



FACULTY OF ENGINEERING & TECHNOLOGY
DEPARTMENT OF BIOTECHNOLOGY

Dr. Simranjit Singh
Assistant Professor
Department of Biotechnology
Rama University, Kanpur

Comparative Genomics

Lecture 1 - Comparative genomics- Ana Marques

Comparison of DNA sequences.

Lecture 2 - Comparative transcriptomics- Chris Ponting

Comparison of RNA/protein

Lecture 3 - Disease genomics- Caleb Webber

Using genomics to understand phenotype and disease

Practical - Ana Marques and Steve Meader

Using web-based data-mining tools to compare disease associated loci between human and mouse.

Overview:

1-Genome(s);

2-Genomics: comparative, functional and evolutionary;

3-Protein-coding genes and evolution;

The genome

The genome contains all the biological information required to build and maintain any given living organism.

The genome contains the organisms molecular history.

Decoding the biological information encoded in these molecules will have enormous impact in our understanding of biology.



Some history

1866- Gregor Mendel suggested that the traits were inherited.

1869-Friedrich Miescher isolated DNA.

1919-Phoebus Levene identified the nucleotides and proposed they were linked through phosphate groups.

1943- Avery, MacLeod and McCarty showed that DNA and not protein is the carrier of genetic information.

1953- Based on a X-ray diffraction taken by Rosalind Franklin and Raymond Gosling and the Erwin Chargaff discovery that DNA bases are paired James D. Watson and Francis Crick suggested the double helix structure for the DNA.

1957- Crick laid out the central dogma of molecular biology (DNA->RNA->protein).

1961 - Nirenberg and colleagues “cracked” the genetic code

1975- Sanger sequencing

1976/79- First viral genome – MS2/fX174 (chromosomal walking- size ~5 kb)

1982 -First shotgun sequenced genome – Bacteriophage lambda (~50 kb)

1995 - First prokaryotic genome – *H. influenzae*

1996 - First unicellular eukaryotic genome – Yeast

1998 - The first multicellular eukaryotic genome – *C.elegans*

2000 - *Drosophila melanogaster* - fruitfly

2000 - *Arabidopsis thaliana*

2001- Human Genome

~50 years

1865

Mendel discovers laws of genetics

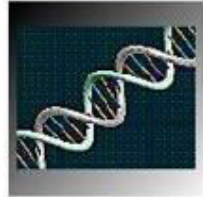


1900

Rediscovery of Mendel's genetics

1944

DNA identified as hereditary material



1953

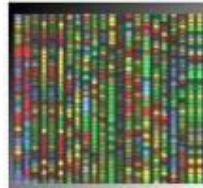
DNA structure

1960's

Genetic code

1977

Advent of DNA sequencing



1975-79

First human genes isolated

1986

DNA sequencing automated

1990

Human genome project officially begins



1995

First whole genome

1999

First human chromosome



2003

'Finished' human genome sequence



The Human genome project promised to revolutionise medicine and explain every base of our DNA.

Large MEDICAL GENETICS focus

Identify variation in the genome that is disease causing

Determine how individual genes play a role in health and disease

This was a huge technical undertaking so further aims of the project were...

- Develop and improve technologies for: DNA sequencing, physical and genetic mapping, database design, informatics, public access
- Genome projects of 5 model organisms e.g. *E. coli*, *S. cerevisiae*, *C. elegans*, *D. melanogaster*, *M. musculus*.



Provide information about these organisms



As test cases for refinement and implementation of various tools required for the HGP

- Train scientists for genomic research and analysis
- Examine and propose solutions regarding ethical, legal and social implications of genomic research (ELSI)

The 2 Human genome project

PUBLIC - Watson/Collins

- Human Genome Project
- Officially launched in 1990
- Worldwide effort - both academic and government institutions
- Assemble the genome using maps
- 1996 Bermuda accord

PRIVATE - Craig Venter

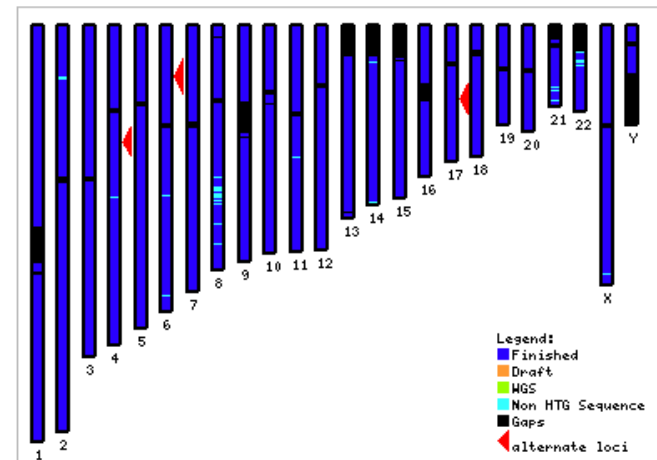
- 1998 Celera Genomics
- Aim to sequence the human genome in 3 years
- 'Shotgun' approach - no use of maps for assembly
- Data release NOT to follow Bermuda principles

It cost 3 billion dollars and took 10 years to complete (5 less than initially predicted).

- Currently 3.2 Gb
- Approx 200 Mb still in progress
 - Heterochromatin
 - Repetitive
- Most recent human genome uploaded February 2009

Finally, it has not escaped our notice that the more we learn about the human genome, the more there is to explore.

“We shall not cease from exploration. And the end of all our exploring will be to arrive where we started, and know the place for the first time.”—T. S. Eliot⁴⁵⁰

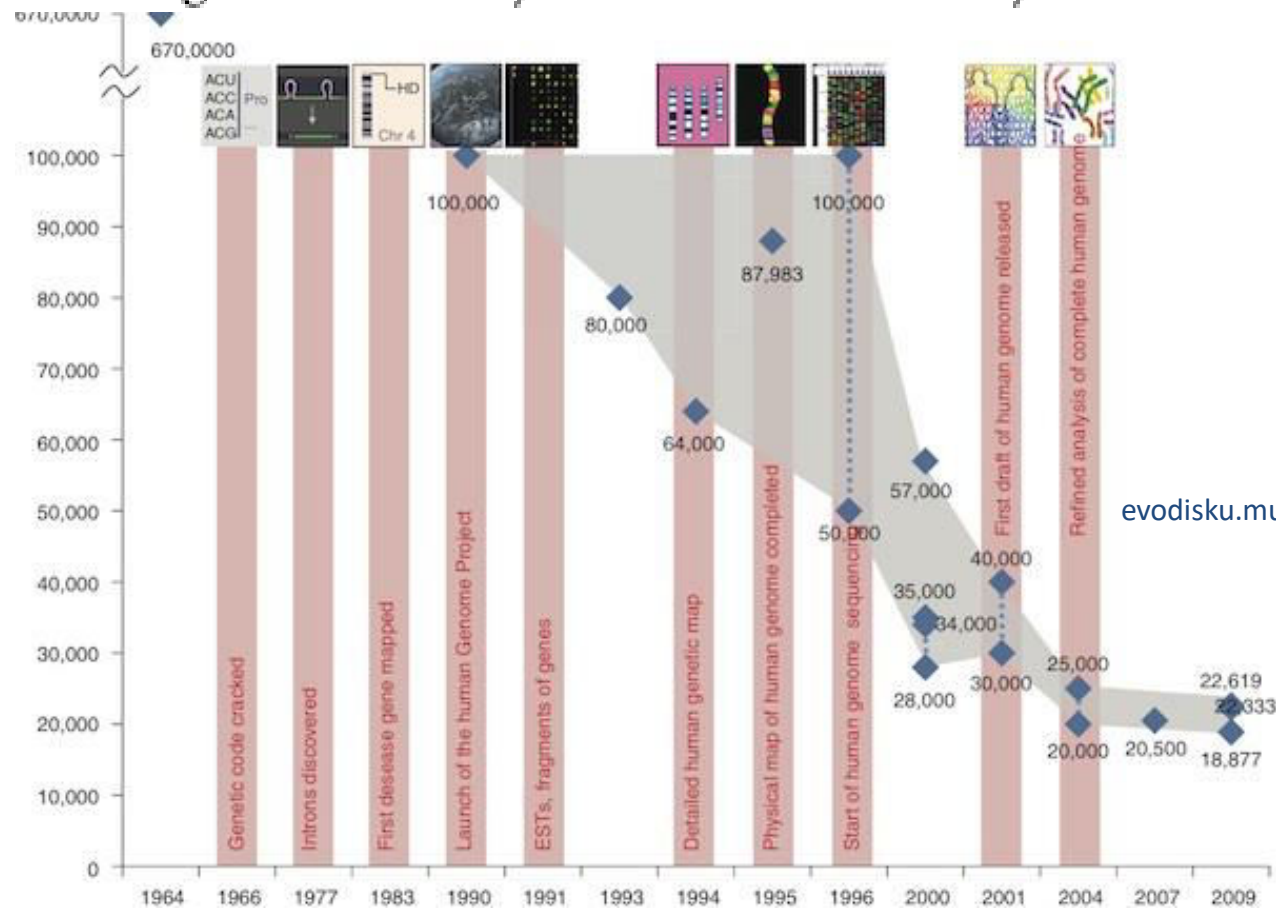


From sequence to function

The scientific program outlined above focuses on how the genome sequence can be mined for biological information. In addition, the sequence will serve as a foundation for a broad range of functional genomic tools to help biologists to probe function in a more systematic manner. These will need to include improved techniques and databases for the global analysis of: RNA and protein expression, protein localization, protein–protein interactions and chemical inhibition of pathways. New computational techniques will be needed to use such information to model cellular circuitry. A full discussion of these important directions is beyond the scope of this paper.

The functional genome

- There appear to be about 30,000–40,000 protein-coding genes in the human genome—only about twice as many as in worm or fly.

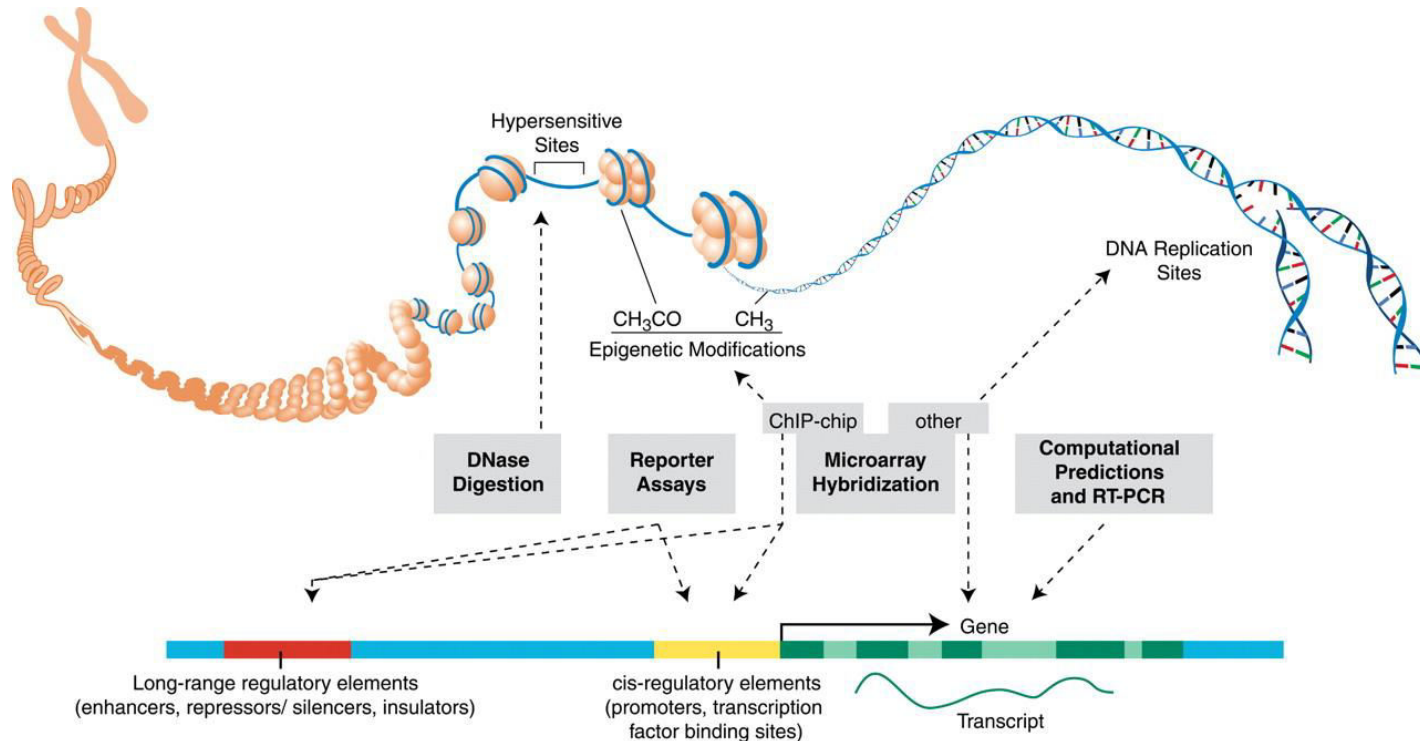


Protein-coding do not explain complexity/diversity.

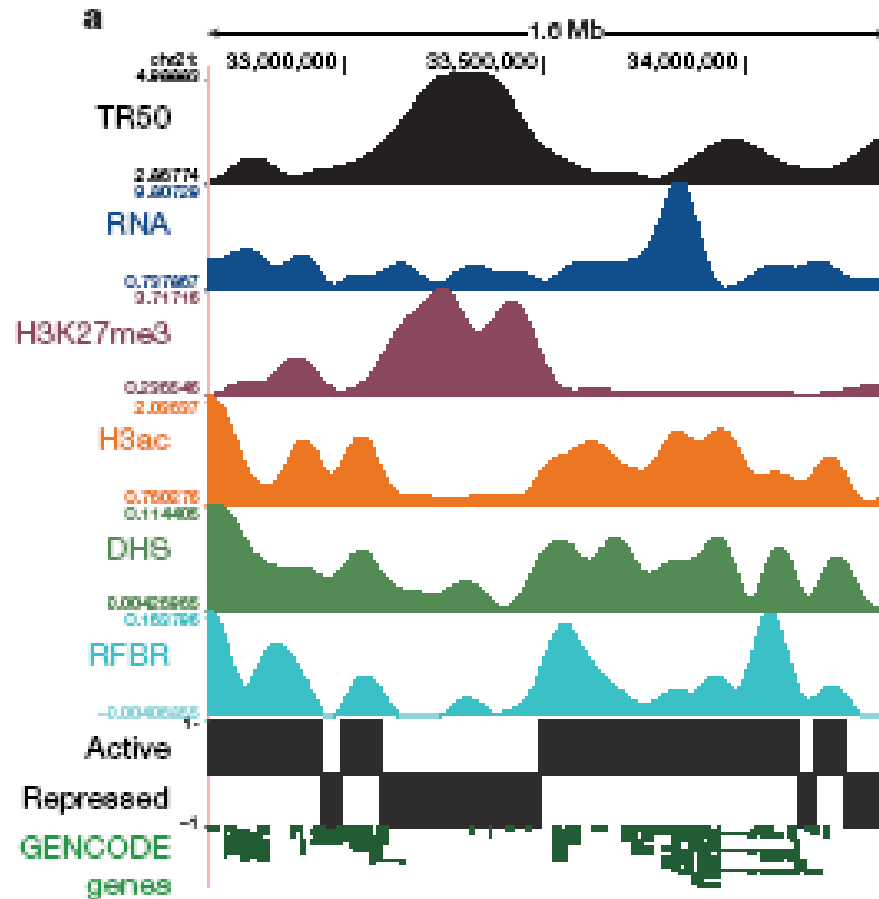
Identification and analysis of functional elements in 1% of the human genome by the ENCODE pilot project

The ENCODE Project Consortium*

35 Research groups threw everything at 30Mb (1%) of human DNA sequence. >200 experimental datasets (transcription, histone-modifications, chromatin structure, regulatory binding sites, replication timing, population variation and more.)



The functional genome map



Estimating the fraction of the genome that is functional

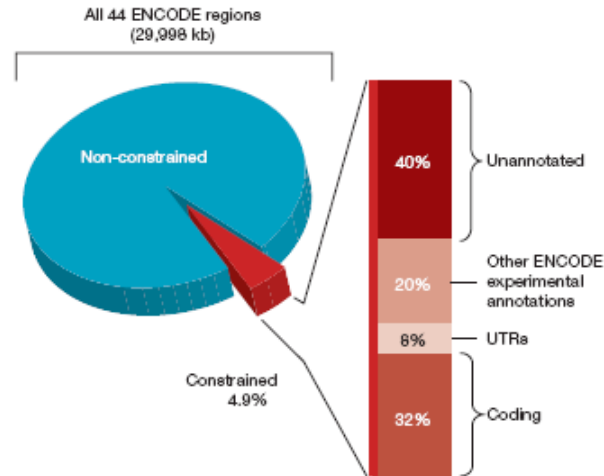


Figure 10 | Relative proportion of different annotations among constrained sequences. The 4.9% of bases in the ENCODE regions identified as constrained is subdivided into the portions that reflect known coding regions, UTRs, other experimentally annotated regions, and unannotated sequence.

- Only about 1.2% of the genome encodes protein sequence
- Most of it is composed of decaying transposons
- 5% appears “constrained” = likely functional
- >70% appears transcribed but unconstrained (lots fast evolving?)

2nd generation sequencing

Genome wide annotation of functional elements made easy!

Platform	Library/ template preparation	NGS chemistry	Read length (bases)	Run time (days)	Gb per run	Machine cost (US\$)	Pros	Cons	Biological applications
Roche/454's GS FLX Titanium	Frag, MP/ emPCR	PS	330*	0.35	0.45	500,000	Longer reads improve mapping in repetitive regions; fast run times	High reagent cost; high error rates in homo- polymer repeats	Bacterial and insect genome <i>de novo</i> assemblies; medium scale (<3 Mb) exome capture; 16S in metagenomics
Illumina/ Solexa's GA _{II}	Frag, MP/ solid-phase	RTs	75 or 100	4 [‡] , 9 [§]	18 [‡] , 35 [§]	540,000	Currently the most widely used platform in the field	Low multiplexing capability of samples	Variant discovery by whole-genome resequencing or whole-exome capture; gene discovery in metagenomics
Life/APG's SOLiD 3	Frag, MP/ emPCR	Cleavable probe SBL	50	7 [‡] , 14 [§]	30 [‡] , 50 [§]	595,000	Two-base encoding provides inherent error correction	Long run times	Variant discovery by whole-genome resequencing or whole-exome capture; gene discovery in metagenomics

Applications

1-Genome sequencing and genome assembly (Panda genome, 2009)

2-Genome re-sequencing (Craig Venter, James Watson...1000 genomes project)

3- Transcriptome sequencing (unbiased)

4- Metagenomics

5-ChIP-seq

7-RIP-seq

...seq.

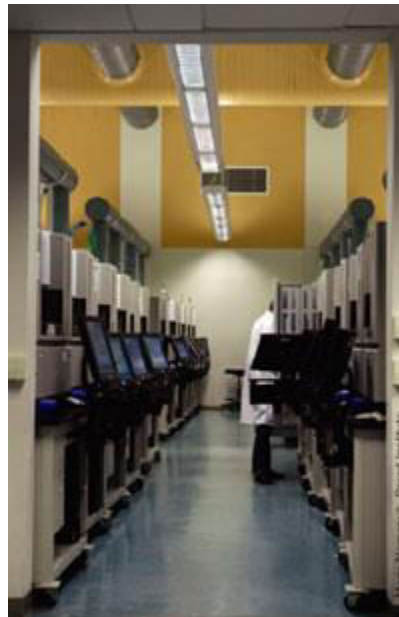


DO SOMETHING EXCITING

Get off your ass and engage your passions.

Single molecule sequencing.

Potential to answer questions that remain open (somatic variation/ single cell transcription...)



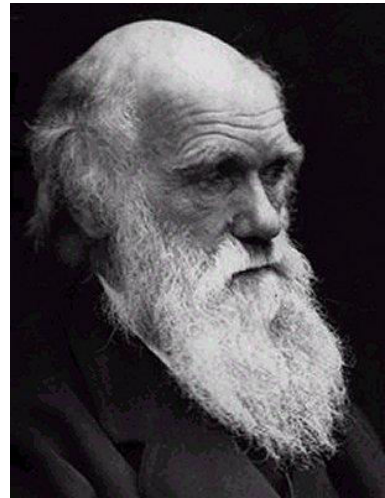
**Next generation sequencing has (and will continue to) changed the way we do and understand biology!
More data but what should we do with it?**

How we use this data to understand physiology, behaviour, disease and variation between species/individuals we need to:

- The evolutionary history of every genetic element (every base)
- Evolutionary forces shaping the genome
- Structural and sequence variation in the population and between species.

Comparative genomics studies differences between genome sequences pin-pointing changes over time. Comparison of the number/type changes against the background “neutral” expected changes provides a better understanding of the forces that shaped genomes and traits.

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



MUTATION

1. Small scale mutations

Nucleotide substitutions

Small Insertions / Deletions (Indels)

ACGTGTC → **ATGTGTC**

ACGTGTC → **AGTGTC**

MUTATION

1. Small scale mutations

Nucleotide substitutions

ACGTGTC → **ATGTGTC**

Small Insertions / Deletions (Indels)

ACGTGTC → **AGTGTC**

2. Large scale mutations (> 1kb)

QuickTime™ and a decompressor are needed to see this picture.

QuickTime™ and a decompressor are needed to see this picture.

How do changes accumulate in the genome?

In 1965 Pauling and colleagues showed that for any given protein the rate of molecular evolution is approximately constant in all lineages.

QuickTime™ and a
decompressor
are needed to see this picture.

1968, proposed that most mutations
accumulated in genomes are neutral.

QuickTime™ and a
decompressor
are needed to see this picture.

The Neutral Theory.

Aim: Identify regions of the genome that are not evolving neutrally!

LOCI X-
Neutral

	↓	↓	↓↓	↓	↓	
Species 1	CGACATTA	AATAGGCGC	AGGACCAG	ATACCAG	ATCAAAGC	CAGGCGCA
Species 2	CGACGTTA	AATTGGCGC	AGTATCAG	ATACCCG	ATCAAAGC	CAGACGCA

LOCI Y

	↓					↓		
Species 1	CATGGG	TCATCA	CTCTAG	CTGTAC	GTCTACT	TCATCAT	CGCGCT	TACG
Species 2	CATGAG	TCATCA	CTCTAG	CTGTAC	GTCTACT	TCATCAT	CGCGTT	TACG

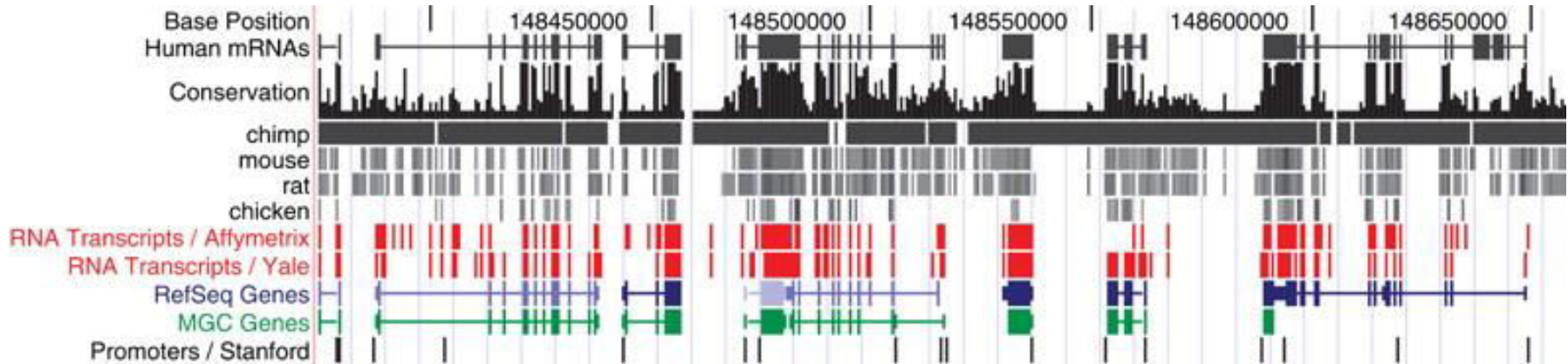
Sequence that is conserved over long evolutionary distances is likely to be under selective constraint

Conservation is often a good predictor of functionality

Conservation highlights exons

Regulatory Element?

Novel exon?

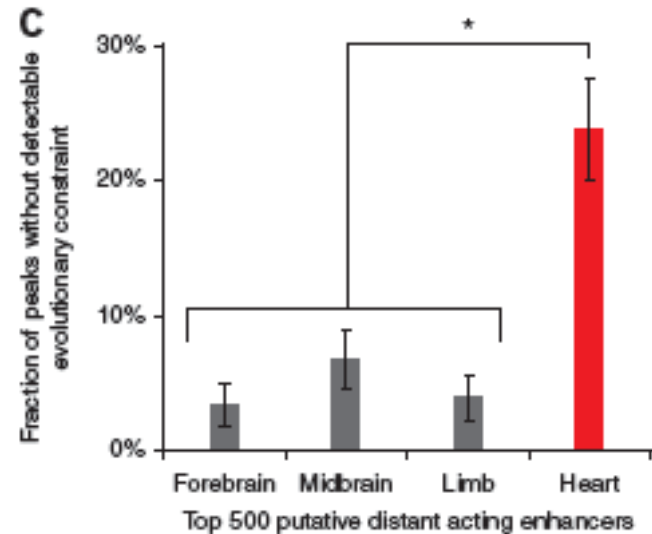
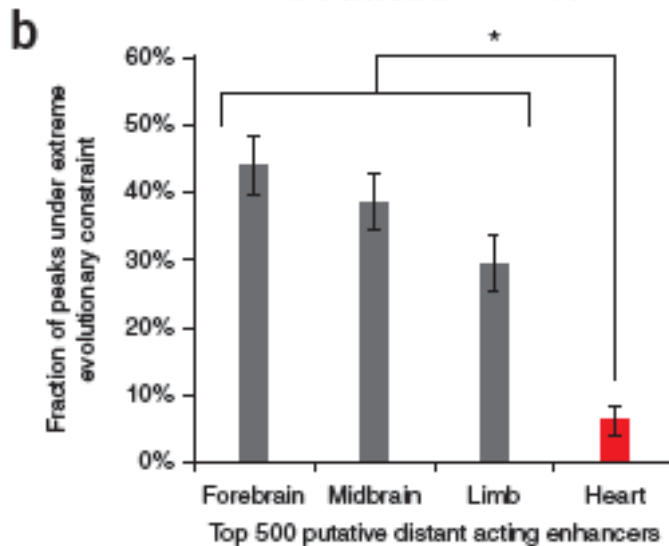
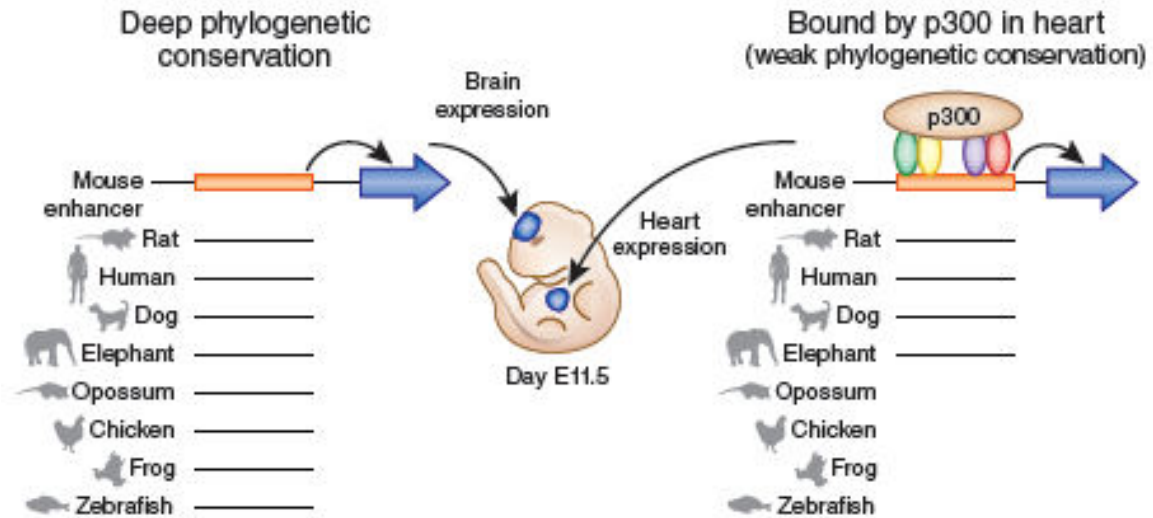


BUT...

Conservation is not synonymous of function

Not all functional sequence is conserved across long evolutionary distance.

Heart Enhancers

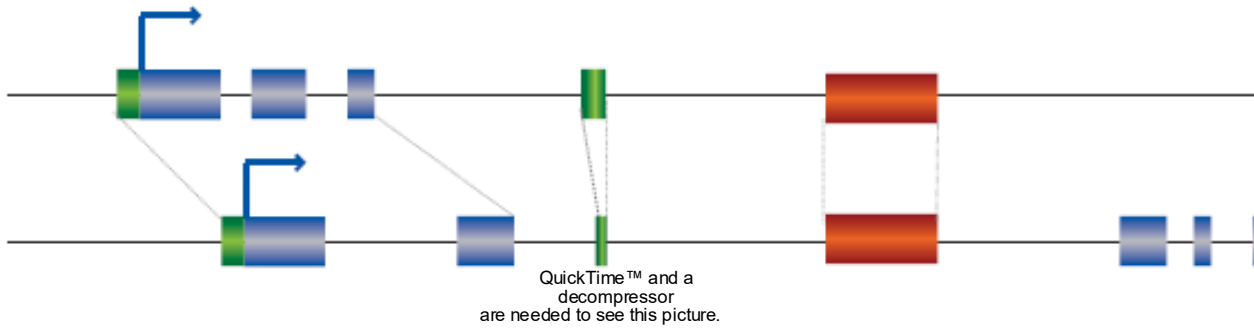


Conservation is not synonymous of function

Long Intergenic ncRNA

QuickTime™ and a decompressor are needed to see this picture.

QuickTime™ and a decompressor are needed to see this picture.

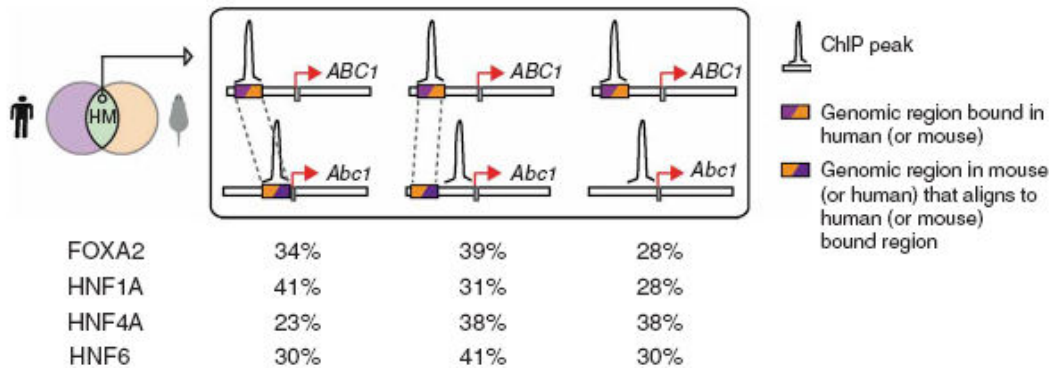


Sequence conservation doesn't imply function conservation

Despite conservation of binding preferences and binding sites only a small proportion of TF binding events is conserved across species

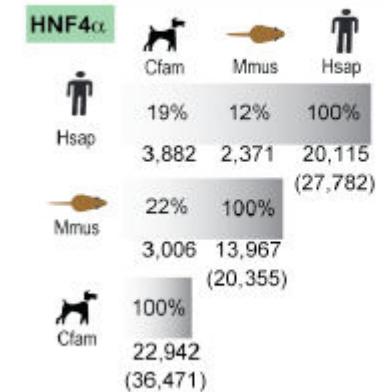
a	Regulator	PFAM category	HS bound	MM bound	Intersection	P value	HS binding sequence	MM binding sequence
	FOXA2	Forkhead	151	574	68	1.0E-45		
	HNF1A	POU-homeodomain	251	224	45	1.0E-29		
	HNF4A	Nuclear receptor	1,251	654	387	1.0E-136		
	HNF6	CUT-homeodomain	157	324	41	1.0E-27		

c

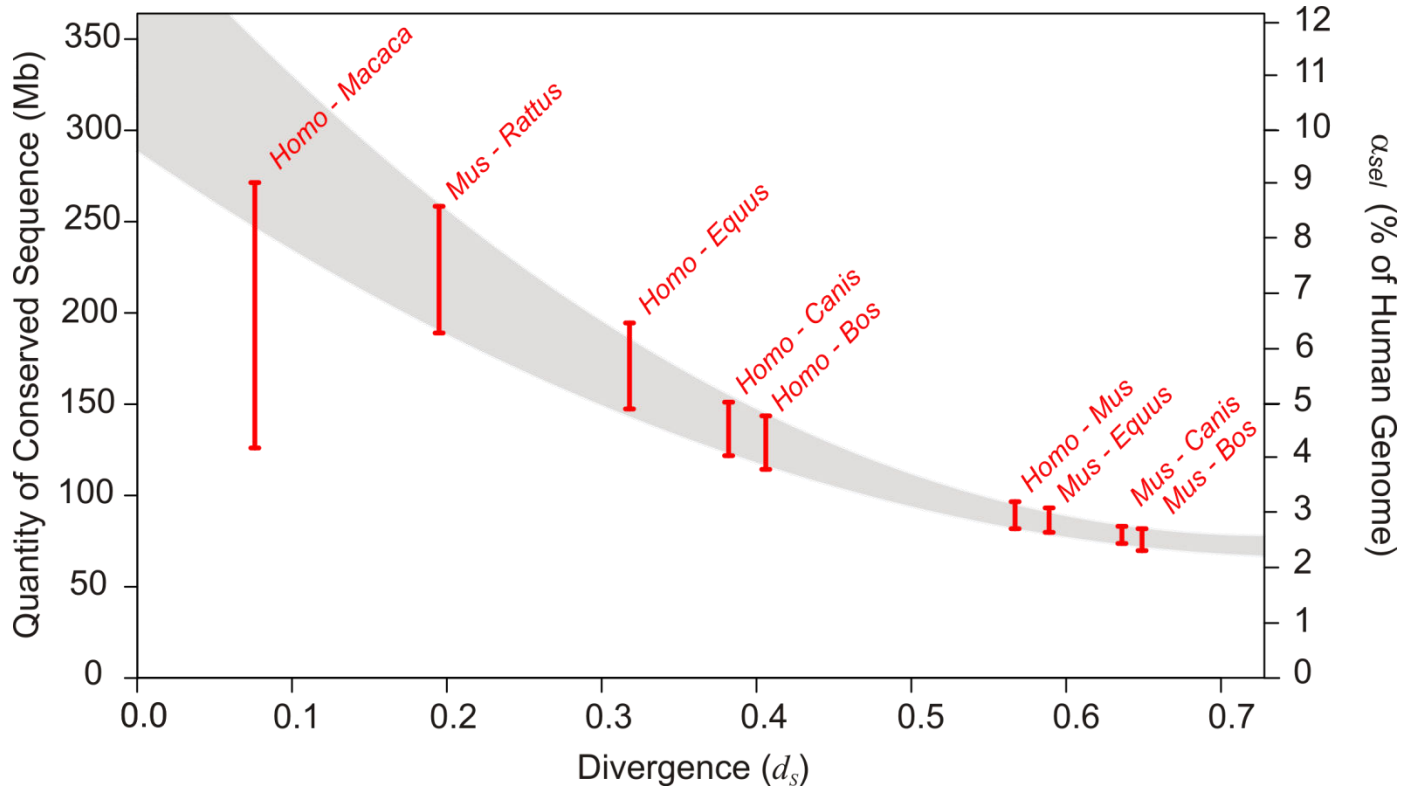


A

Stringent SWEMBL cutoff (R = 0.01)



Sequence conservation doesn't imply function conservation



Massive turnover of functional sequence in mammalian genomes

Lessons from comparative genomics: Changes of protein coding repertoires and contributions to phenotypic differences

