

Current Topics in Genome Analysis  
October 21, 2003

# Computational Techniques in Comparative Genomics



Elliott H. Margulies, Ph.D.  
Genome Technology Branch  
National Human Genome Research Institute  
elliott@nhgri.nih.gov

## Outline

- Fundamental concepts of comparative genomics
- Alignment and visualization tools
- Information available through genome browsers
- Gene prediction and identification
- Identifying regulatory sequences
- Insights from Human-Mouse sequence comparisons
- Multi-species sequence analysis

# Genome Wide Sequence Data

'Finished'



Human

Draft assemblies on the way



Dog



Chicken



Tetraodon



Chimpanzee



Zebrafish

Draft assemblies Available



Pufferfish



Mouse



Sea Squirt



Rat

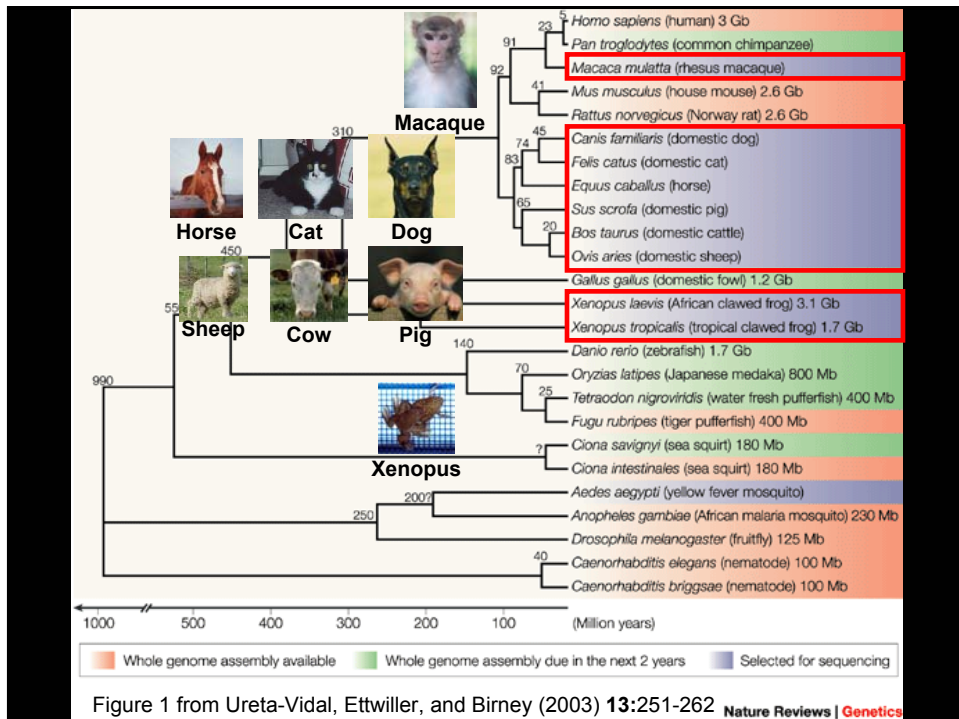


Figure 1 from Ureta-Vidal, Ettwiller, and Birney (2003) 13:251-262 Nature Reviews | Genetics

---

***What can be more curious than that the hand of a man, formed for grasping, that of a mole for digging, the leg of the horse, the paddle of the porpoise, and the wing of the bat, should all be constructed on the same pattern, and should include similar bones, in the same relative positions?***

*Charles Darwin, The Origin of Species (1859)*

---

**Similarity in Morphology**



**Similarity in Genetics**

## **Comparative Genomics**

- Find sequences that have diverged less than we expect  
*These sequences are likely to have a functional role*
- Our expectation is related to the time since the last common ancestor



Human



Chimpanzee



Horse



Rat



Platypus



Zebrafish

***Evolutionary Distance***

# Comparative Genomics

*Similarity in Identity*



*Similarity in Function*

My **name** is **Elliott**. I am **presenting** a **seminar** about **comparative genomics**.

Mon **nom** est **Elliott**. Je vous **présente** un **séminaire** à propos de **génomique comparative**.

Meine **name** ist **Elliott**. Ich **präsentiere** ein **Seminar** über **Comparativ Genomics**.

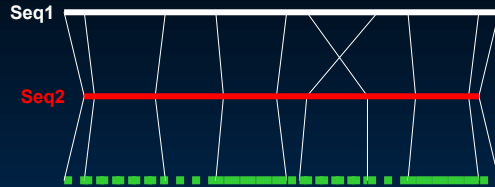
## Outline

- Fundamental concepts of comparative genomics
- Alignment and visualization tools
- Information available through genome browsers
- Gene prediction and identification
- Identifying regulatory sequences
- Insights from Human-Mouse sequence comparisons
- Multi-species sequence analysis



# PipMaker vs. VISTA

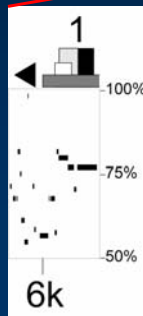
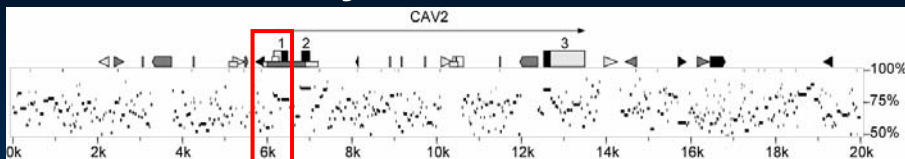
- Visualization
- Alignment Strategy
  - VISTA: *avid*
  - PipMaker: *blastz*
- East Coast – West Coast



## PipMaker

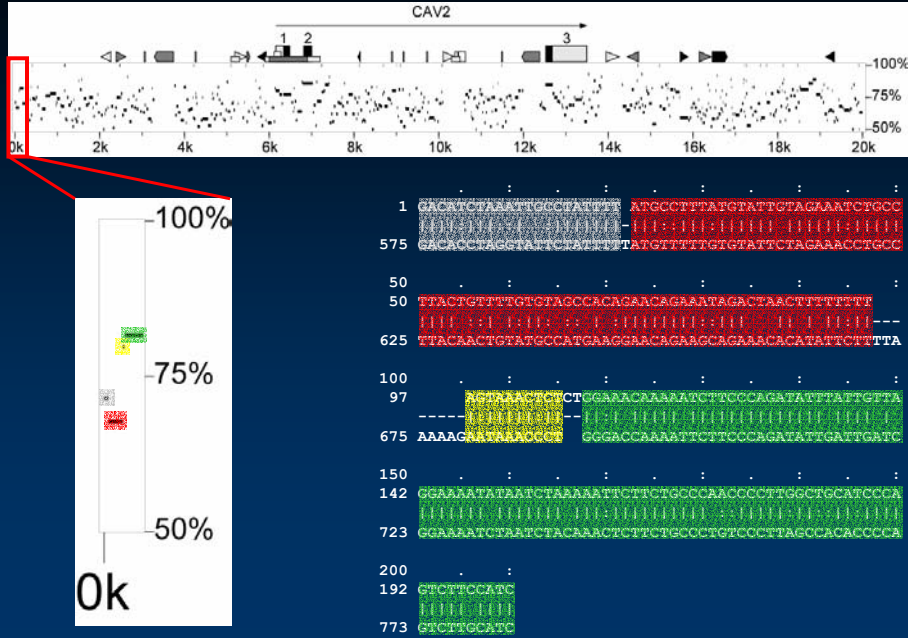
<http://bio.cse.psu.edu/pipmaker/>

- Percent Identity Plot



- X-axis is the reference sequence
- Horizontal lines represent gap-free alignments

<http://bio.cse.psu.edu/pipmaker/>

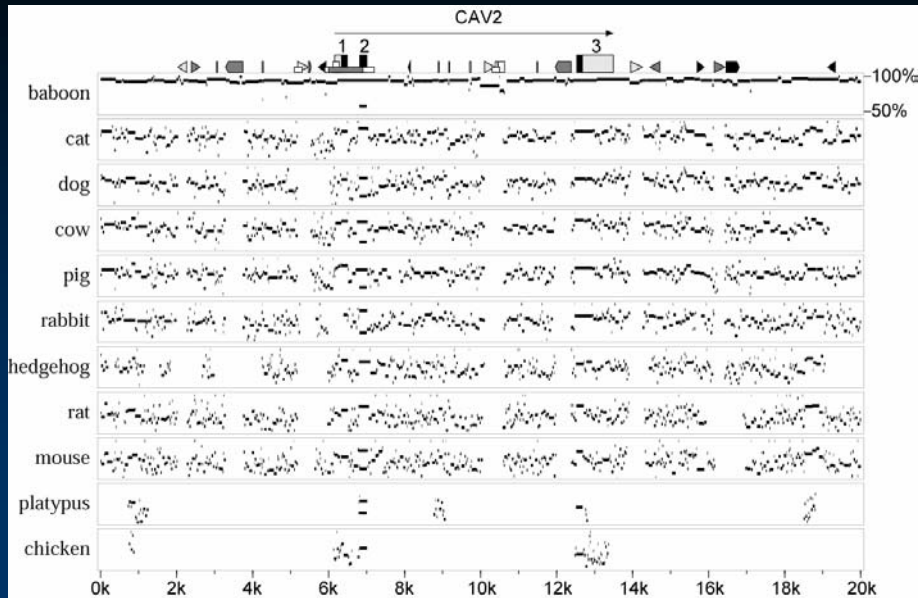


## PipMaker

<http://bio.cse.psu.edu/pipmaker/>

- **P**ercent **I**dentify **P**lot
- Available in 3 Flavors:
  - **Regular**
    - No Additional Options
  - **Advanced** [go to submission page](#)
    - Different Alignment Strategies
    - Additional Output Summaries
  - **MultiPipMaker**
    - Multi-species display

# MultiPipMaker



# PipTools

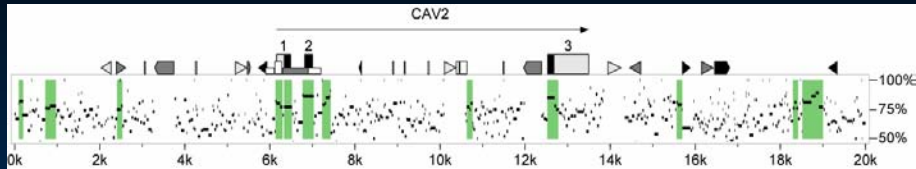
- “Show me all alignments that are at least X% identity and Y% length”
  - `strong-hits`
- Coordinate Conversion
  - `transform-pos`
  - `shift-pos`
  - `where-hit`
- Coordinate Extraction from GenBank Files
  - `genbank2exons`
  - `genbank2repeats`
- Laj – Interactive Alignment Viewer – [Open Laj](#)



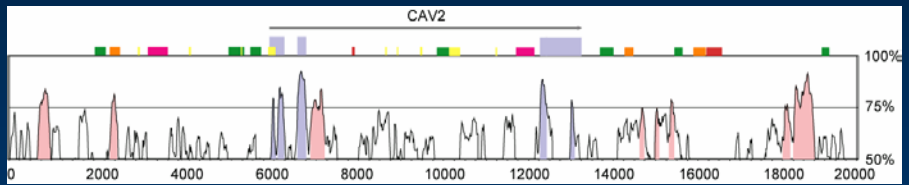


# What's Your Preference?

## PipMaker



## VISTA



## Summary of Alignment Tools

- PipMaker (`blastz`)
- VISTA (`avid`)
- Lagan and mLagan (glocal alignments)  
– <http://lagan.stanford.edu/>
- **Box 1** from:

Ureta-Vidal, Ettwiller, and Birney (2003) Comparative Genomics: Genome-Wide Analysis in Metazoan Eukaryotes *Nature Reviews Genetics* 4: 251-262

## Outline

- Fundamental concepts of comparative genomics
- Alignment and visualization tools
- Information available through genome browsers
- Gene prediction and identification
- Identifying regulatory sequences
- Insights from Human-Mouse sequence comparisons
- Multi-species sequence analysis

## Genome Browsers

UCSC Genome Bioinformatics  
<http://genome.ucsc.edu>

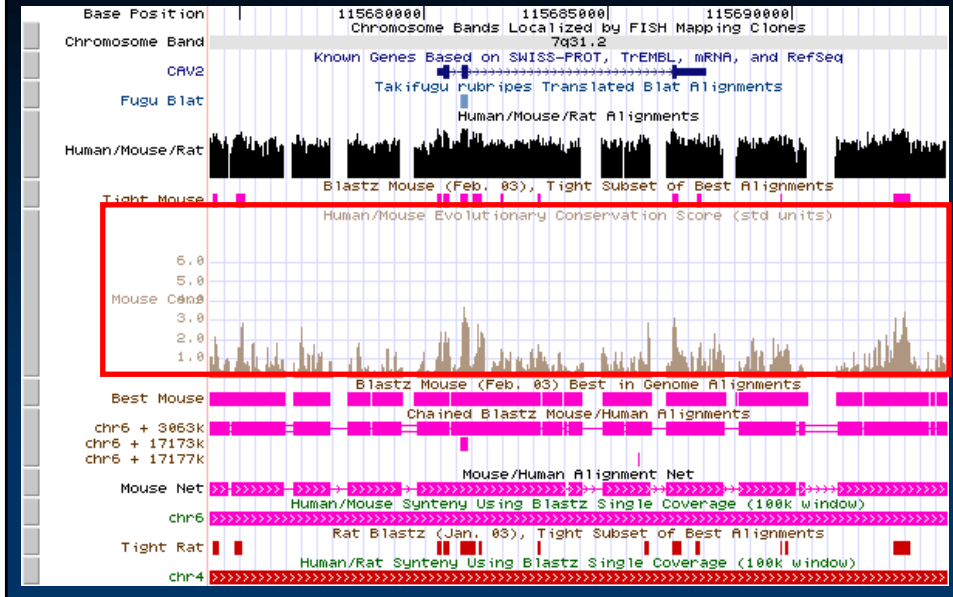
 project **Ensembl**

<http://www.ensembl.org>

 NCBI Map Viewer

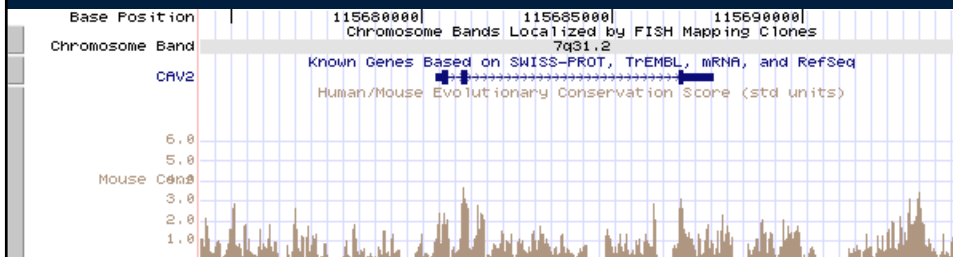
<http://www.ncbi.nlm.nih.gov/mapview/>

# Comparative Sequence Tracks



# Mouse Conservation Score at UCSC

Probability that the observed conservation would occur by chance under *Neutral Evolution*



## Neutral Evolution

- No selective pressure/advantage to keep or change the DNA sequence
- Rate of variation should correlate with:
  - Mutation rate
  - Amount of time since the last common ancestor
- The neutral rate can vary across the genome

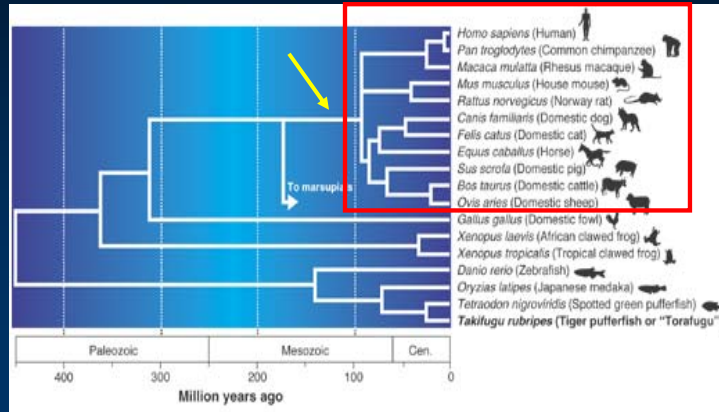
## Types of Neutrally Evolving DNA

- 4-Fold Degenerate Sites
  - Third position of codons which can be any base and code for the same amino acid

First	Second				Last
	U	C	A	G	
U	Phe	Ser	Tyr	Cys	U
	Phe	Ser	Tyr	Cys	C
	Leu	Ser	Stop	Stop	A
	Leu	Ser	Stop	Trp	G
C	Leu	Pro	His	Arg	U
	Leu	Pro	His	Arg	C
	Leu	Pro	Gln	Arg	A
	Leu	Pro	Gln	Arg	G
A	Ile	Thr	Asn	Ser	U
	Ile	Thr	Asn	Ser	C
	Ile	Thr	Lys	Arg	A
	Met	Thr	Lys	Arg	G
G	Val	Ala	Asp	Gly	U
	Val	Ala	Asp	Gly	C
	Val	Ala	Glu	Gly	A
	Val	Ala	Glu	Gly	G

# Types of Neutrally Evolving DNA

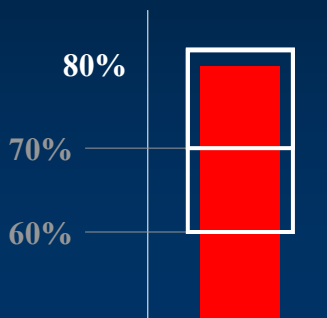
- **Ancestral Repeats**
  - Ancient Relics of Transposons Inserted Prior to the Eutherian Radiation



## Mouse Conservation Score at UCSC

Probability that the observed conservation would occur by chance under **Neutral Evolution**

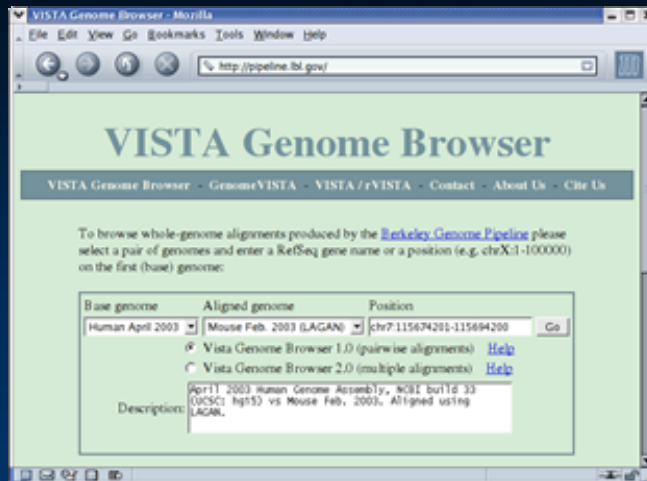
- Calculated from 50-base windows
- Score is weighted on the surrounding neutral rate



L-Score	Frequentist Probability	Bayesian Probability
1	0.1	0.32
2	0.01	0.75
3	0.001	0.94
4	0.0001	0.97
5	0.00001	0.98
6	0.000001	0.99

# Berkley Genome Pipeline

<http://pipeline.lbl.gov/>



Link to [Static View](#)

## Analysis of Comparative Sequence Data

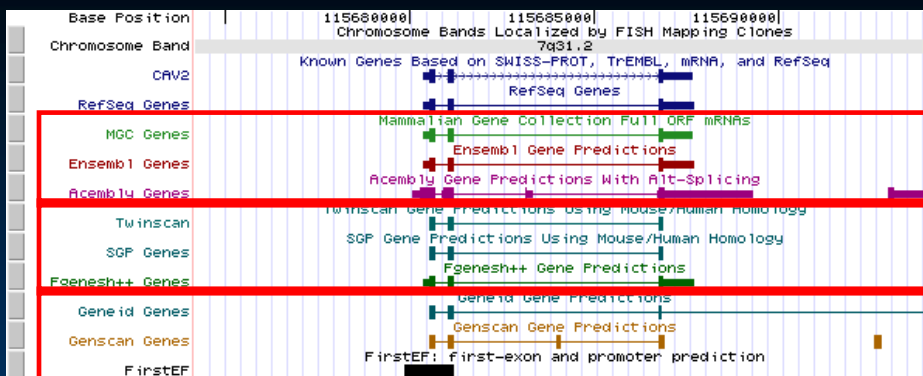
Sequence Conservation  Identification of Functional Elements

- Coding Sequences (i.e. Genes)
  - Relatively EASY to identify
  - Basic understanding of the 'language'
  - Complementary datasets available (ESTs, cDNAs)
- Non-Coding Functional Sequences
  - HARD to identify
  - Very little idea of what to look for
  - Virtually no complementary datasets

## Outline

- Fundamental concepts of comparative genomics
- Alignment and visualization tools
- Information available through genome browsers
- **Gene prediction and identification**
- Identifying regulatory sequences
- Insights from Human-Mouse sequence comparisons
- Multi-species sequence analysis

## Approaches to Gene Prediction



- **Evidence-Based**
  - MGC
  - Acembly
  - Ensembl
- **Ab Initio**
  - Genscan
  - Geneid
  - FirstEF
- **Dual-Genome**
  - Twinscan
  - SGP
  - Fgenesh++



## Additional Gene Prediction Resources

- Fugu BLAT Track at UCSC
- EXOFISH – <http://www.genoscope.cns.fr/proxy/cgi-bin/exofish.cgi>
- SLAM – <http://baboon.math.berkeley.edu/~syntenic/slam.html>
  - Cawley et al. (2003) *Nucleic Acids Research* 31:3507-3509
- **Box 1** from:
  - Ureta-Vidal et al. (2003) *Nature Reviews Genetics* 4:251-262

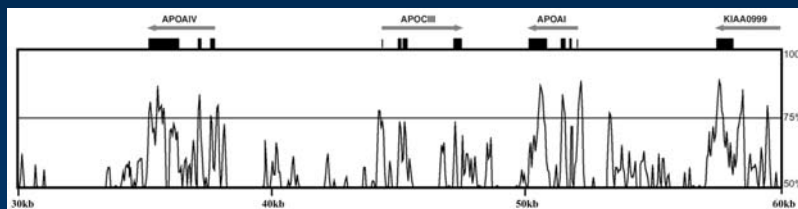
## Gene Finding based on Sequence Conservation

*Science* 294:169-173

### An Apolipoprotein Influencing Triglycerides in Humans and Mice Revealed by Comparative Sequencing

Len A. Pennacchio,<sup>1</sup> Michael Olivier,<sup>2\*</sup> Jaroslav A. Hubacek,<sup>3</sup>  
Jonathan C. Cohen,<sup>3</sup> David R. Cox,<sup>2</sup> Jean-Charles Fruchart,<sup>4</sup>  
Ronald M. Krauss,<sup>1</sup> Edward M. Rubin<sup>1†</sup>

### APOAI/CIII/AIV Gene Cluster



Adapted from Figure 1, Pennacchio et al. *Science* 294:169-173

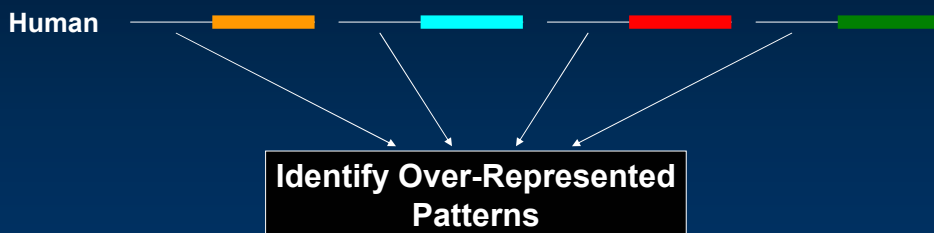
## Outline

- Fundamental concepts of comparative genomics
- Alignment and visualization tools
- Information available through genome browsers
- Gene prediction and identification
- **Identifying regulatory sequences**
- Insights from Human-Mouse sequence comparisons
- Multi-species sequence analysis

## Motif Finding

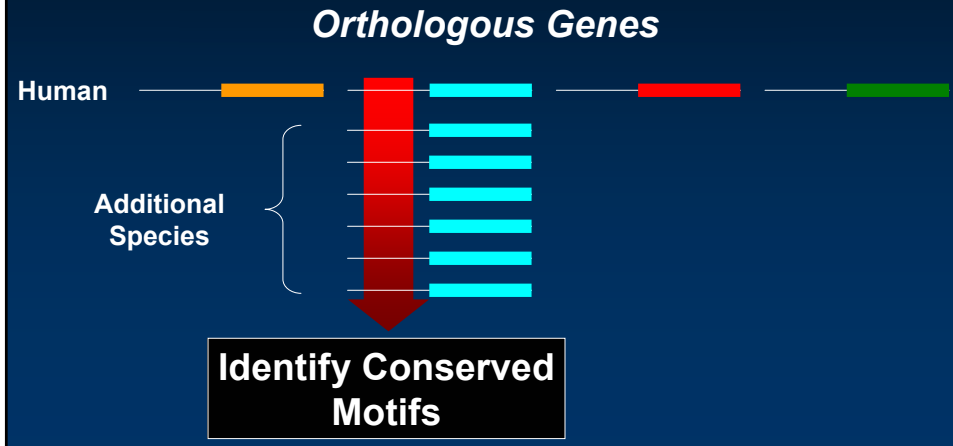
- **Identify Transcription Factor Binding Sites**
- **What sequences should be searched?**

*Coordinately Regulated Genes*



## Phylogenetic Footprinting

- **FootPrinter** – <http://bio.cs.washington.edu/software.html>
- Takes the phylogeny into account



## Summary of Phylogenetic Footprinting Tools

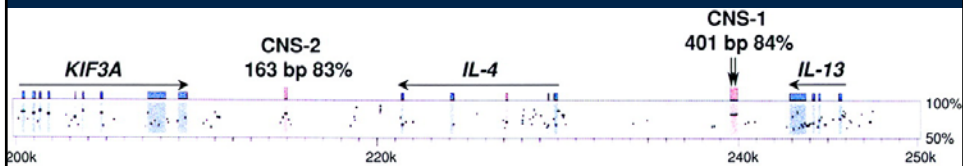
- **FootPrinter** – <http://bio.cs.washington.edu/software.html>
  - Blanchette and Tompa (2003) *Nucleic Acids Research* **31**:3840–3842
- **rVISTA** – <http://www-gsd.lbl.gov/vista/rVistaInput.html>
  - Loots et al. (2002) *Genome Research* **12**: 832–839
- **List of motif-finding algorithms:**
  - [Box 1](#) of Ureta-Vidal et al. (2003) *Nature Reviews Genetics* **4**:251-262
- **Bayesian Approaches (and home of the Gibbs sampler)**
  - <http://www.wadsworth.org/resnres/bioinfo/>
- **Example of motif-finding limited by mouse conservation:**
  - Wasserman et al. (2000) *Nature Genetics* **26**:225-228

# Prioritize Non-Coding Sequence Conservation

*Science* (2000) 288:136-140

## Identification of a Coordinate Regulator of Interleukins 4, 13, and 5 by Cross-Species Sequence Comparisons

G. G. Loots,<sup>1,2</sup> R. M. Locksley,<sup>3</sup> C. M. Blakespoor,<sup>1</sup> Z. E. Wang,<sup>3</sup> W. Miller,<sup>4</sup> E. M. Rubin,<sup>1\*</sup> K. A. Frazer<sup>1\*</sup>



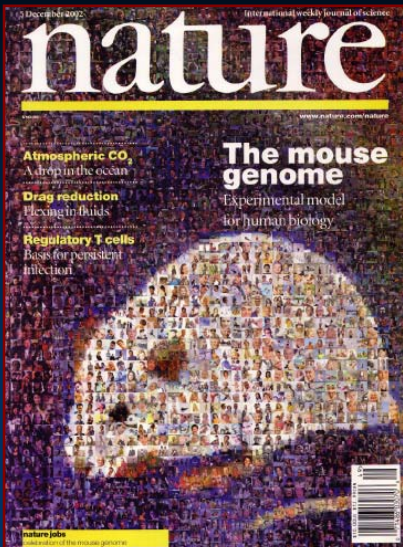
### Good review on Comparative Genomics and Regulatory Sequences:

Pennacchio and Rubin, Genomic Strategies to Identify Mammalian Regulatory Sequences, *Nature Reviews Genetics* 2:100-109

## Outline

- Fundamental concepts of comparative genomics
- Alignment and visualization tools
- Information available through genome browsers
- Gene prediction and identification
- Identifying regulatory sequences
- Insights from Human-Mouse sequence comparisons
- Multi-species sequence analysis

## Insights from Human-Mouse Sequence Comparisons

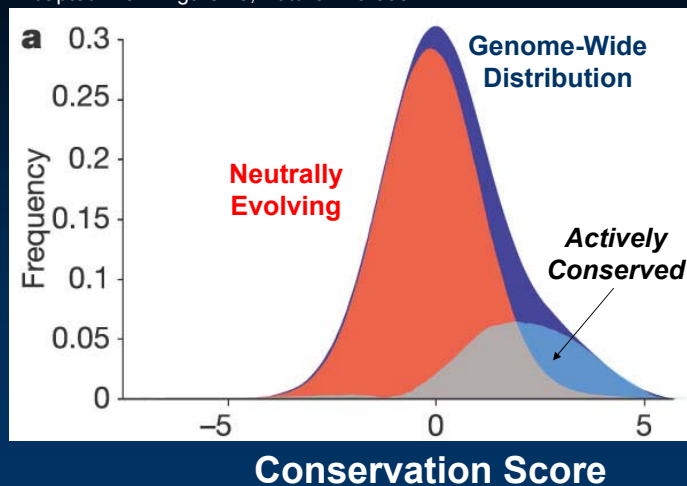


*Nature 420:520, 2002*

- Similar gene content and linear organization
  - ~340 syntenic blocks
- Difference in genome size
  - Mouse genome is 14% smaller
- Sequence Conservation
  - ~40% in Alignments
  - ~5% Under Selection
    - ~1.5% Protein Coding
    - ~3.5% Non-Coding
- Also see January 2003 issue of *Genome Research*

## Actively Conserved Sequence

Adapted From Figure 28, *Nature 420:553*



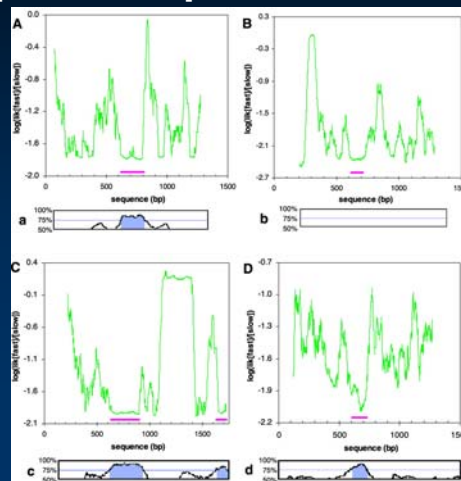
## Outline

- Fundamental concepts of comparative genomics
- Alignment and visualization tools
- Information available through genome browsers
- Gene prediction and identification
- Identifying regulatory sequences
- Insights from Human-Mouse sequence comparisons
- Multi-species sequence analysis

## Phylogenetic Shadowing

Boffelli et al. (2003) *Science* 299:1391-1394.

- Identifying sequence *differences* between multiple primate species

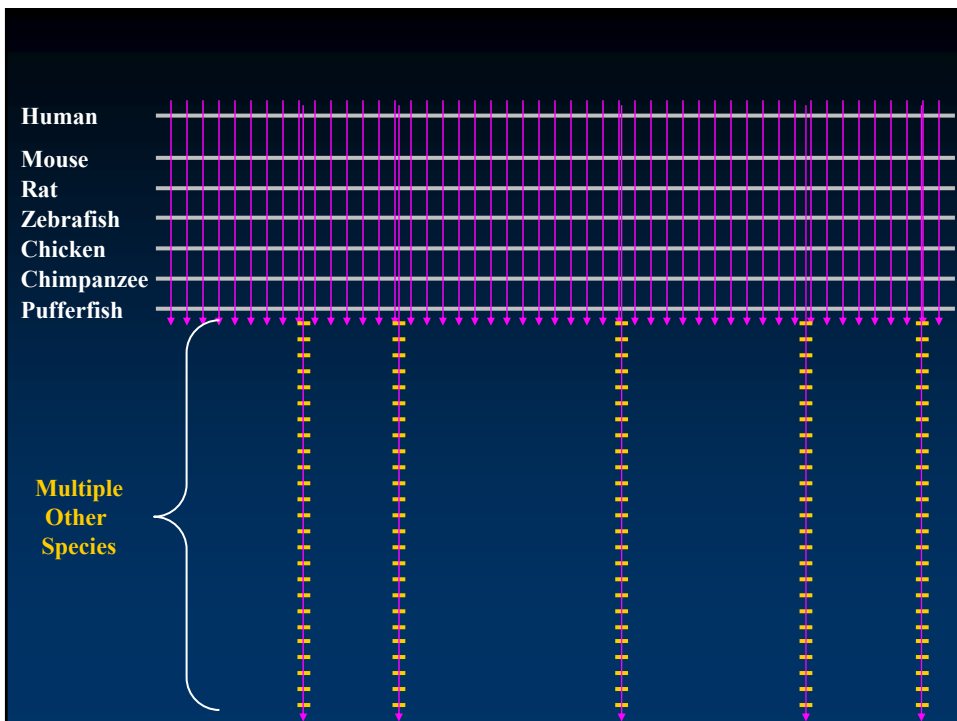


# Comparative analyses of multi-species sequences from targeted genomic regions

J. W. Thomas<sup>1,\*</sup>, J. W. Touchman<sup>1,2,\*</sup>, R. W. Blakesley<sup>1,2</sup>, G. G. Bouffard<sup>1,2</sup>, S. M. Beckstrom-Sternberg<sup>1,2</sup>, E. H. Margulies<sup>1</sup>, M. Blanchette<sup>3</sup>, A. C. Siepel<sup>3</sup>, P. J. Thomas<sup>2</sup>, J. C. McDowell<sup>2</sup>, B. Maskeri<sup>2</sup>, N. F. Hansen<sup>2</sup>, M. S. Schwartz<sup>3</sup>, R. J. Weber<sup>3</sup>, W. J. Kent<sup>3</sup>, D. Karolchik<sup>3</sup>, T. C. Bruen<sup>3</sup>, R. Bevan<sup>3</sup>, D. J. Cutler<sup>4</sup>, S. Schwartz<sup>5</sup>, L. Elnitski<sup>5</sup>, J. R. Idol<sup>1</sup>, A. B. Prasad<sup>1</sup>, S.-Q. Lee-Lin<sup>1</sup>, V. V. B. Maduro<sup>1</sup>, T. J. Summers<sup>1</sup>, M. E. Portnoy<sup>1</sup>, N. L. Dietrich<sup>2</sup>, N. Akhter<sup>2</sup>, K. Ayele<sup>2</sup>, B. Benjamin<sup>2</sup>, K. Cariaga<sup>2</sup>, C. P. Brinkley<sup>2</sup>, S. Y. Brooks<sup>2</sup>, S. Granite<sup>2</sup>, X. Guan<sup>2</sup>, J. Gupta<sup>2</sup>, P. Haghighi<sup>2</sup>, S.-L. Ho<sup>2</sup>, M. C. Huang<sup>2</sup>, E. Karlins<sup>2</sup>, P. L. Laric<sup>2</sup>, R. Legaspi<sup>2</sup>, M. J. Lim<sup>2</sup>, Q. L. Maduro<sup>2</sup>, C. A. Masiello<sup>2</sup>, S. D. Mastrian<sup>2</sup>, J. C. McCloskey<sup>2</sup>, R. Pearson<sup>2</sup>, S. Stantripop<sup>2</sup>, E. E. Tongson<sup>2</sup>, J. T. Tran<sup>2</sup>, C. Tsurgeon<sup>2</sup>, J. L. Vogt<sup>2</sup>, M. A. Walker<sup>2</sup>, K. D. Wetherby<sup>2</sup>, L. S. Wiggins<sup>2</sup>, A. C. Young<sup>2</sup>, L.-H. Zhang<sup>2</sup>, K. Osoegawa<sup>6</sup>, B. Zhu<sup>6</sup>, B. Zhao<sup>6</sup>, C. L. Shu<sup>6</sup>, P. J. De Jong<sup>6</sup>, C. E. Lawrence<sup>7</sup>, A. F. Smit<sup>8</sup>, A. Chakravarti<sup>4</sup>, D. Haussler<sup>3,9</sup>, P. Green<sup>10</sup>, W. Miller<sup>5</sup> & E. D. Green<sup>1,2</sup>

*Nature* 424:788-793

<http://genome.ucsc.edu/cgi-bin/hgGateway?org=Zoo>



## Multi-Species Weighted Conservation Score

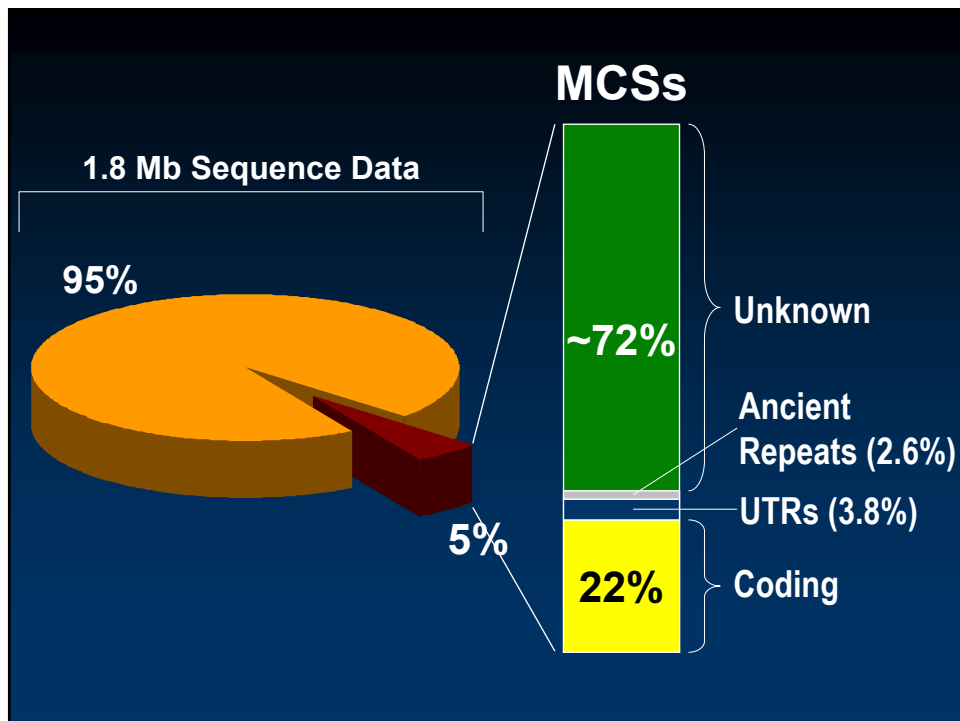
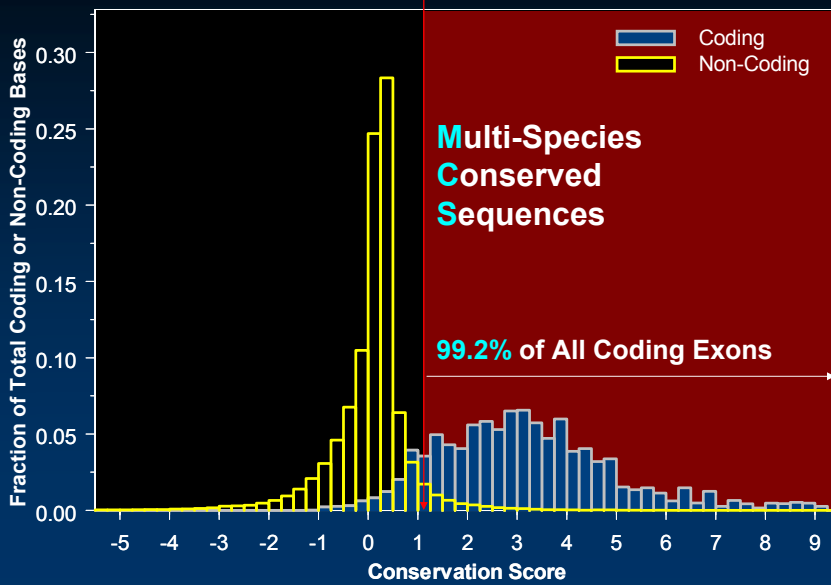
- Takes into Account the Different Divergence Rates of Each Species
  - “A Chicken Alignment Will Contribute More Than a Baboon Alignment”
- Based On the Substitution Rates at Bases under Neutral Selection
  - Calculated from 4-Fold Degenerate Positions

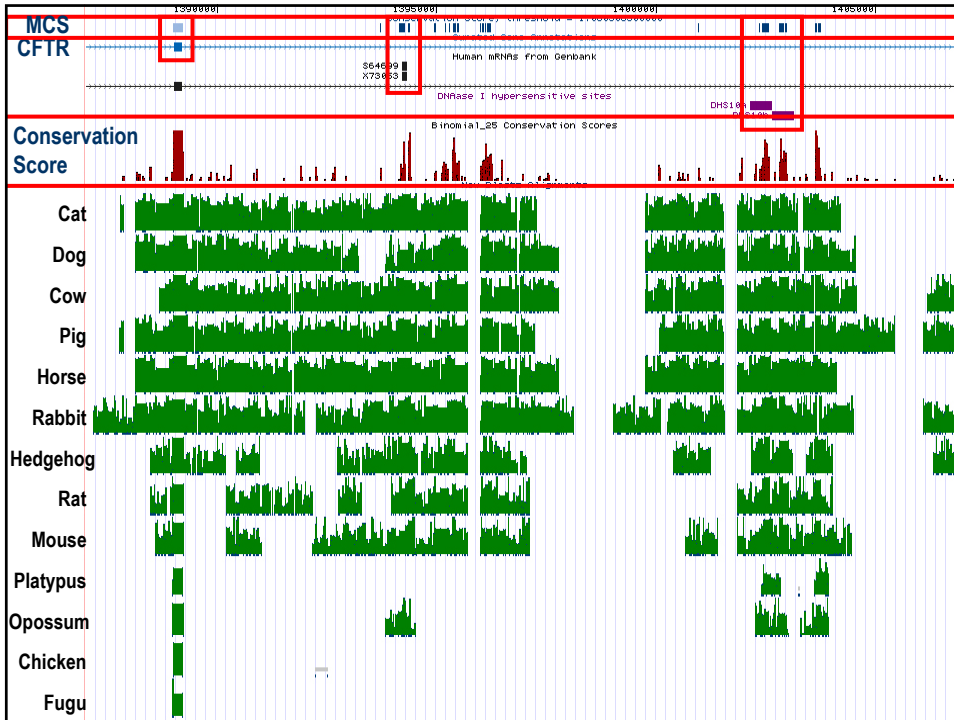
Human	GCGGGGGCCTTCGGACCGCGCGGCG	
Cat	iiiiiiiiiiiiimmmiiiiiii	i = identity
Chicken	m+miiiiiiimmmmm++iiiiiiim	m = mismatch
Chimpanzee	iiiiiiiiiiiiiiiiiiiiiiiiiiii	+ = insertion
Baboon	iiiiiiiiiiiiimmmiiiiiiiiiiii	- = unalignable
Dog	iiiiiiiiiiii++++++iiiiiiiiiii	
Cow	immmiiiiimmmmmiiiiiiiiiiiiiii	
Pig	iiimiiiiiiiiimmmmmiiiiiiiiiiii	
Rat	im++++++mmmmmmiiiiiiiiiiiiiiii	
Mouse	immmiiiiimmmmmmmmmmmmmmmmm	
Fugu	-----	
Tetraodon	-----	
Zebrafish	-----	

*Weighted Conservation Score*

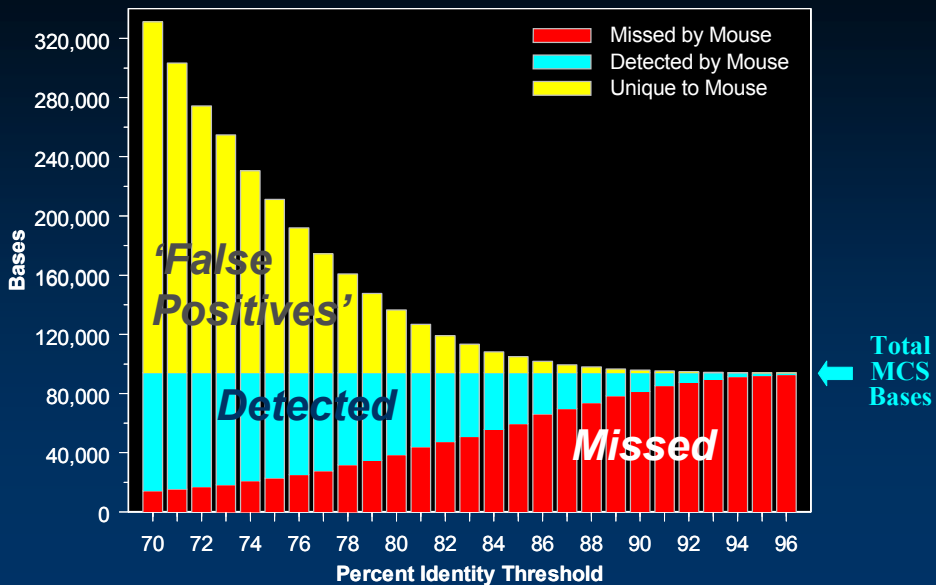


# Multi-Species Conservation Score Distribution

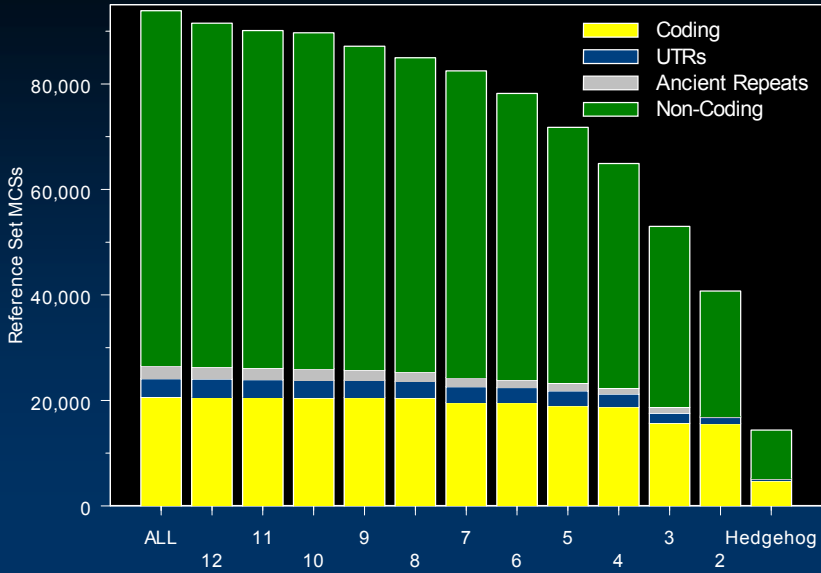




## MCS Overlap with Mouse Alignments



## Subsets of Species Perform Better



## More Species is Better

