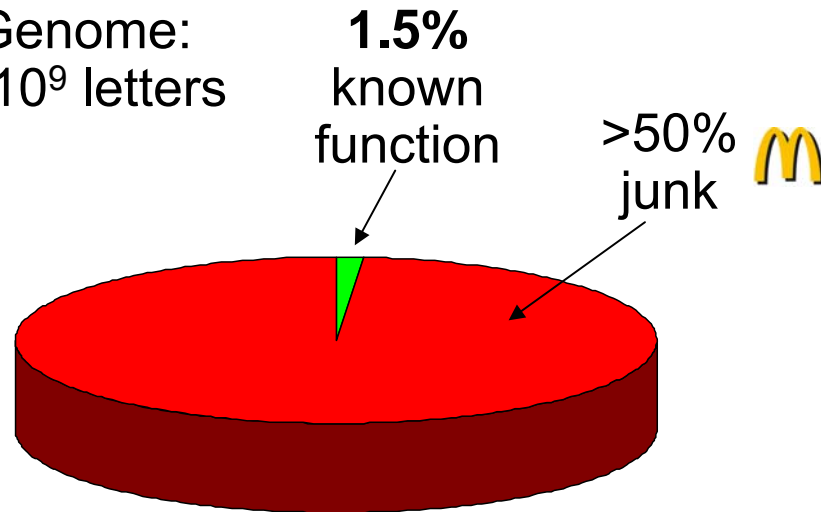
Species	UCSC Assembly	LF-SINE Detected	Species	UCSC Assembly	LF-SINE Detected
<i>Homo sapiens</i>	hg17	Yes	<i>Danio rerio</i>	danRer2	No
<i>Pan troglodytes</i>	panTro1	Yes	<i>Tetraodon nigroviridis</i>	tetNig1	No
<i>Macaca mulatta</i>	rheMac1	Yes	<i>Takifugu rubripes</i>	fr1	No
<i>Mus musculus</i>	mm6	Yes	<i>Ciona intestinalis</i>	ci1	No
<i>Rattus norvegicus</i>	rn3	Yes	<i>Strongylocentrotus purpuratus</i>	strPur1	No
<i>Canis familiaris</i>	canFam1	Yes	<i>Drosophila melanogaster</i>	dm2	No
<i>Bos taurus</i>	bosTau1	Yes	<i>Anopheles gambiae</i>	anoGam1	No
<i>Monodelphis domestica</i>	monDom1	Yes	<i>Caenorhabditis elegans</i>	ce2	No
<i>Gallus gallus</i>	galGal2	Yes	<i>Saccharomyces cerevisiae</i>	sacCer1	No
<i>Xenopus tropicalis</i>	xenTro1	Yes			

- ✓ 100 diverged copies in a Gigabase
- ⚡ 60 highly similar copies in a Megabase

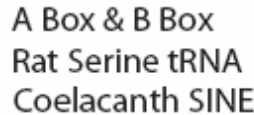
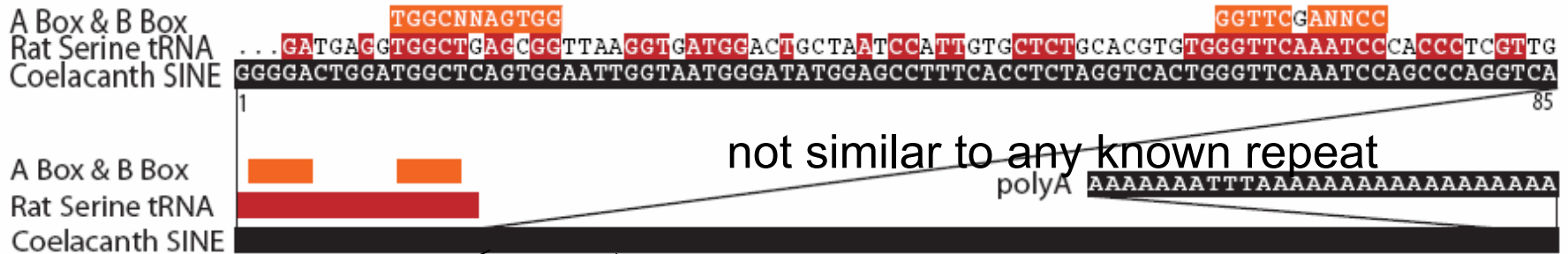
# Repeats / Mobile Elements ("selfish DNA")



Human  
Genome:  
 $3 \times 10^9$  letters



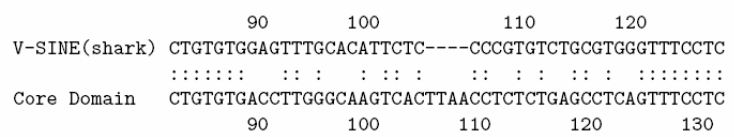
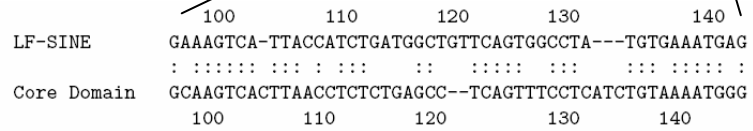
# The LF SINE (for Lobefin Fish / "Living Fossil")



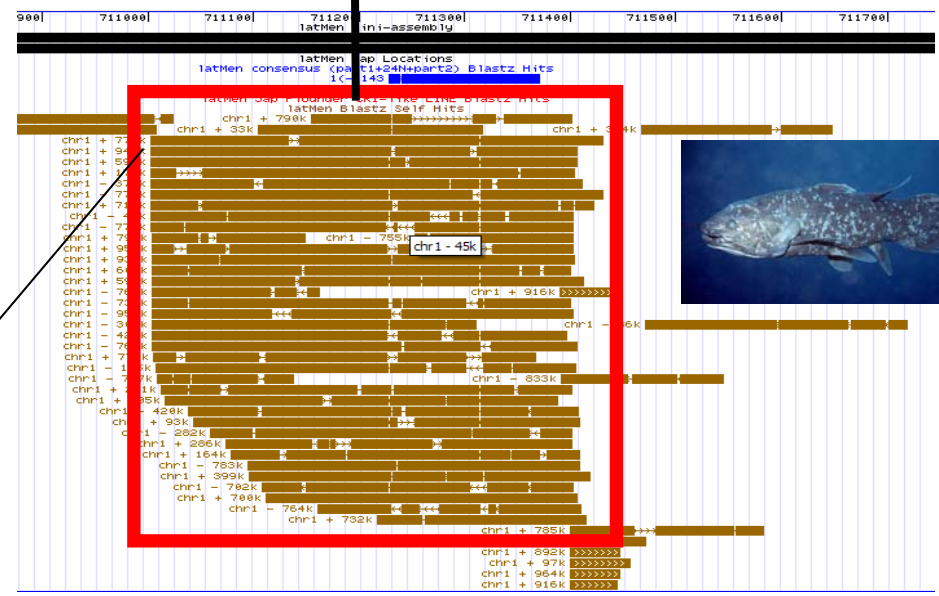
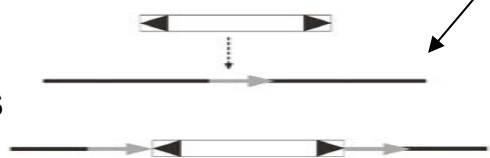
out

back

Reconstruction

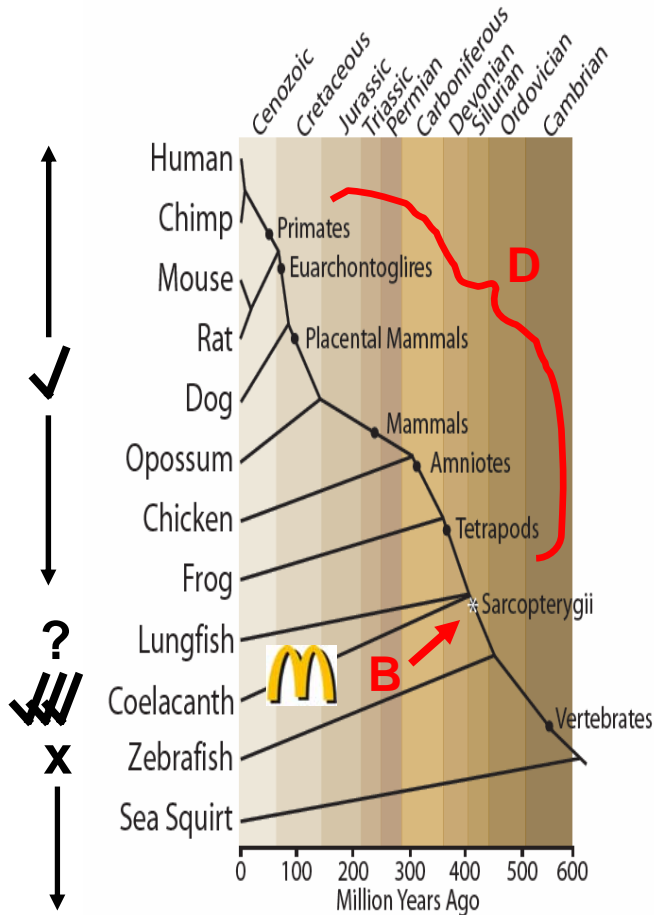


target site duplications



# >360My Old and Going Strong

Upto 80%id between Coelacanth SINE and some human instances, inc uc.338.

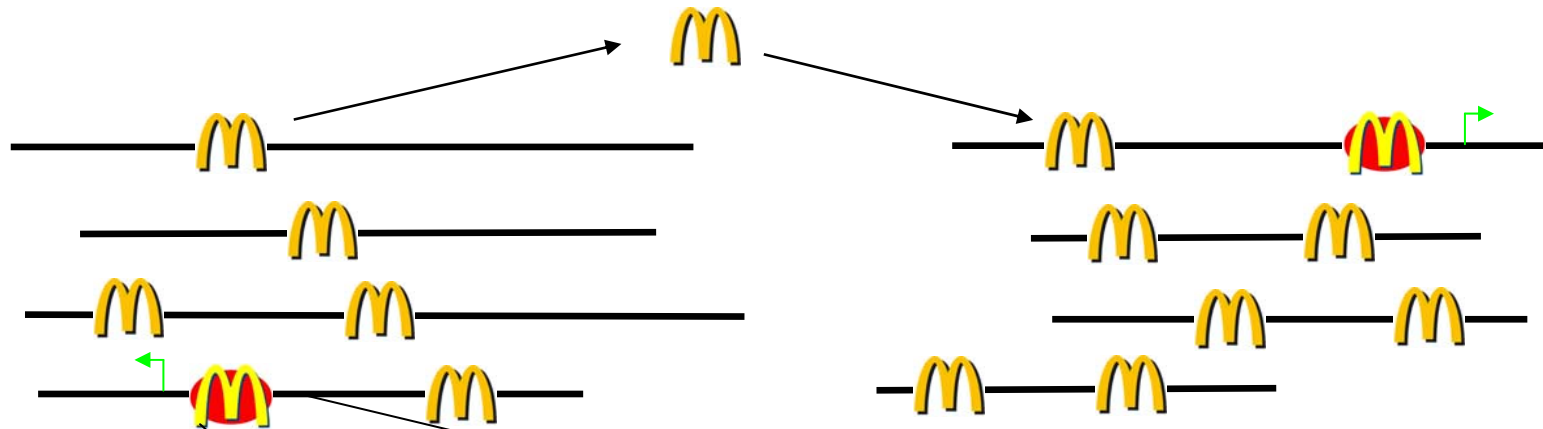


Species	UCSC Assembly	LF-SINE Detected	Species	UCSC Assembly	LF-SINE Detected
<i>Homo sapiens</i>	hg17	Yes	<i>Danio rerio</i>	danRer2	No
<i>Pan troglodytes</i>	panTro1	Yes	<i>Tetraodon nigroviridis</i>	tetNig1	No
<i>Macaca mulatta</i>	rheMac1	Yes	<i>Takifugu rubripes</i>	fr1	No
<i>Mus musculus</i>	mm6	Yes	<i>Ciona intestinalis</i>	ci1	No
<i>Rattus norvegicus</i>	rn3	Yes	<i>Strongylocentrotus purpuratus</i>	strPur1	No
<i>Canis familiaris</i>	canFam1	Yes	<i>Drosophila melanogaster</i>	dm2	No
<i>Bos taurus</i>	bosTau1	Yes	<i>Anopheles gambiae</i>	anoGam1	No
<i>Monodelphis domestica</i>	monDom1	Yes	<i>Caenorhabditis elegans</i>	ce2	No
<i>Gallus gallus</i>	galGal2	Yes	<i>Saccharomyces cerevisiae</i>	sacCer1	No
<i>Xenopus tropicalis</i>	xenTro1	Yes			





# Cis-reg & Ultra elements from Mobile Elements



Co-option event,  
probably due to  
favorable genomic  
*context*

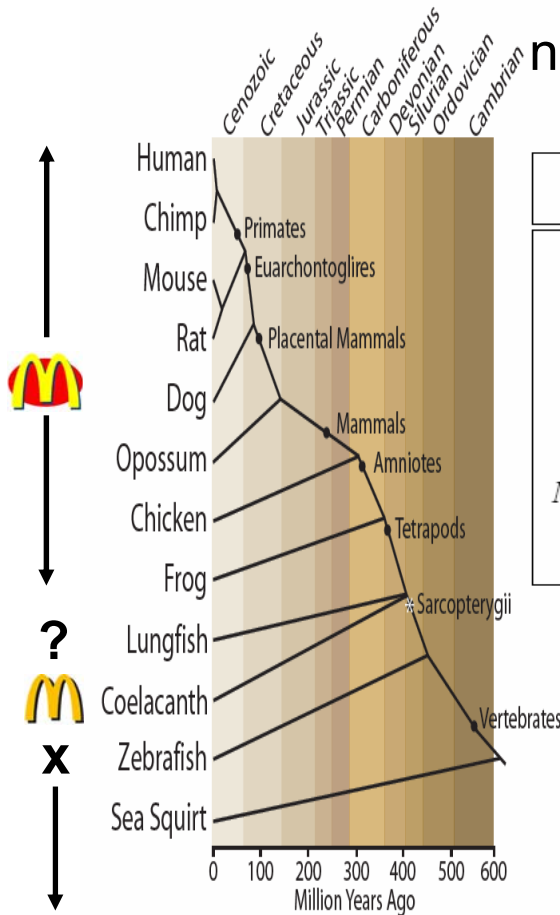


All other copies  
are destined to  
decay over time  
at a neutral rate

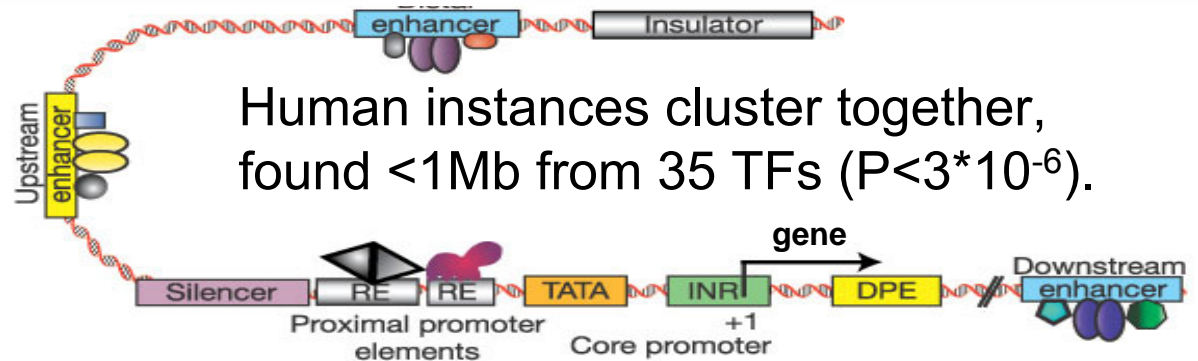
[Yass is a small town in  
New South Wales, Australia.]

# Exapted Into Which Cellular Roles?

No evidence for Transcription (Tx) as small RNAs,  
no orientation preference in introns, not in antisense Tx.

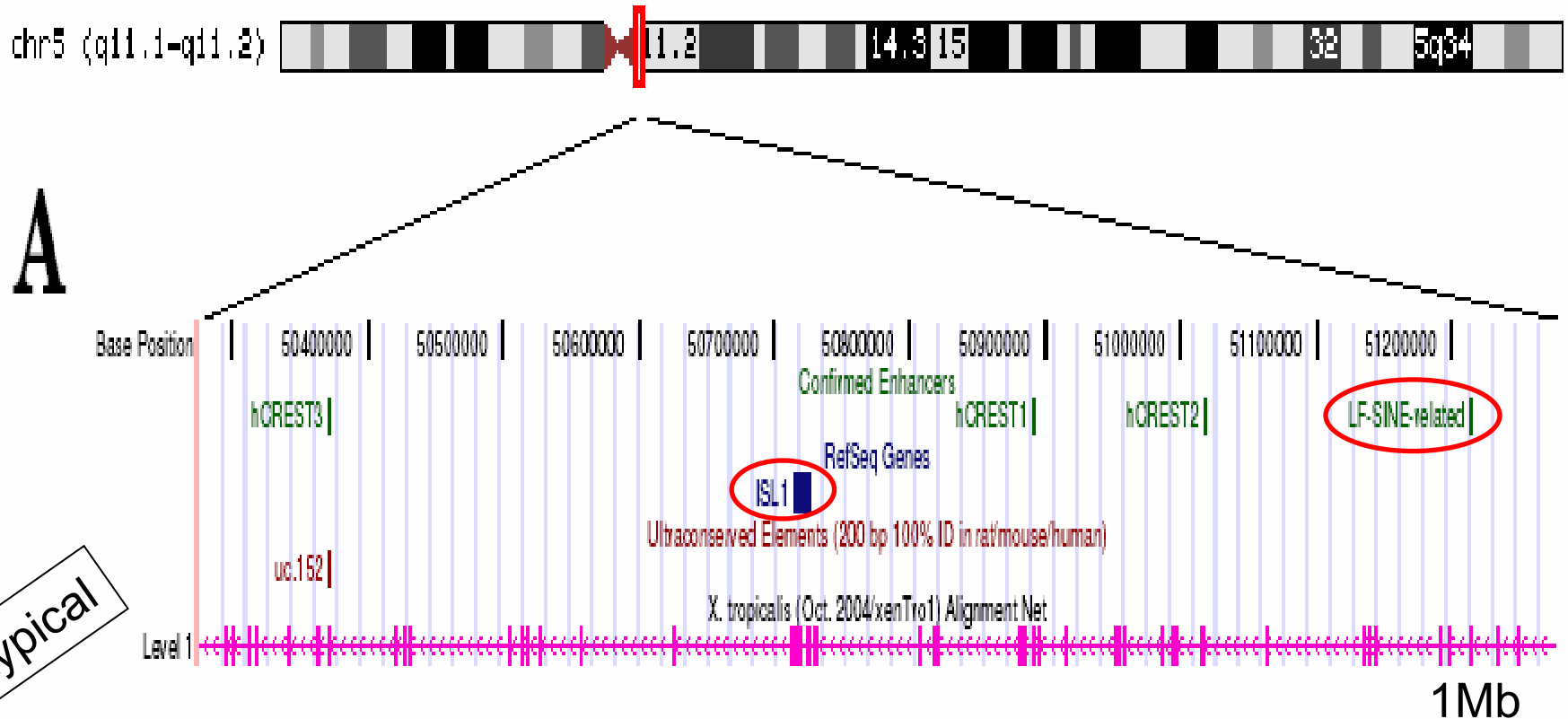


Organism	5' UTR	3' UTR	Exonic		Intronic	Intergenic	Total
			Alt-Spliced	Total			
<i>Homo sapiens</i>	1	0	12	13	68	163	245
<i>Pan troglodytes</i>	-	-	-	-	-	-	210
<i>Macaca mulatta</i>	-	-	-	-	-	-	229
<i>Canis familiaris</i>	-	-	-	-	-	-	235
<i>Bos taurus</i>	-	-	-	-	-	-	169
<i>Mus musculus</i>	0	1	7	8	25	57	91
<i>Rattus norvegicus</i>	-	-	-	-	-	-	87
<i>Monodelphis domestica</i>	-	-	-	-	-	-	394
<i>Gallus gallus</i>	0	1	2	3	244	451	699
<i>Xenopus tropicalis</i>	0	0	1	2	10	14	26



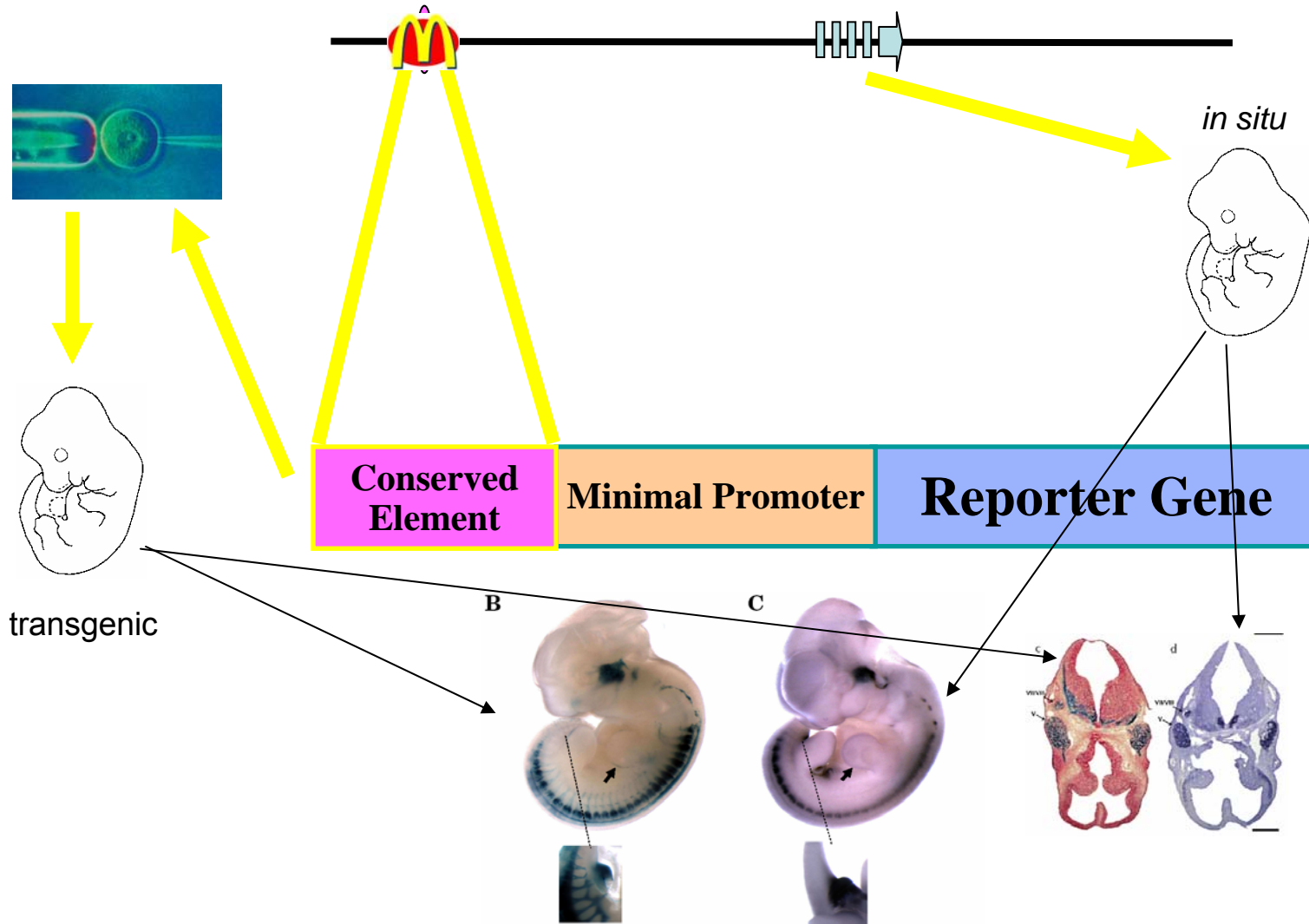
Human instances cluster together,  
found <1Mb from 35 TFs ( $P < 3 \times 10^{-6}$ ).

# Instance 500kb Downstream of ISL1



ISL1 is a neuro-developmental gene, also expressed in testis.  
Three previously known enhancers are conserved across vertebrates.

# Repeat made Regulatory Region



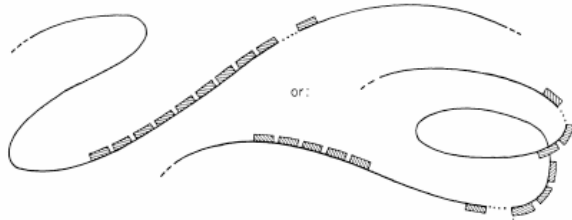
# From junk DNA to pathway recruitments?

JUNE 1971]

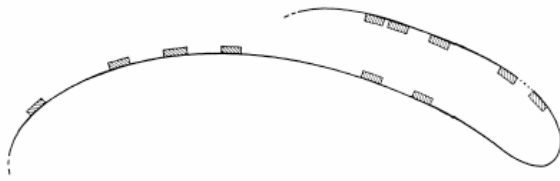
REPETITIVE AND NON-REPETITIVE DNA

127

1) A portion of the genome containing a new saltatory replication ( [ ] ) :



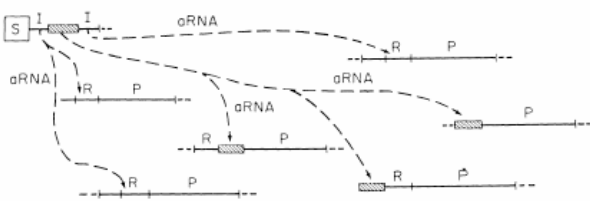
2) "Diffusion" of sequences throughout the genome by subsequent chromosomal rearrangements:



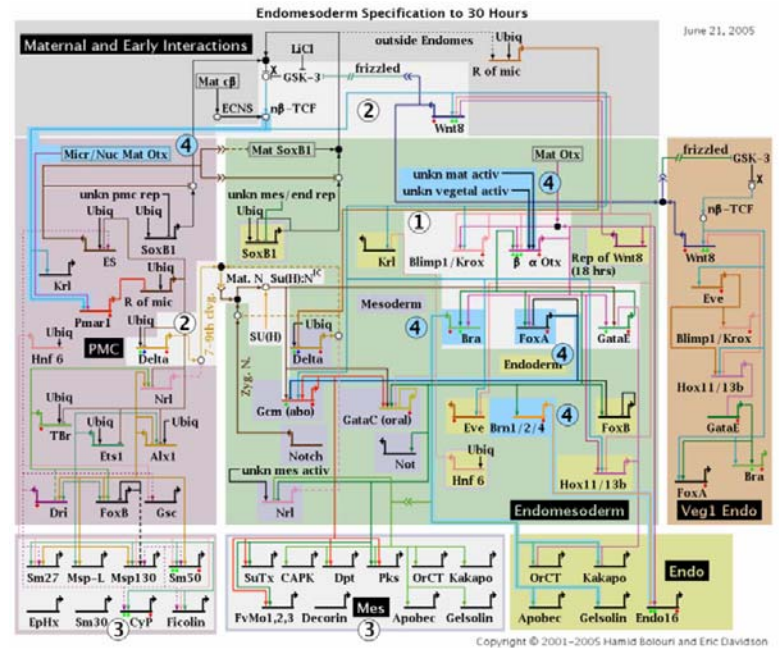
3) Among some local arrangements which might thus arise could be these:



4) In this way new regulative pathways could arise, for example:



[Britten & Davidson, 1971]



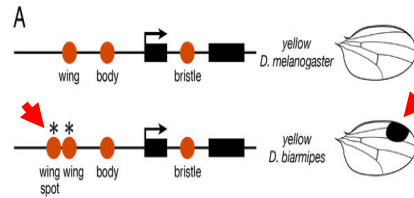
[Davidson & Erwin, 2006]

# The Co-Optionome

quantify co-option

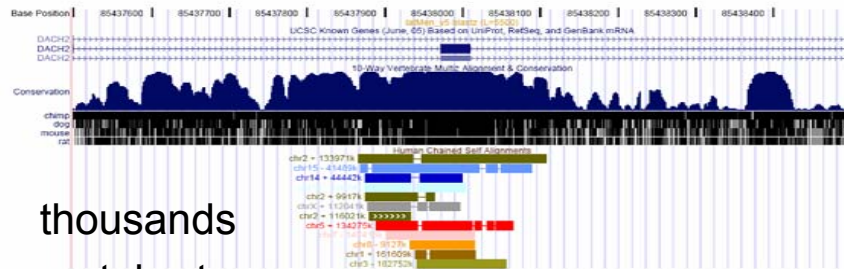
characterized repeats

functional elements



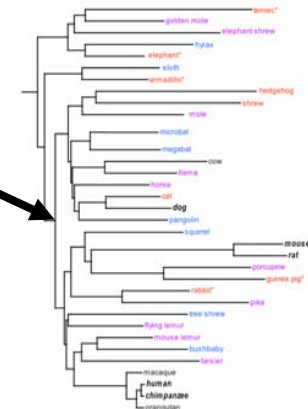
genomic DNA

repeat detectors



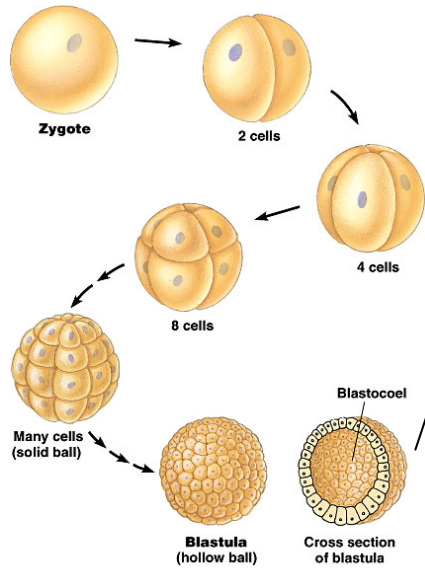
thousands  
vertebrate  
candidates

reconstructing  
Boreoeutherian  
ancestor

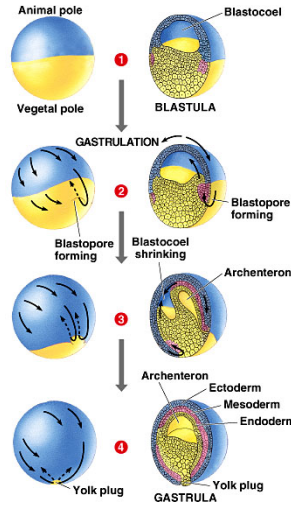


DeuSINE, MER121, ...

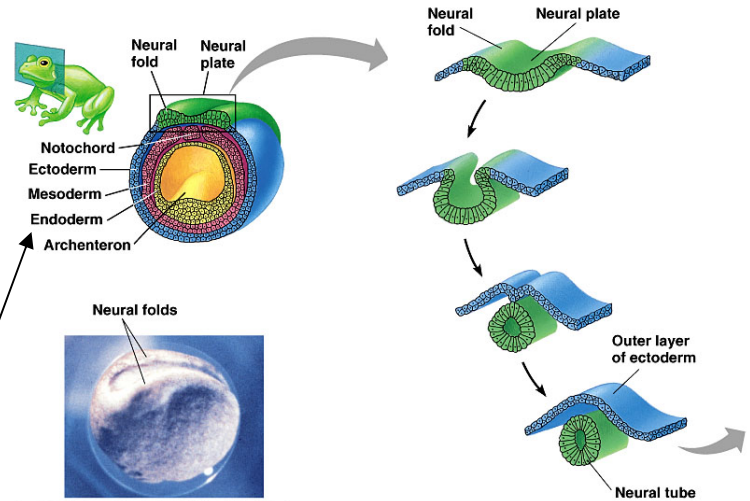
# Objective: Marry Development & Genomics



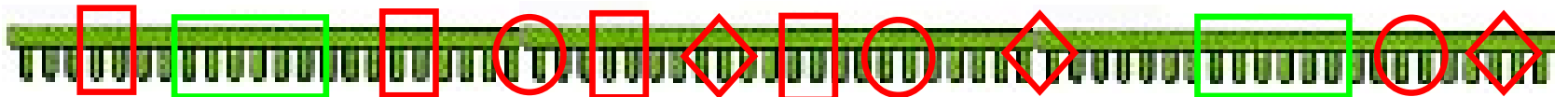
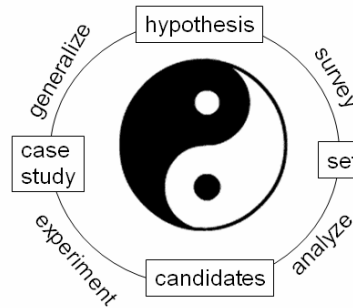
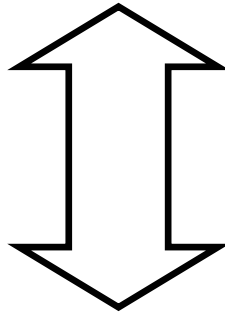
Copyright © 2003 Pearson Education, Inc., publishing as Benjamin Cummings.



Copyright © 2003 Pearson Education, Inc., publishing as Benjamin Cummings.



Copyright © 2003 Pearson Education, Inc., publishing as Benjamin Cummings.



# Bejerano Lab: Research Interests

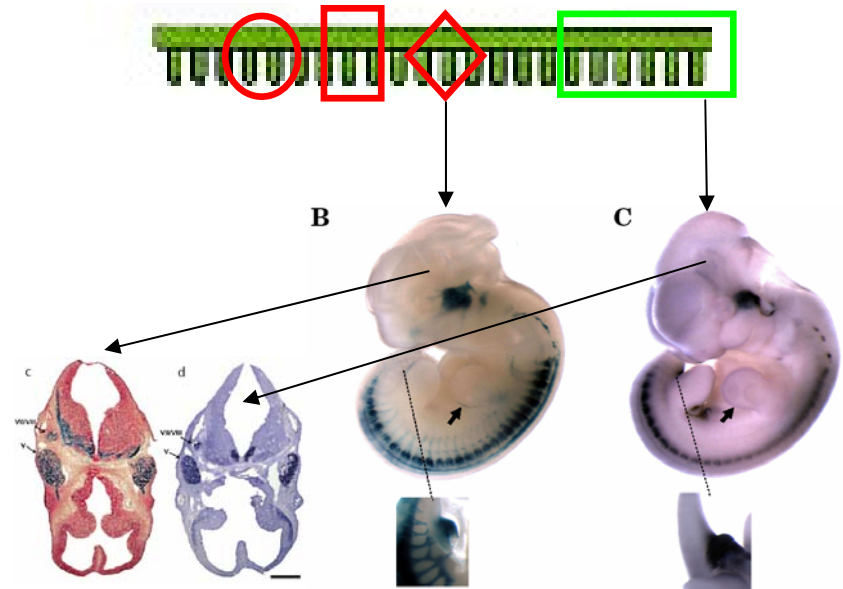
Many thousands of human conserved elements  
congregate en-masse near developmental genes.

Origins & Evolution

Functions & Encoding

Contribution to  
Human Disease

**Break regulatory code**





# Bejerano Lab: Research Interests

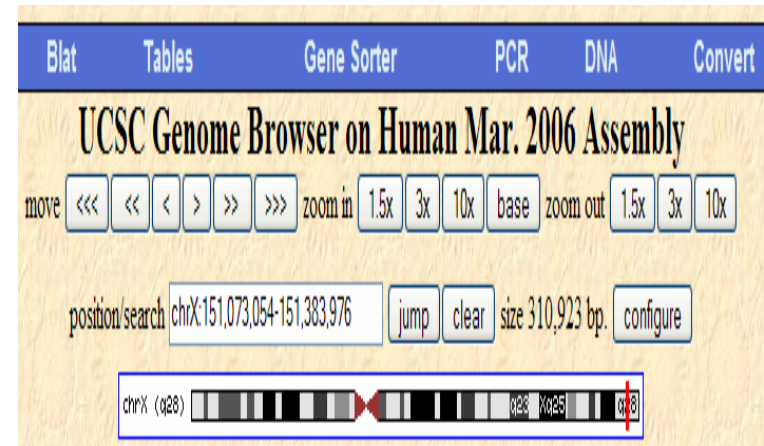
Many thousands of human conserved elements  
congregate en-masse near developmental genes.

Origins & Evolution

Functions & Encoding

Contribution to  
Human Disease

Break regulatory code  
**Provide tools to community**



thousands and thousands of  
page requests served *daily*

# Bejerano Lab: Research Interests

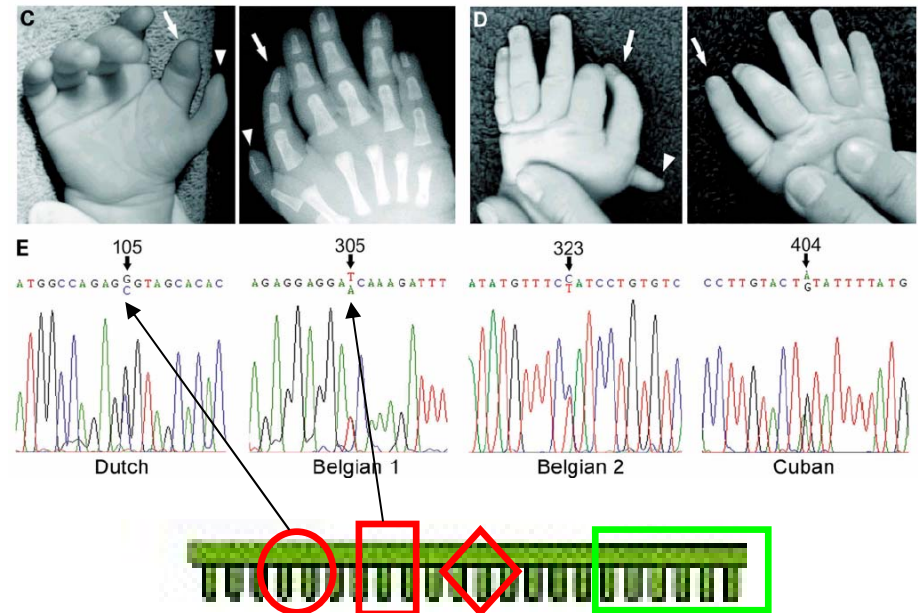
Many thousands of human conserved elements  
congregate en-masse near developmental genes.

Origins & Evolution

Functions & Encoding

Contribution to  
Human Disease

Break regulatory code  
Provide tools to community  
**Improve human health**





# Kudos

---

## UC Santa Cruz

David Haussler

Sofie Salama, Jim Kent, Craig Lowe,  
Bryan King, Adam Siepel, Jakob Pedersen  
Katie Pollard, Courtney Onodera  
Rachel Harte, Genomics/Browser Group

## Lawrence Berkeley Labs

Eddy Rubin

Nadav Ahituv

## McGill U.

Mathieu Blanchette

## Penn State U.

Webb Miller's group

## U. Queensland

John Mattick's group

Genome Sequencing Consortia  
All GenBank contributors

Gill Bejerano

bejerano@stanford.edu

