



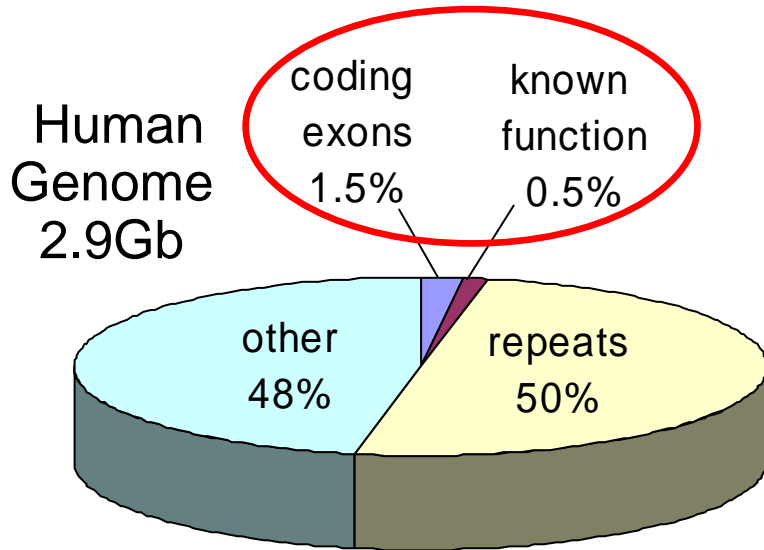
# Genomics

PHASE TWO: INTERPRETATION

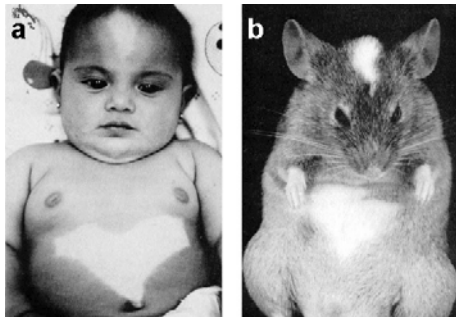
SLEEMAN The New Ledger



# Conserved Non Coding Regions



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



40% DNA alignable  
95% coding genes shared

of the 5%  $\frac{1}{3}$  coding  
 $\frac{2}{3}$  non coding

# What To Do?

Assay!

Which Assay?

Which elements?



UCSC Genome Browser Home - Mozilla Firefox

http://genome.ucsc.edu

UCSC Genome Bioinformatics

Genomes - Blat - Tables - Gene Sorter - PCR - Proteome - FAQ - Help

Genome Browser  
ENCODE  
Blat  
Table Browser  
Gene Sorter  
In Silico PCR  
Proteome Browser  
Utilities  
Downloads  
Release Log  
Custom Tracks  
Mirrors  
Archives  
Credits  
Publications  
Training

<http://www.nature.com/naturemethods>

CSHL PRESS

PUBLISHED IN ASSOCIATION WITH  
COLD SPRING HARBOR LABORATORY PRESS

## Computational screening of conserved genomic DNA in search of functional noncoding elements

Gill Bejerano<sup>1</sup>, Adam C Siepel<sup>1</sup>, W James Kent<sup>1</sup> & David Haussler<sup>1,2</sup>

<sup>1</sup>Center for Biomolecular Science and Engineering, <sup>2</sup>Howard Hughes Medical Institute, University of California Santa Cruz, Santa Cruz, California 95064, USA. Correspondence should be addressed to G.B. (jill@soe.ucsc.edu).

The sequencing of the mouse genome allowed, for the first time, the large-scale estimation of the extent of sequence conservation within our own genome. In particular, it suggested that in mammals there is at least twice as much conserved genomic DNA as there is protein coding DNA<sup>1</sup>. The abundance of conserved noncoding regions holds even for so-called ultraconserved elements at the very tip of the mammalian

PROTOCOL

one *million* elements  
genome wide...

[Bejerano *et al.*, *Nature Methods* 2005]

# DNA Replication is Imperfect

Small Scale: bases are mutated, deleted, inserted

Medium Scale: regions are duplicated, deleted, inverted

Large Scale: whole chromosomes are duplicated, merged

*junk*

*functional*

...ACGTACGACTGACTAGCATCGACTACGA...

**duplication**

*functional*

*functional*

...ACGTACGACTGACTAGCATCGACTACGA.....TCTGACTAGCATCGACTACGA...

**divergence**

*functional'*

CG

...ACGTACGACTGACTAGCATCGACTACGA.....TCTGACTAGCATCGACTACGA...

*functional''*

AA

# Computational Approach

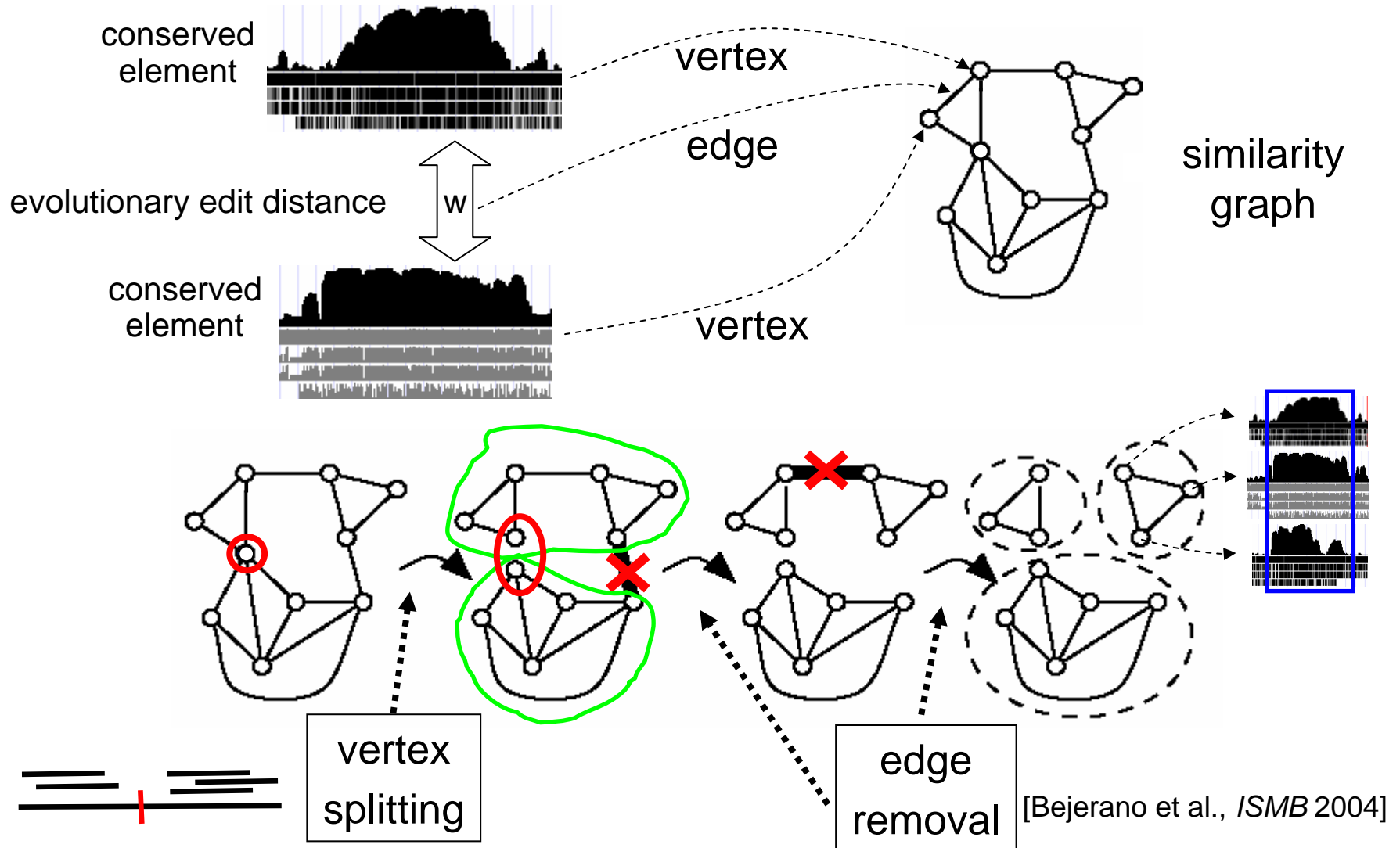


Group them into *paralog families* of human functional regions of common origins:

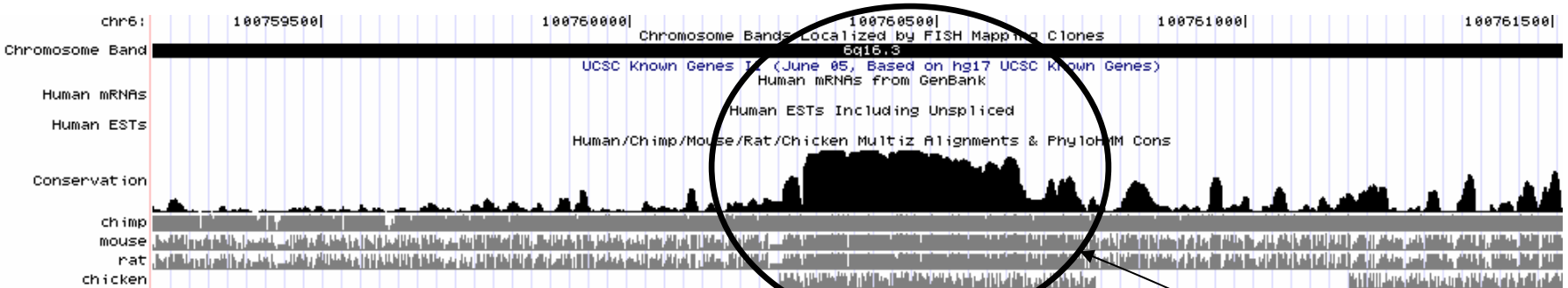
- Annotated members induce function on all.
- Examine core, substitutions in family.
- Test for “*guilt by association*”.

[Bejerano et al., *ISMB* 2004]

# Conserved Element Clustering



# Functional Annotation by Families



After removing from top 5% Human *all* annotated regions, and more:  
700,000 elements, covering 3.5% Human Genome

metrics

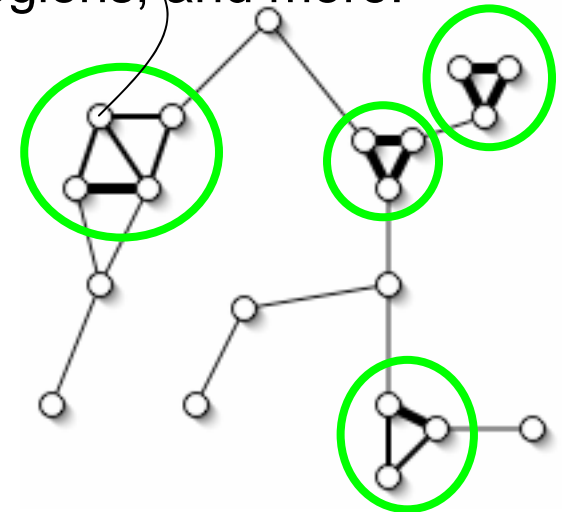
## Puzzling News:

96% of the 700,000  
are unique(!)

## Good News:

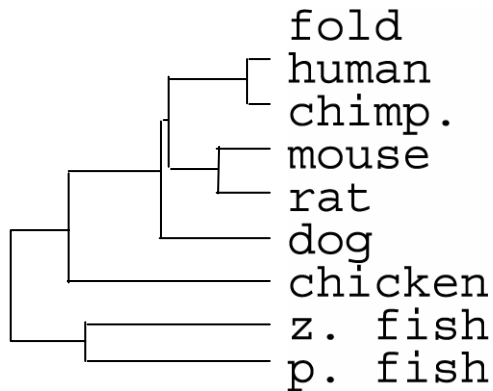
We still find  
**12,027 families**

**novel** putative ncRNAs, cis-regulatory elements, etc.



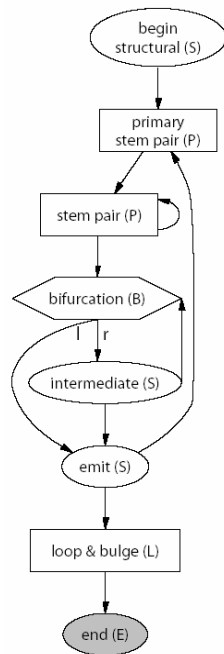
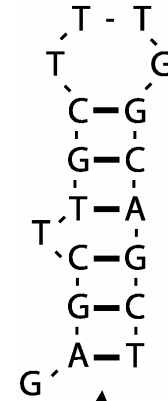
[Bejerano et al., *ISMB* 2004]

# Families of Novel ncRNAs

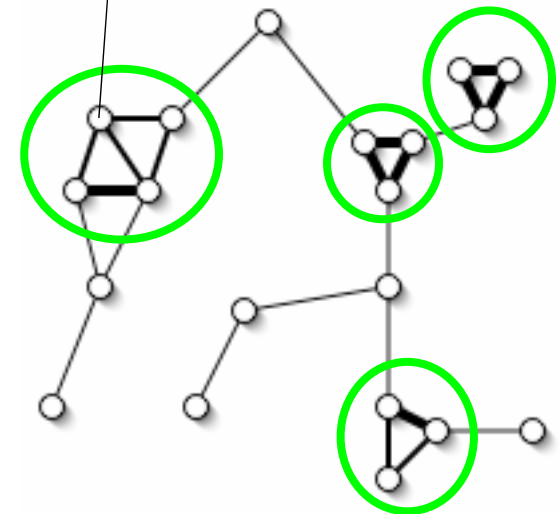


```

((((((((.....)))))))))
GAGCTTGCTTTGGCAGCT
GAGCTTGCTTTGGCAGCT
GAGTTTGCTTTGGCAGCT
AAGCTTACTTACGTAGCT
GAGCATACTAAGGTGGCT
GGGCTTACGCTGGTGGCC
GGGCTTACAATTGTGGCC
GGGCTTAAAAATTTGGCC
    
```

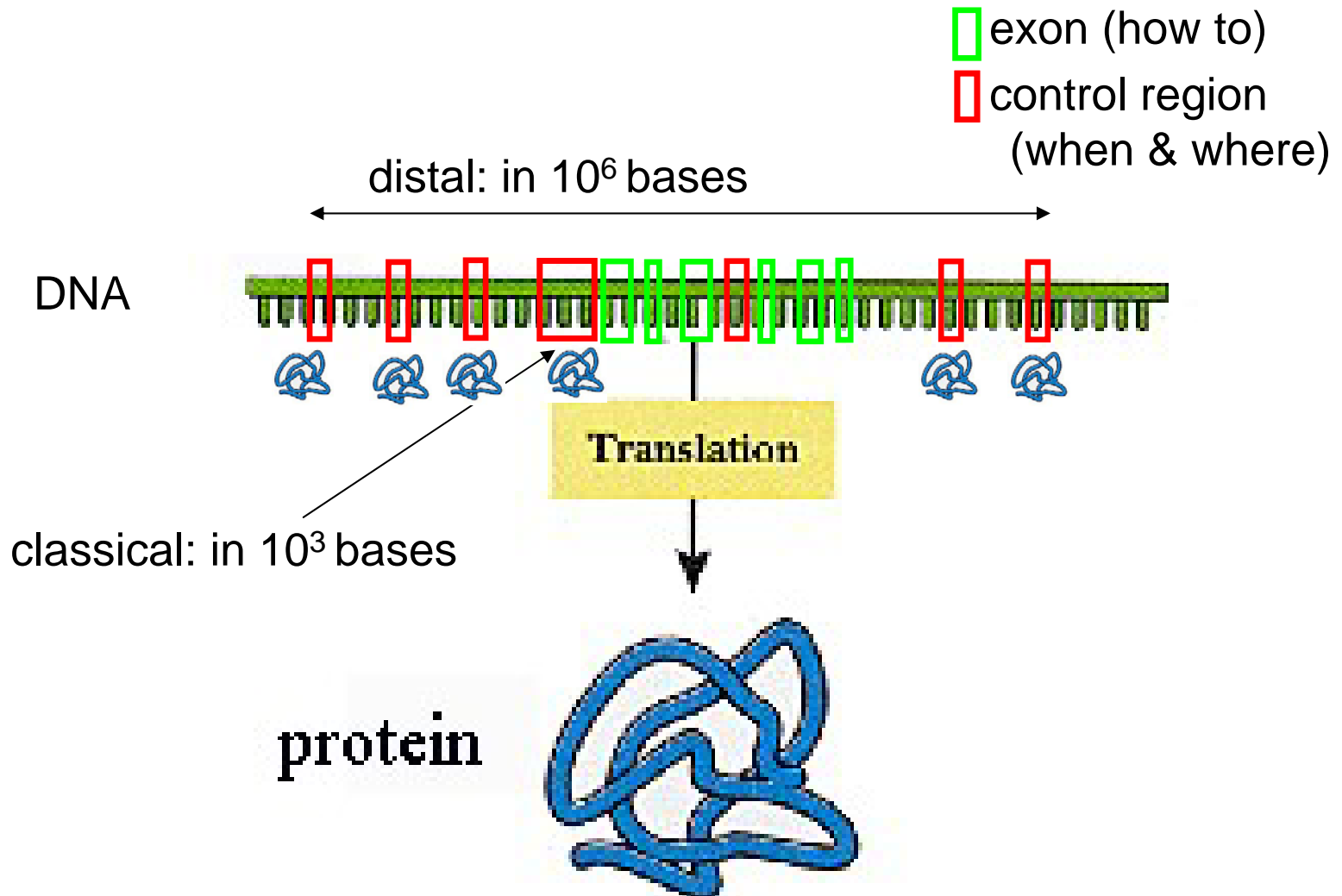


stochastic context  
free grammar model



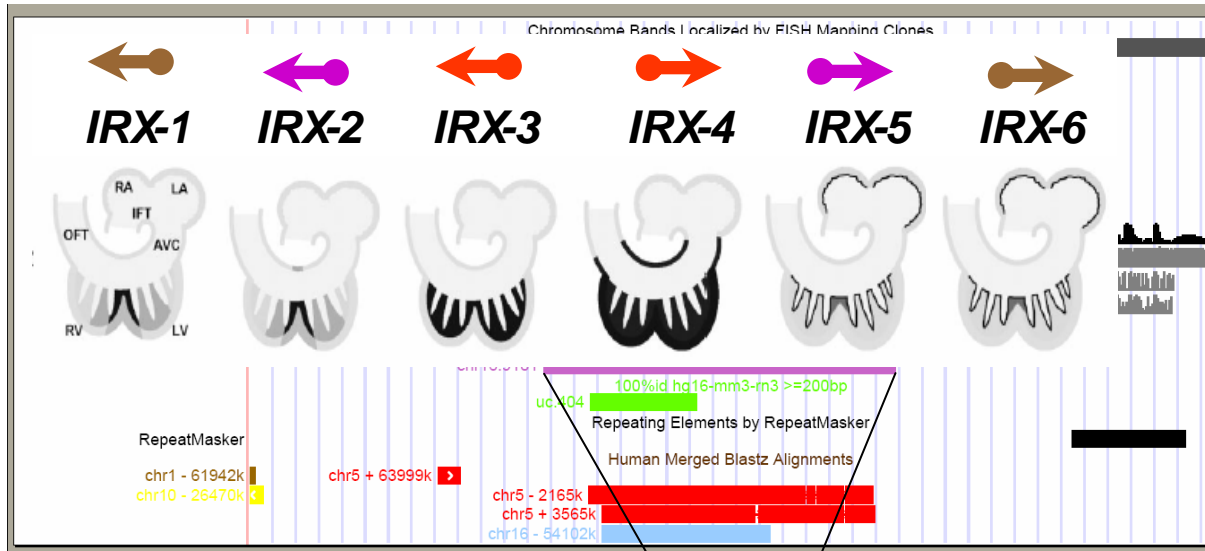
[Pedersen, Bejerano et al., *PLoS CompBio*, 2006]

# Genes, Proteins and Gene Control

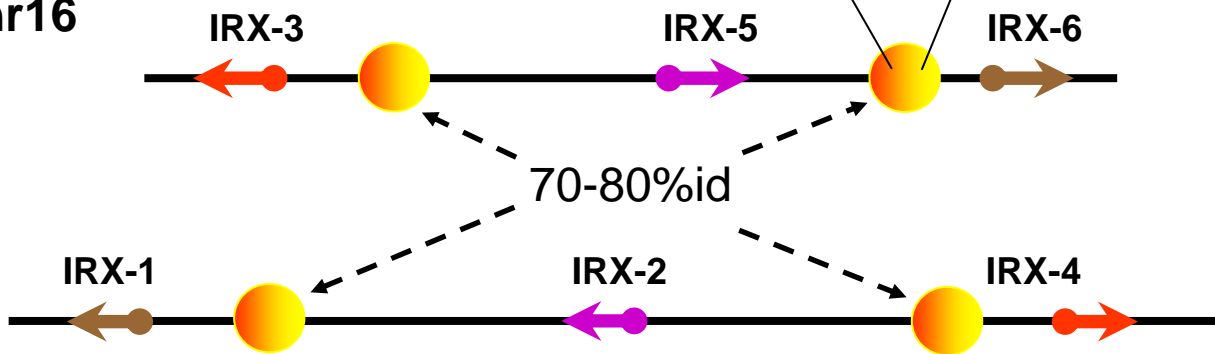


# Cis-Regulatory Families

early  
body  
patterning

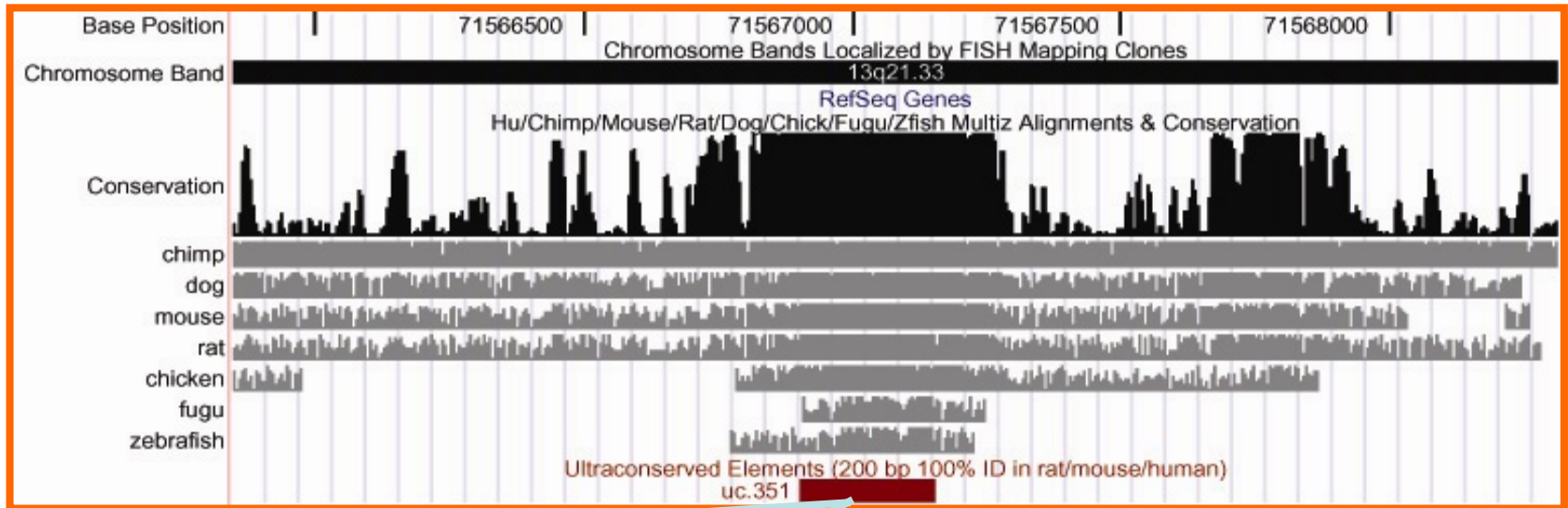


Hs.chr16



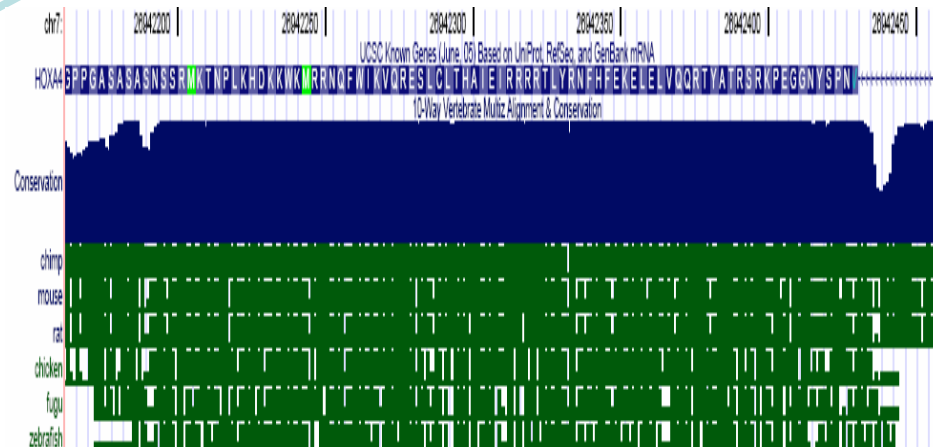
Hs.chr5

# Ultraconserved Elements



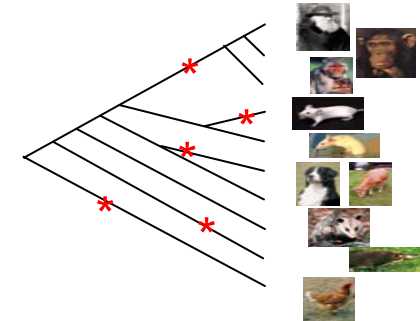
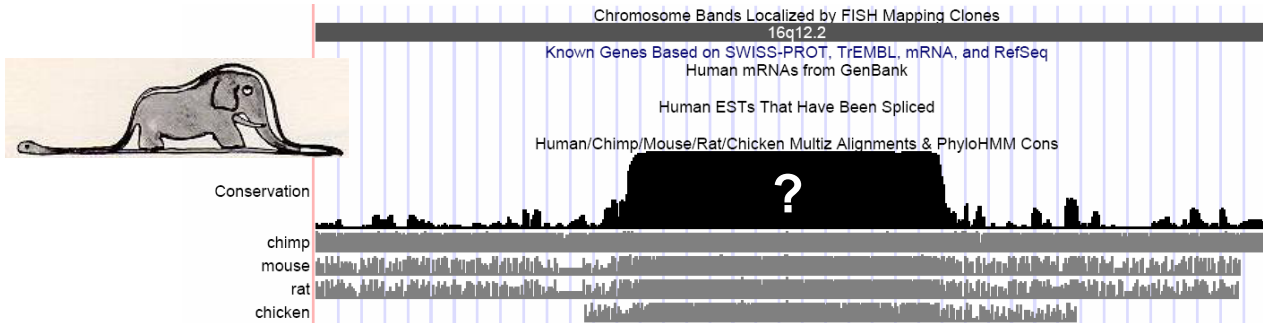
Human/Chimp/Mouse/Rat/Chicken/Fugu/Zebrafish Multiz Alignments & PhyloHMM Cons

Gaps	human	chimp	mouse	rat	chicken	fugu	zebrafish
	CAATGTTAAGAAGCAATTTAAATGAAATTCATTTAAATGTCAGCTCTCATGCTTGCGATTTTCTATAATTGGTATTTAAAT	CAATGTTAAGAAGCAATTTAAATGAAATTCATTTAAATGTCAGCTCTCATGCTTGCGATTTTCTATAATTGGTATTTAAAT	CAATGTTAAGAAGCAATTTAAATGAAATTCATTTAAATGTCAGCTCTCATGCTTGCGATTTTCTATAATTGGTATTTAAAT	CAATGTTAAGAAGCAATTTAAATGAAATTCATTTAAATGTCAGCTCTCATGCTTGCGATTTTCTATAATTGGTATTTAAAT	CAATGTTAAGAAGCAATTTAAATGAAATTCATTTAAATGTCAGCTCTCATGCTTGCGATTTTCTATAATTGGTATTTAAAT	CAATGTTCAAAAGCAATTTAAATGAATTCATTTAAATGTCAAATCTCATGGTTATGATTTTCTATAATTGGCTATTTAAAT	CAATGTTCAAAAGCAATTTAAATGAAATTCATTTAAATGTCAAATCTCATGGTTAGGATTTTCTATAATTGGCTATTTAAAT



[Bejerano et al., *Science* 2004]

# No known function requires *this* much conservation

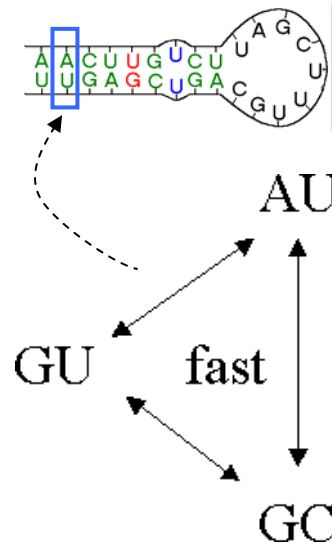


## The genetic code

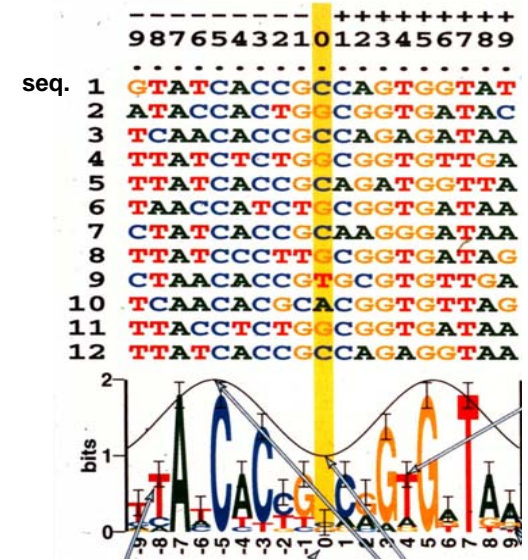
	U	C	A	G	
U	Phe	<b>CUU</b> <b>CUC</b> <b>CUA</b> <b>CUG</b>	Trp	Phe	U C A G U C A G U C A G
Leu	<b>Leucine</b>		Trp	Phe	
C			Leu	Trp	Phe
A			Ile	Trp	Phe
G		Val	Glu	Trp	

## CDS

## ncRNA



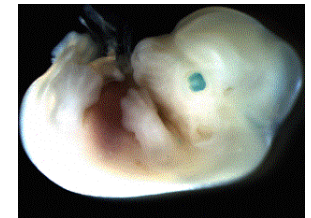
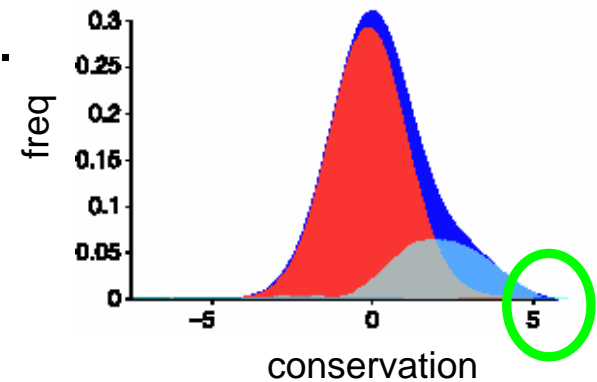
## TFBS



# Ultraconserved Elements

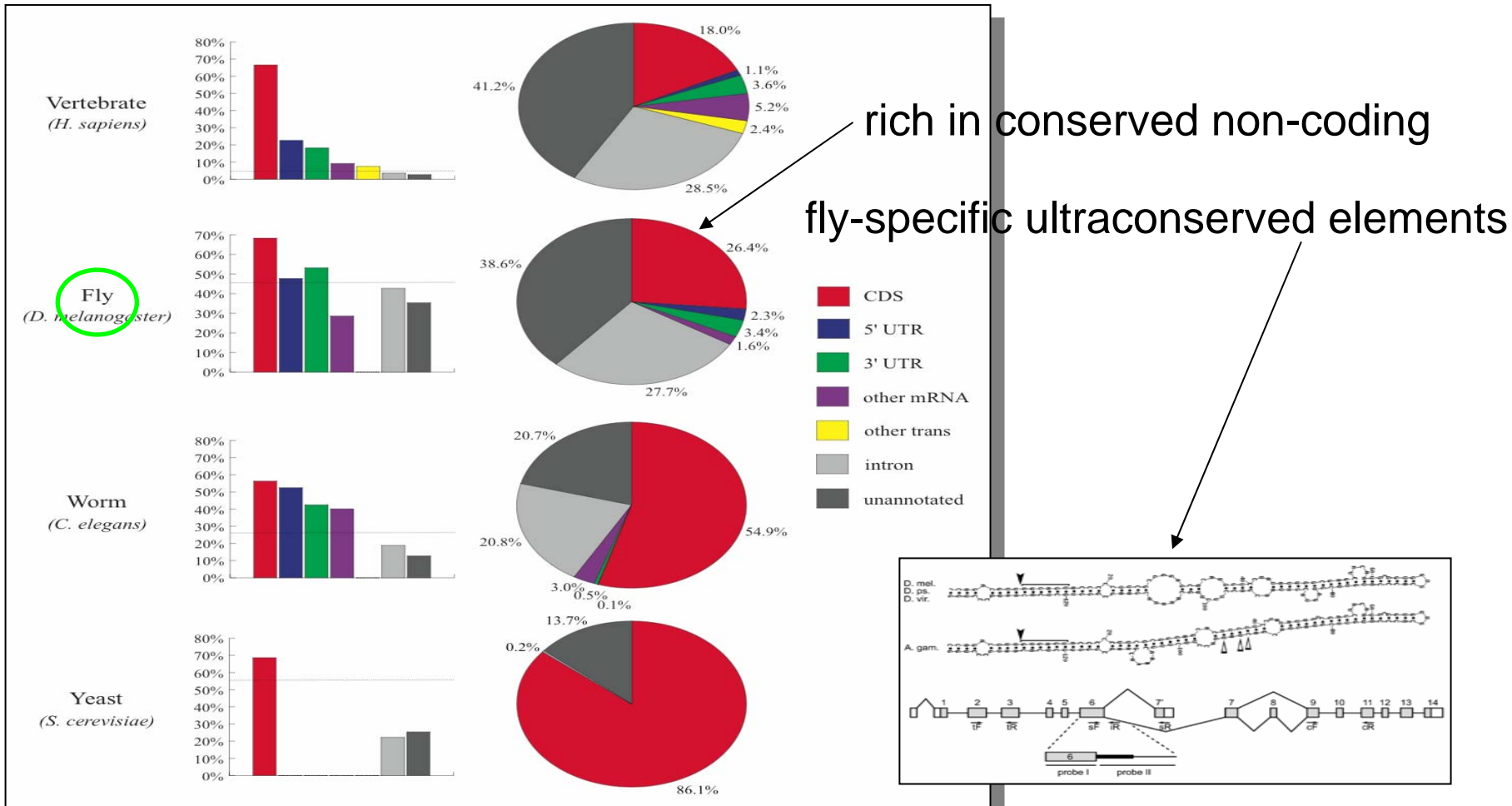
481 regions *perfectly* conserved over 200bp or more, between human, mouse and rat ( $P < 10^{-22}$  neutrally evolving)

- 20-fold fewer SNPs than human average.
- Most do *not* overlap coding DNA.
- The **non-exonic** tend to cluster spatially, in or near DNA binding proteins. Dozens validated as enhancers by now.
- The **exonic** tend to overlap alternatively spliced exons, in RNA binding proteins.
- The ultras cannot be found beyond vertebrates.
- The tip of a continuum of very slowly evolving elements.



[Bejerano et al., *Science* 2004  
Chicken Consortium, *Nature* 2004]

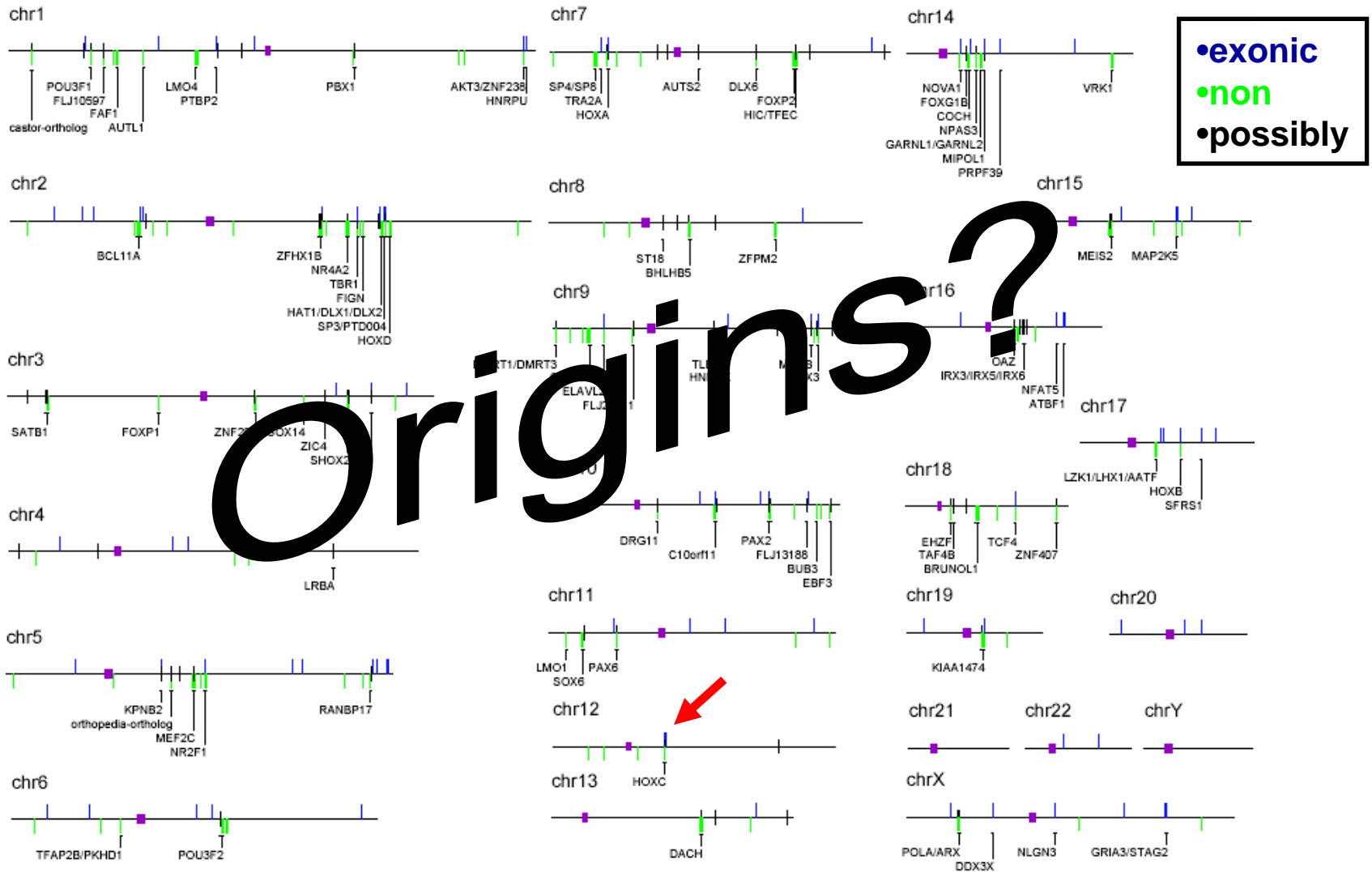
# Similar Phenomena in Flies



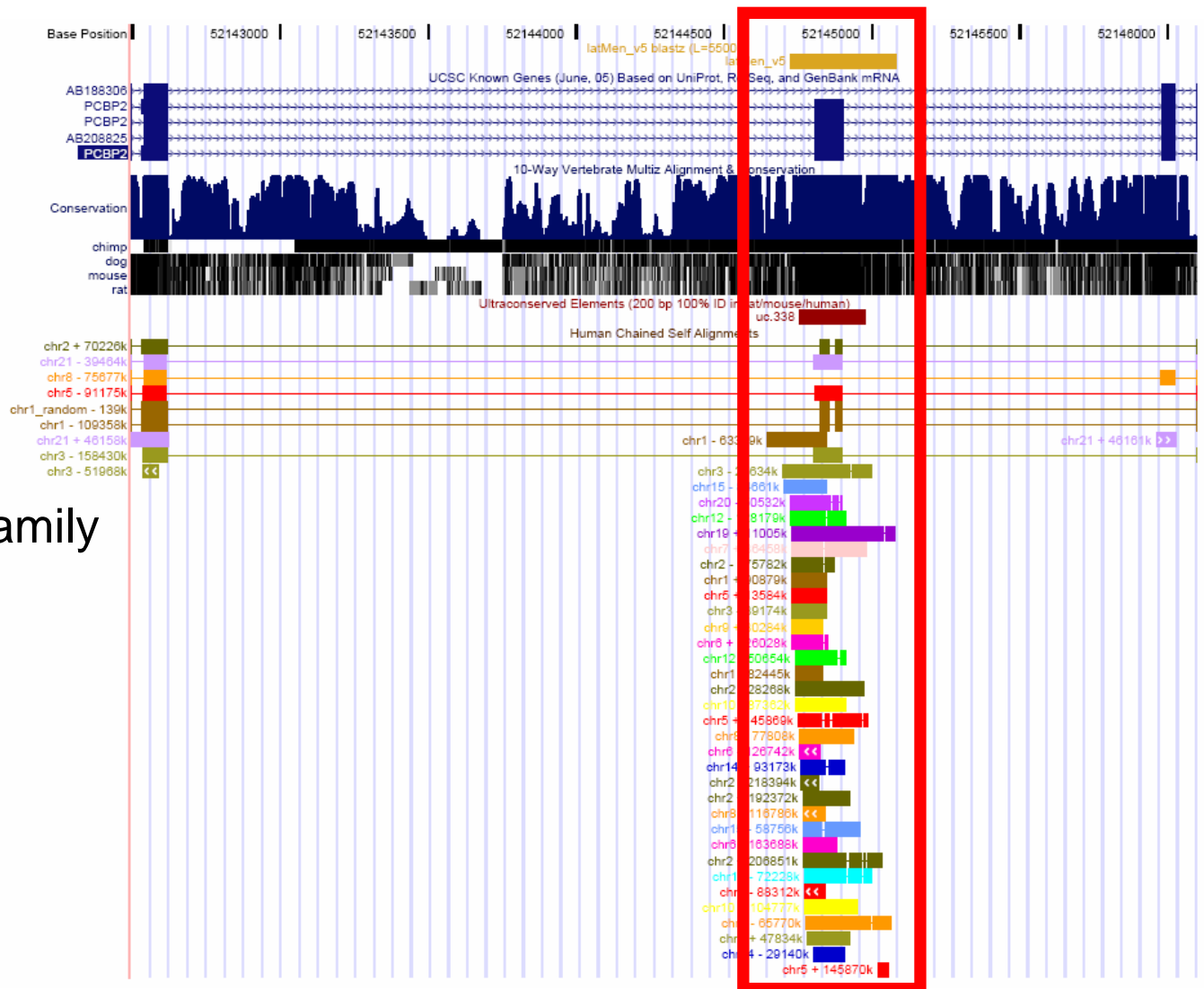
[Siepel, Bejerano et al., *Genome Research* 2005]

[Glazov, ..., Bejerano, Mattick, *Genome Research* 2005]

# Mammalian Ultras - Genomic Distribution



# Ultraconserved Element uc.338



uc.338 paralog family

# Coelacanth Homologs Closer than Human Ones

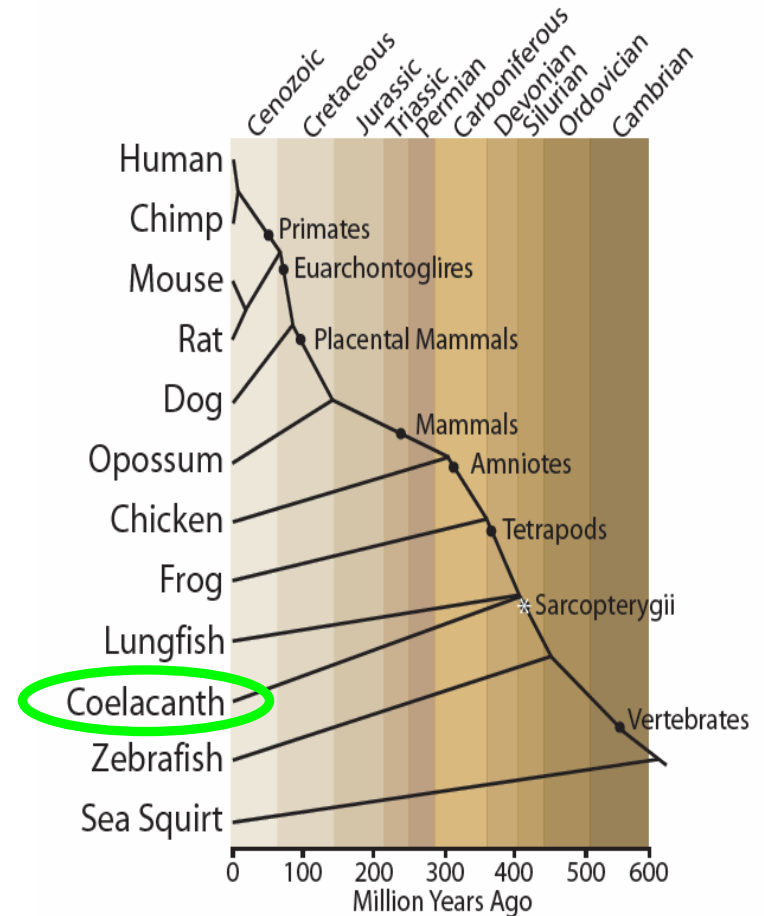
```
Searching...10...20...30...40...50...60...70...80...90...100% done
```

Sequences producing High-scoring Segment Pairs:							High	Smallest	
							Score	Sum	Probability
								P (N)	N
EM_HUM:AC023509	AC023509.42	Homo sapiens	12	BAC	RP11-793H...	1315	3.6e-51	1	
EM_MUS:AC137156	AC137156.3	Mus musculus		BAC clone	RP23-11...	1297	2.3e-50	1	
EM_MUS:AF236844	AF236844.1	Mus musculus		alpha-CP2					
EM_OV:AC151571	AC151571.1	Latimeria menadoensis		cl					
EM_PAT:AR237520	AR237520.1	Sequence 52 from patent							
EM_OV:AC150308	AC150308.1	Latimeria menadoensis		cl					
EM_HUM:AC006127	AC006127.1	Homo sapiens chromosome							
EM_HUM:AF254822	AF254822.1	Homo sapiens SMARCA4 is							
EM_OV:AC150309	AC150309.1	Latimeria menadoensis		cl					
EM_OV:AC150283	AC150283.1	Latimeria menadoensis		cl					
EM_OV:AC150284	AC150284.1	Latimeria menadoensis		cl					
EM_OV:AC147788	AC147788.1	Latimeria menadoensis		cl					
EM_HUM:AC005076	AC005076.2	Homo sapiens BAC clone							
EM_OV:AF131253	AF131253.1	Latimeria chalumnae rhoc							
EM_MUS:MMNUABPRO	L19661.1	Mus musculus (clones pB1001, pB...				498	1.2e-14	1	
EM_RO:BC078909	BC078909.1	Rattus norvegicus poly(rC) bind...				498	1.3e-14	1	



[Bejerano, Lowe et al., *Nature*, 2006]

# Coelacanth “the Living Fossil” Fish



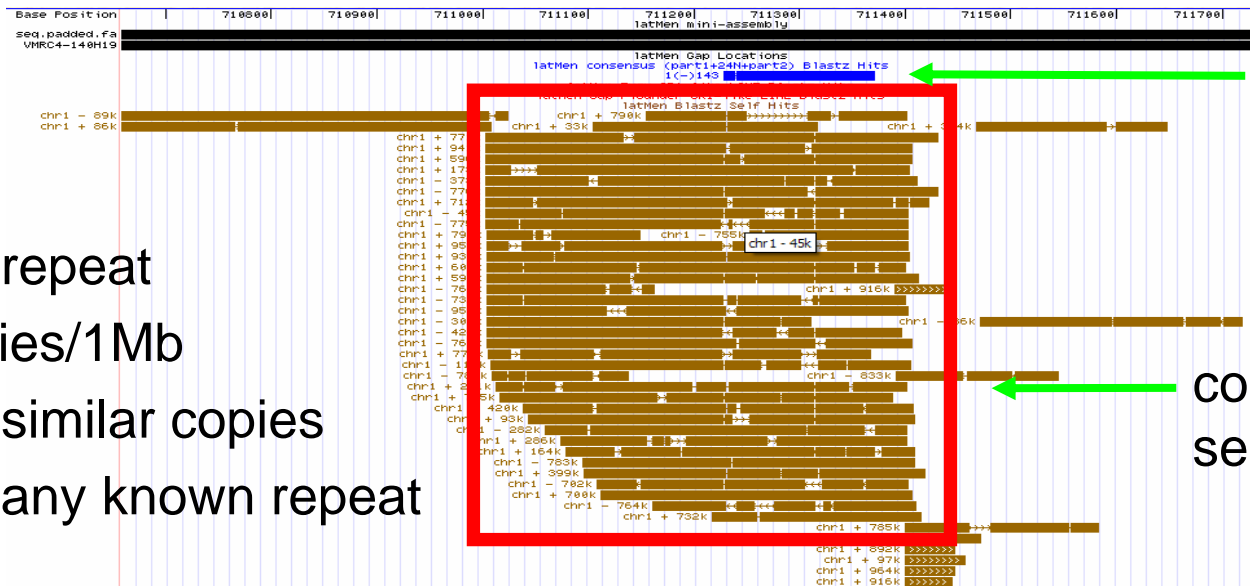
Appeared in Fossil Record >360Mya; Peaked 240Mya;  
Disappeared 80Mya; Rediscovered (by science) in 1938.



# Coelacanth Repeat

1 Mb sequenced.  
4 genomic regions.  
Instances in all.

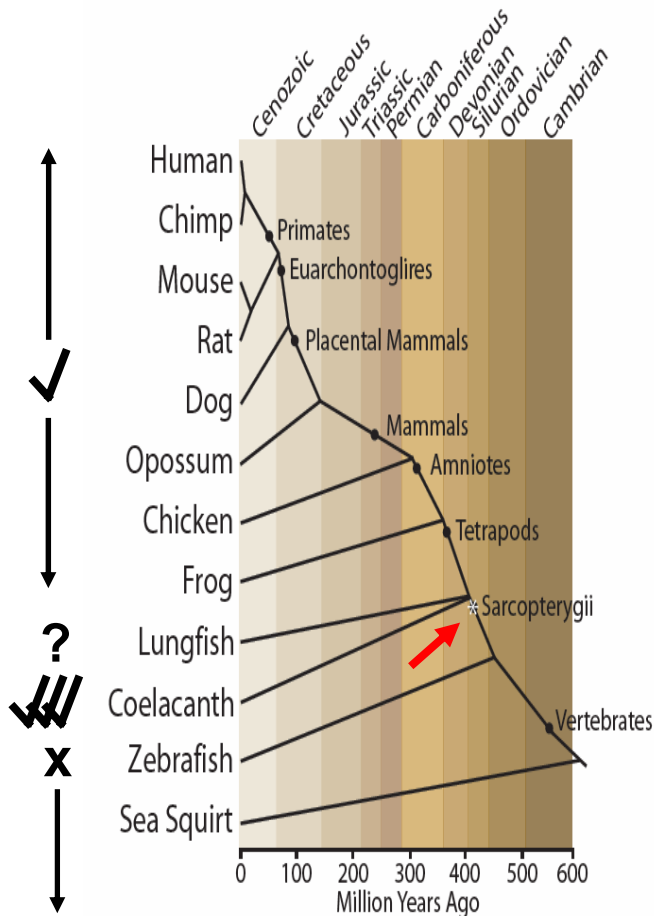
Region Of Homology	BACs	Non-overlapping Length (bp)	uc.338 like Instances
Protocadherin Cluster	AC150310	609,120	21
	AC150283		
	AC150284		
	AC150309		
	AC150308		
Hox Cluster	AC151571	187,392	21
Hox Cluster	AC147788	168,364	4
Genomic	AC140159	92,794	13
<b>Total</b>	<b>8</b>	<b>1,057,670</b>	<b>59</b>



481bp repeat  
59 copies/1Mb  
Highly similar copies  
Unlike any known repeat

# Coelacanth Repeat Distribution

Reconstruct ancestral coelacanth repeat.  
Search all available genomes.



Species	UCSC Assembly	repeat Detected	Species	UCSC Assembly	repeat 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			

Similar distribution in GenBank, Trace Archives.

# Interim Summary

---



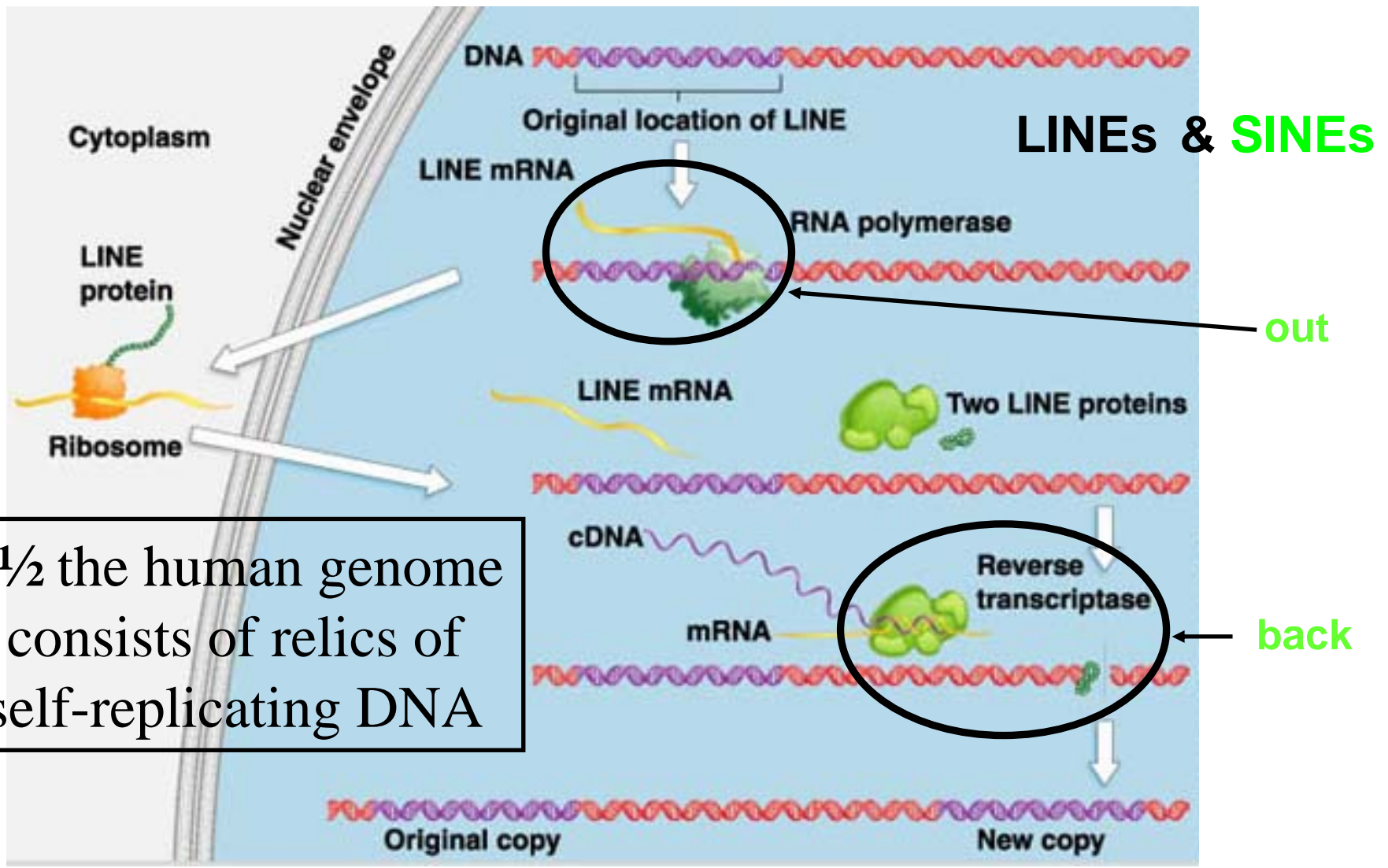
480bp recently active repeat  
80%id to 360bp incl. uc.338  
10<sup>5</sup> copies/genome

---

Unlike any known repeat



# Repeats *are* a Major Force in Vertebrate Evolution

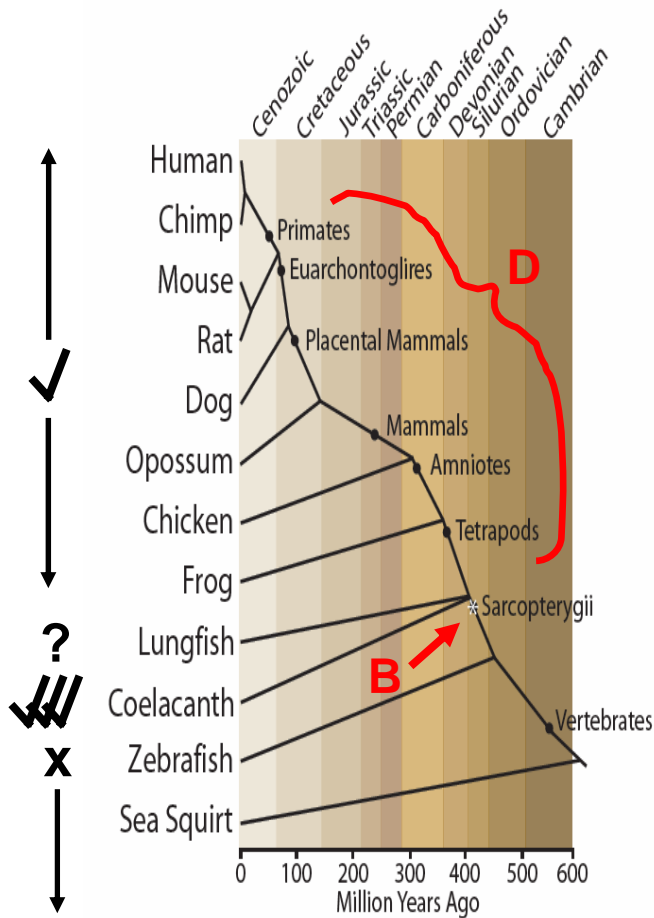


>1/2 the human genome consists of relics of self-replicating DNA



# >360My Old and Going Strong

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

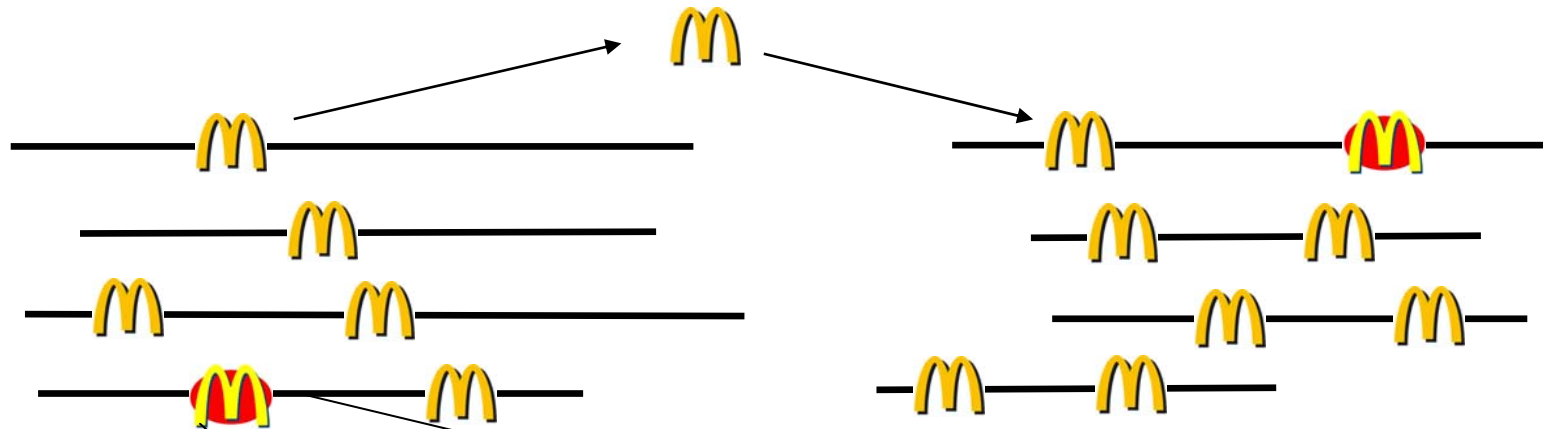


Species	UCSC Assembly	LF-SINE Detected	Species	UCSC Assembly	LF-SINE Detected
<i>Homo sapiens</i>	hg17	Yes	<i>Danio rerio</i>	danRer2	No
<i>Pan troglodytes</i>	panTro1	Yes	<i>Tetraodon nigroviridis</i>	tetNig1	No
<i>Macaca mulatta</i>	rheMac1	Yes	<i>Takifugu rubripes</i>	fr1	No
<i>Mus musculus</i>	mm6	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			



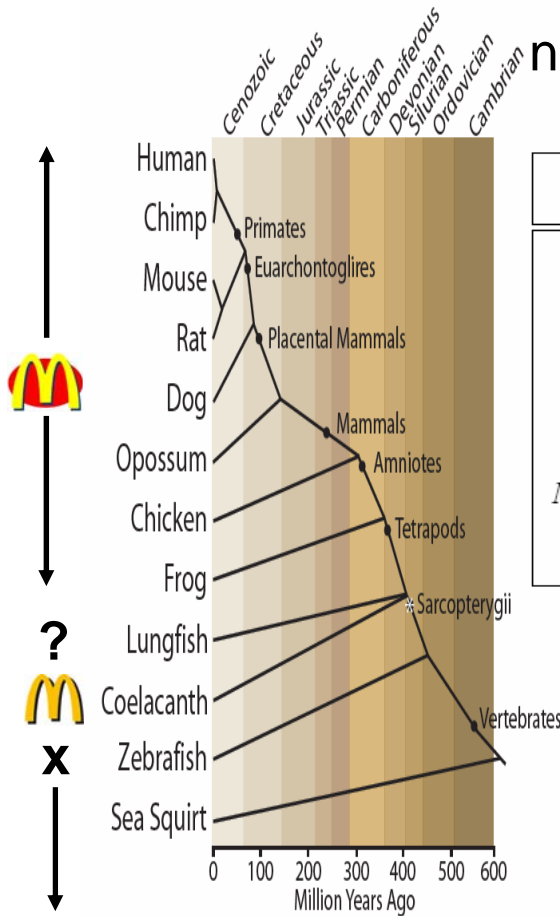
# Co-Option (Exaptation) of 🍔 Mobile Elements

---

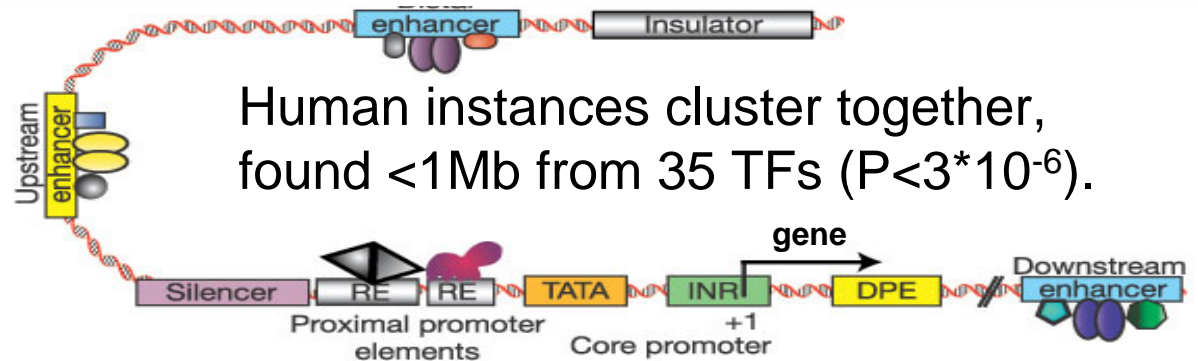


# Exapted Into Which Cellular Roles?

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

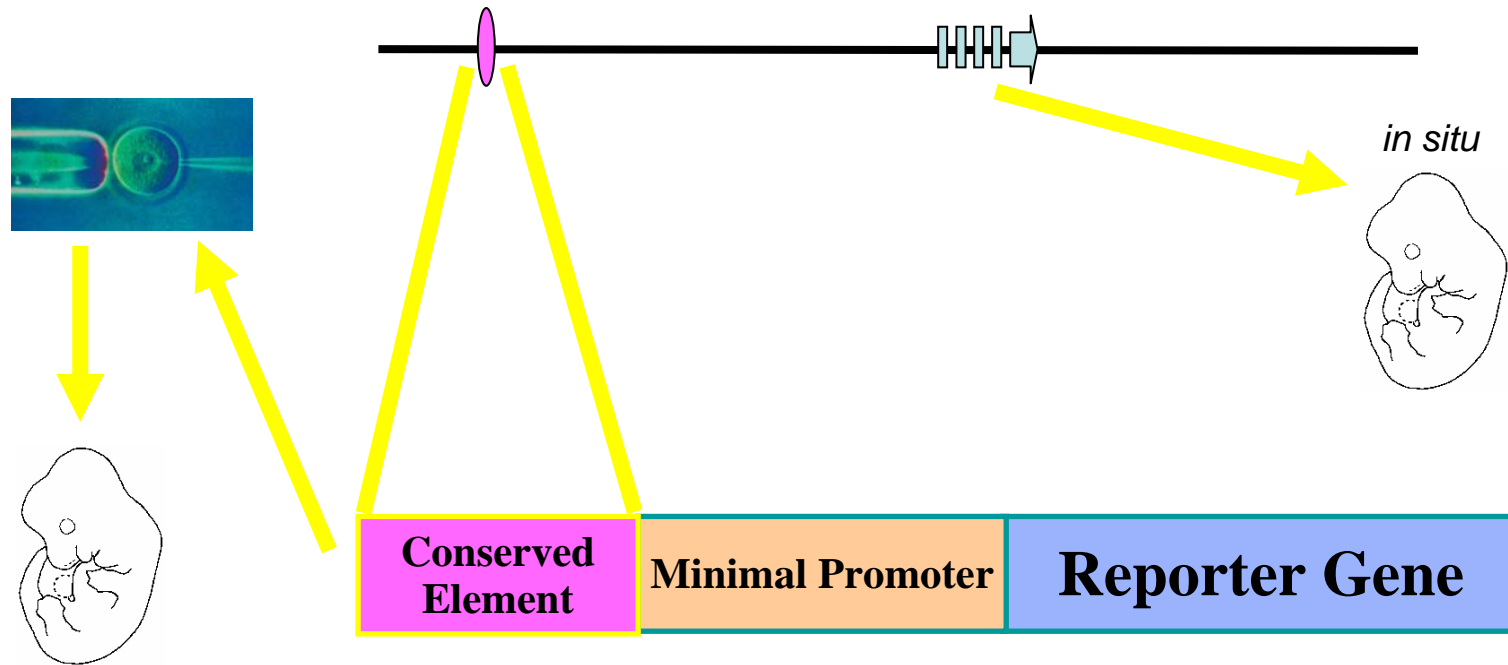


Organism	5' UTR	3' UTR	Exonic		Intronic	Intergenic	Total
			Alt-Spliced	Total			
<i>Homo sapiens</i>	1	0	12	13	68	163	245
<i>Pan troglodytes</i>	-	-	-	-	-	-	210
<i>Macaca mulatta</i>	-	-	-	-	-	-	229
<i>Canis familiaris</i>	-	-	-	-	-	-	235
<i>Bos taurus</i>	-	-	-	-	-	-	169
<i>Mus musculus</i>	0	1	7	8	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



# Transient Transgenic Enhancer Assay

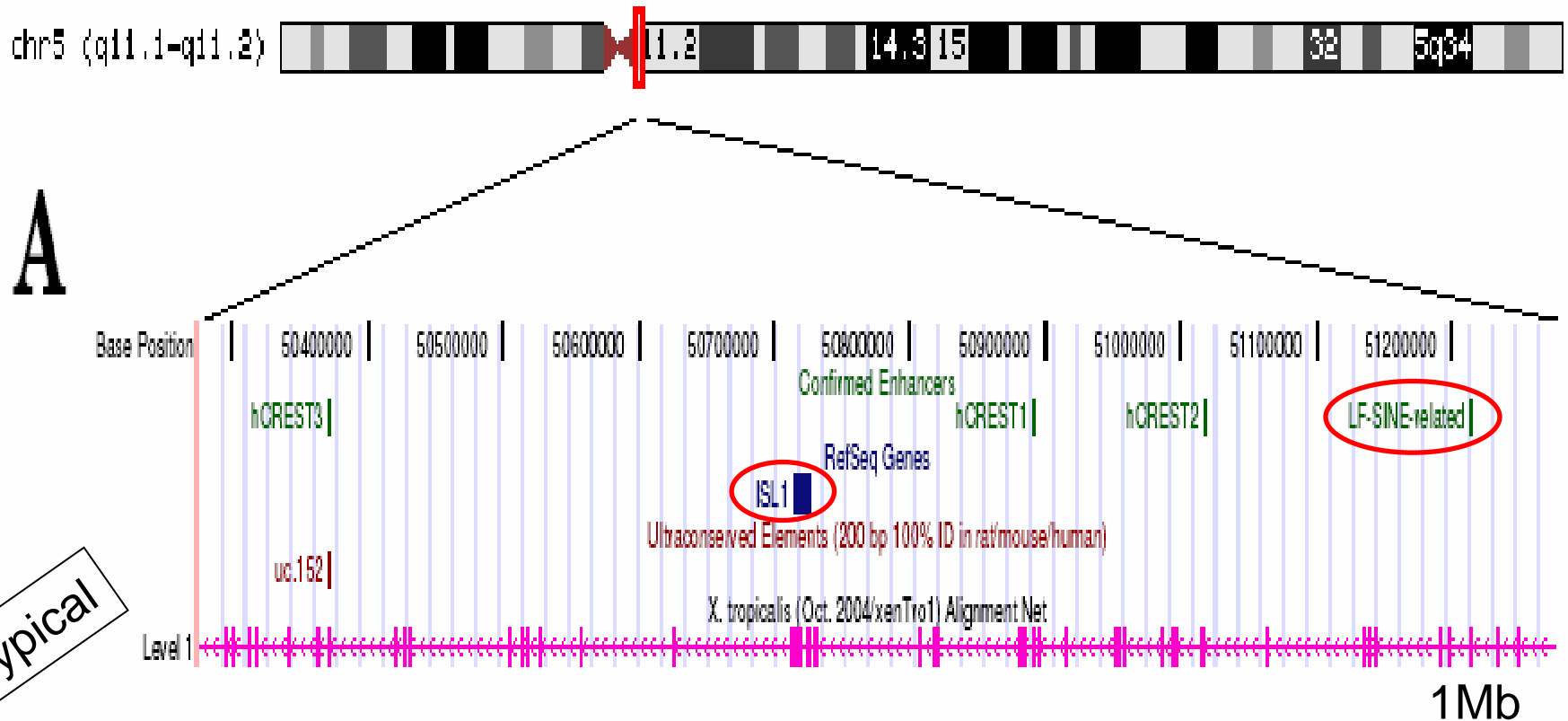
Eddy Rubin's Lab, LBNL



transgenic

Construct is injected into 1 cell embryos  
Taken out at embryonic day 10.5-14.5  
Assayed for reporter gene activity

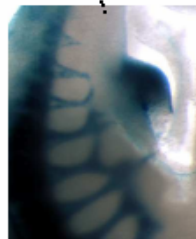
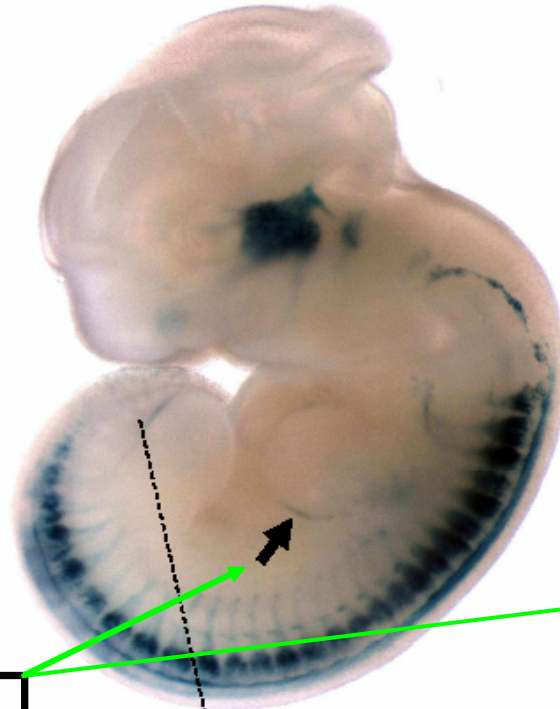
# Instance 500kb Downstream of ISL1



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

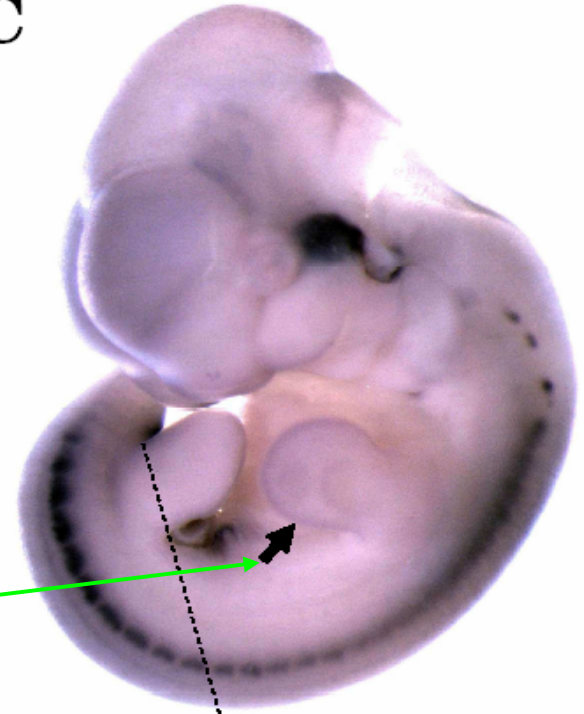
# LF SINE Transgenic (B) vs. Mouse *Isl1 in situ* (C)

**B**



Matched staining  
in dorsal apical  
ectodermal ridge  
(part of limb bud)

**C**

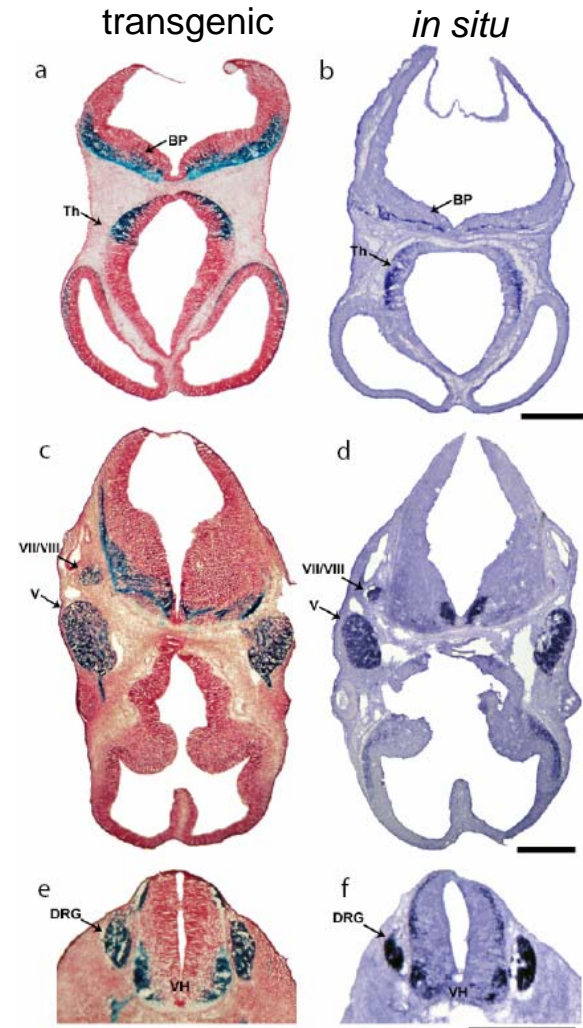


Matched staining in  
genital eminence

Nadav Ahituv, Eddy Rubin

# Matched Level Sections

Corresponding expression patterns in:  
(a, b) the developing thalamus (Th)  
and basal plate (BP) in the brain.  
(c, d) the trigeminal (V) ganglion and  
facio-acoustic (VII/VIII) ganglia  
in the head region.  
(e, f) the dorsal root ganglion (DRG),  
and the lateral region of the  
ventral horn (VH) of the spinal cord  
in thoracic sections.



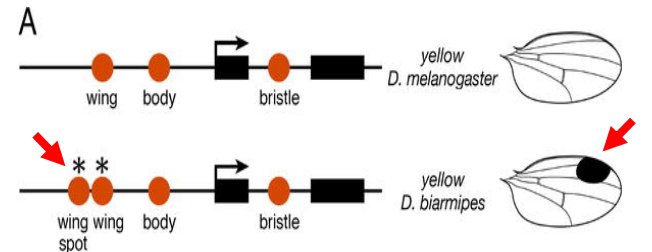
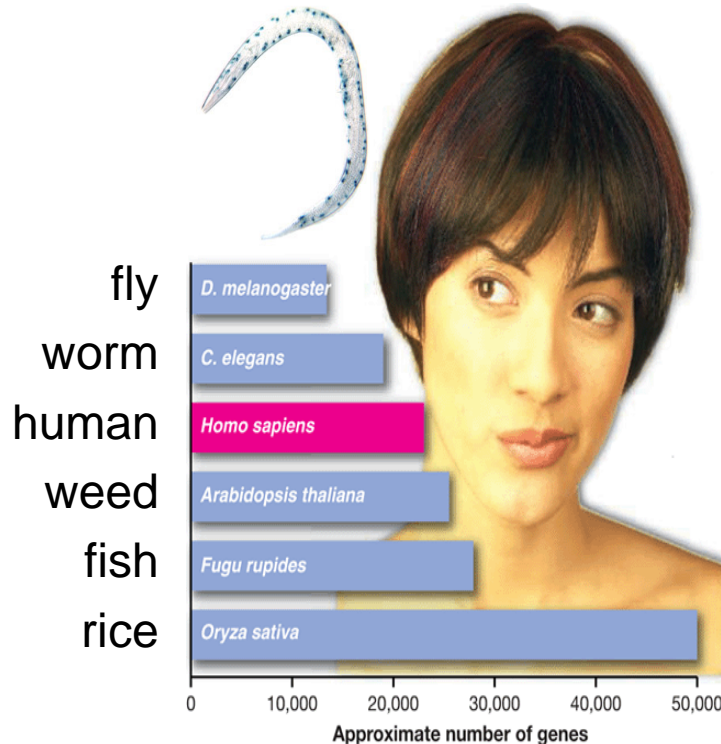
Bryan King, Sofie Salama, Nadav Ahituv, Eddy Rubín

# Mobile Elements Give Birth to Distal Cis-Regulatory Elements: Implications

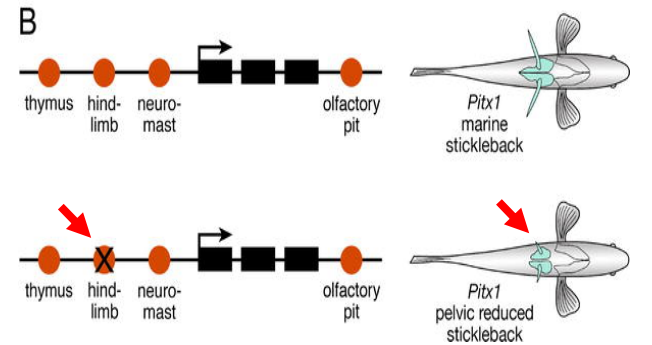
Gene numbers do not correlate with organism complexity.

Many gene families are surprisingly old.

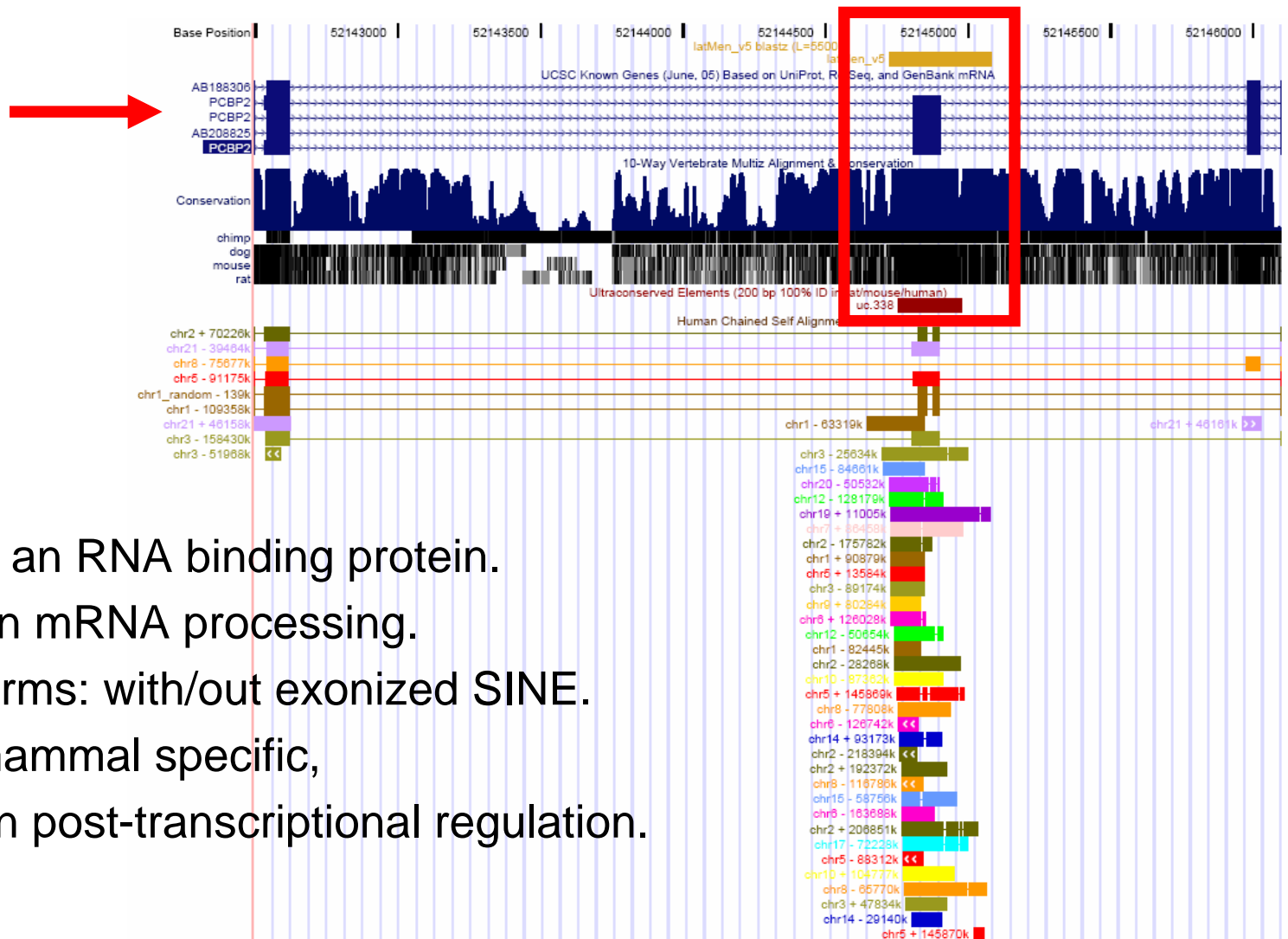
“Regulatory sequence evolution must be the major contribution to the evolution of form.” [Sean Carroll, *PLoS Bio* 2005]



## In/vertebrate Divide

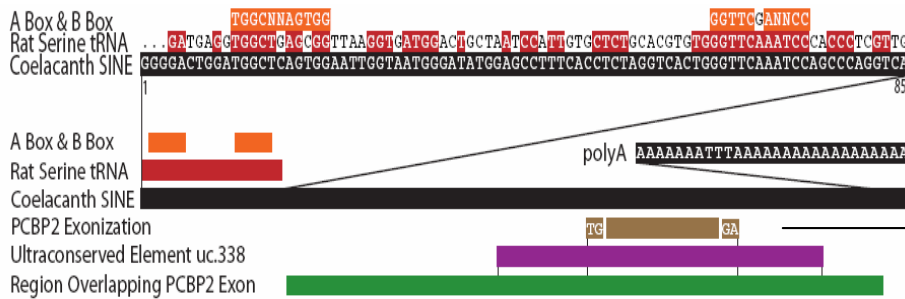


# uc.338 Has Undergone Exonization



PCBP2 is an RNA binding protein.  
Involved in mRNA processing.  
Two Isoforms: with/without exonized SINE.  
Exon is mammal specific,  
involved in post-transcriptional regulation.

# Exonized LF SINE Instances



## 19 different tetrapod exonizations

Gene Name	Sp.	Exon No.	Alt-Spl.	CDS/UTR	Ins. F.S.	Ins. Stop	Trig. NMD	3' S.S. Seq.	ex5'	ex3'	5' S.S. Seq.
PCBP2	H	9	+	C	-	-	-	381 cacagGACAG	289	TAGAGgtgag	
	M		+	C	-	-	-	381 cacagGACAG	289	TAGAGgtgag	
SMARCA4	H	27	+	C	-	-	-	381 cacagAGCAG	289	TAAAGgtgag	
	M		+	C	-	-	-	381 cacagAGCAG	289	TAAAGgtgag	
EEF1B2	H	3	+	C	-	+	+ <sup>a</sup>	376 aacagGTAGA	101	TCATAggtgag	
TCERG1	H	22	+	C	+	+	+ <sup>a</sup>	469 tatagTTAAT	377	AACAGgtaca	
	M		+	C	+	+	+ <sup>a</sup>	469 tatagTTAAT	377	AACAGgtaca	
PTDSR	H	5	+	C	+	+	+	376 aacagATACA	289	TCGGGgtaag	
	M		+	C	+	+	+	376 aacagATACA	289	TCGGGgtaag	
RORA	H	3	+	C	+	+	+	366 cgcagGGCAG	289	AGGTGgtaag	
	M		+	C	+	+	+	366 cgcagGGCAG	289	AGGTGgtaag	
GRID1	H	1	+	5U					289	GTAGGgtaag	
ATF2	H	14	+	C	+	+	+ <sup>a</sup>	450 taaagTGAAT	377	ACCAGgtaca	
	M		+	3U <sup>b</sup>							
FLJ22833	H	4	+	C	+	+	+	381 cacagTCTAG	280	GTAAGgtaat	
ARHGAP6	H	13	+	C	-	+	+	381 cccagAACAA	289	TAAGGgtgag	
	M		+	C	-	+	+	381 cacagAACCA	289	TAAGGgtgag	
KIAA1409	H	34	-	C	-	-	-	381 cacagAACAG	281	GAGAGgttag	
	M		-	C	-	-	-	381 cacagAACAG	281	GAGAGgttag	
	C		-	C	-	-	-	381 cacagAGCAG	281	CAGAGgttag	
NT5C2	H	9	+	C	+	+	+	450 taaagTGAAT	107	GGATGgtaat	
	M		+	C	+	+	+	450 taaagTGAAT	240	GGGAGgtttg	
LRP1B	H	90	+	C	-	-	-	381 tacagGCCAG	289	CTGGGgtgag	
DHX30	H	4	+	C	-	-	-	381 cccagATCGG	289	TCGAGgtaag	
gg-DMTF1	C	12	+	C	+	+	+	450 tttagTGAAT	377	AACAGgtaca	
gg-PPP2R2C	C	2	+	C	+	+	+	381 cacagGAGAG	289	TGGAGgtgag	
gg-SHFM1	C	3	+	3U <sup>b</sup>							
xt-MBNL1	F	4	+	C	-	-	-	381 cacagGCCAG	289	TATGGgtgag	
JGI-49280	F	5	+	C	+	+	+ <sup>a</sup>	450 tctagTGATT	377	AACAGgtaaa	

The 19 proteins are unrelated  
 19/19 exon antisense to SINE  
 17/19 novel coding exon  
 16/17 alternatively-spliced  
 11/17 introduce early stop codon\*  
 post-transcriptional regulation

\* The 6 read-thru exons all match a 93bp internal LF-SINE region, (still!) free of stop codons in all 3 frames (p~0.002)



# Some Challenges

---

- What molecular mechanisms lead to ultra conservation?
- Which ultraconserved elements (ultras) were seeded by transposons? When? Origins of the others?
- How do ultras evolve: all by purifying selection, or is there some reduced mutation or hyper-repair?
- Why the strong association with DNA binding and RNA binding genes? Transcription/splicing enhancers?
- What information do ultras encode and how?
- Are ultras part of delicate self-regulating gene networks that are critical in development?
- Are they related to human disease?
- Are vertebrate and fly ultras related?
- ...



# Kudos

## UC Santa Cruz

David Haussler

Sofie Salama, Craig Lowe, Bryan King

Jim Kent, Adam Siepel, Jakob Pedersen

Katie Pollard, Courtney Onodera

Rachel Harte, Genomics/Browser Group

## Lawrence Berkeley Labs

Eddy Rubin

Nadav Ahituv, Marcelo Nobrega

## Northwestern U.

Jack Kessler, Alex Bassuk

## McGill U.

Mathieu Blanchette

## Penn State U.

Webb Miller's group



## U. Queensland

John Mattick's group

Genome Sequencing Consortia  
All GenBank contributors

