



Current Topics in Genome Analysis  
March 1, 2005



## 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
  - Pair-wise and multi-species methods
  - Combining with transcription factor binding site data
- Motif Identification
- Comparative genomics resources available at UC Santa Cruz -- <http://genome.ucsc.edu>
  - Genome-wide sequence availability
  - Gene prediction and identification Finding orthologous sequences in other species
  - Identifying conserved sequences
- Insights from vertebrate genome sequence comparisons
- Multi-species sequence analysis

# Finishing the euchromatic sequence of the human genome

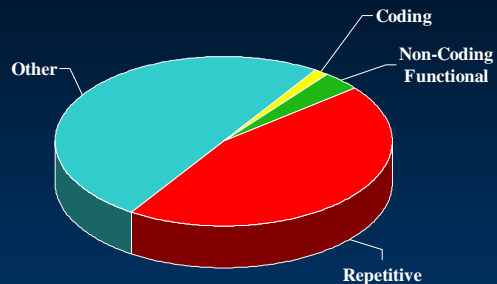
International Human Genome Sequencing Consortium\*

\* A list of authors and their affiliations appears in the Supplementary Information



## Why Compare Genomic Sequences from Different Species?

- Explore evolutionary relationships



- Enhanced gene prediction algorithms
- Identify **functionally constrained** sequence

## Charles Darwin

- Served as *naturalist* on a British science expedition around the world (1831 -- 1836)
- *On the Origin of Species by Means of Natural Selection, or the Preservation of Favoured Races in the Struggle for Life*
- ***The Origin of Species* (1859)**
  - All species evolved from a single life form
  - “Variation” within a species occurs randomly
  - Natural selection
  - Evolutionary change is gradual

## Other Intellectual Foundations

- Darwin (1859)
  - Theories of Evolution
- Mendel (1866) (*rediscovered in 1900*)
  - Genes are units of heredity
- Avery, McCarty & MacLeod (1944)
  - DNA as the “transforming principle”
- Watson & Crick (1953)
  - Structure of DNA
- Sanger (1977)
  - Methods of sequencing DNA

## Rationale

- DNA represents a “blueprint” for structure and physiology of all living things
- All species use DNA
- Mutations in *functional* DNA are less likely to be tolerated

## 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



## What's in a Name?

- Highly conserved sequences
- Sequences under purifying selection
- Functionally constrained sequences
- **ECOR** – **E**volutionary **CO**nserved **R**egion  
– Variant: **ECR**
- **CNS** – **C**onserved **N**on-coding **S**equence
- **CNGs** – **C**onserved **N**on-**G**enic **s**equence
- **MCS** – **M**ulti-species **C**onserved **S**equence

## Outline

- Fundamental concepts of comparative genomics
- **Alignment and visualization tools**
  - Pair-wise and multi-species methods
  - Combining with transcription factor binding site data
- **Motif Identification**
- **Comparative genomics resources available at UC Santa Cruz** -- <http://genome.ucsc.edu>
  - Genome-wide sequence availability
  - Gene prediction and identification
  - Finding orthologous sequences in other species
  - Identifying conserved sequences
- **Insights from vertebrate genome sequence comparisons**
- **Multi-species sequence analysis**

# Sequence Alignments

## 100% Identical

Species 1 CATGGGCAAATTGGCCATTGGCCATGGGGGCCACCGTA  
Species 2 CATGGGCAAATTGGCCATTGGCCATGGGGGCCACCGTA

## 80% Identical

Species 1 CATGGGCAAATTGGCCATTGGCCATGGGGGCCACCGTA  
Species 2 CACGGGCTAATCCGCCAATTGGCTATGGGG-CCCAGCGTA

## 30% Identical

Species 1 CATGGGCAAATTGGCCATTGGCCATGGGGGCCACCGTA  
Species 2 CACGAACTAATCCGCCAATAGCCTATAGCG-CACAGCGAA

# Tools for Aligning Genomic Sequences

Resource: *Genome Research* (2000) 10:577-586

### PipMaker—A Web Server for Aligning Two Genomic DNA Sequences

Scott Schwartz,<sup>1</sup> Zheng Zhang,<sup>1</sup> Kelly A. Frazer,<sup>2</sup> Arian Smit,<sup>3</sup> Cathy Riemer,<sup>1</sup> John Bouck,<sup>4</sup> Richard Gibbs,<sup>4</sup> Ross Hardison,<sup>5</sup> and Webb Miller<sup>1,6</sup>

Departments of <sup>1</sup>Computer Science and Engineering and <sup>2</sup>Biochemistry and Molecular Biology and Center for Gene Regulation, The Pennsylvania State University, University Park, Pennsylvania USA 16802; <sup>3</sup>Genome Sciences Department, Lawrence Berkeley National Laboratory, Berkeley, California USA 94720; <sup>4</sup>Asyst Pharmaceuticals, La Jolla, California USA 92037; <sup>5</sup>Department of Molecular and Human Genetics, Baylor College of Medicine, Houston, Texas USA 77030

**BIOINFORMATICS APPLICATIONS NOTE** Vol. 16 no. 11 2000  
Pages 1046–1047

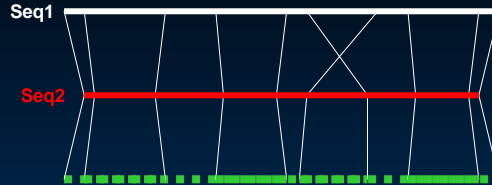
#### VISTA: visualizing global DNA sequence alignments of arbitrary length

Chris Mayor<sup>1</sup>, Michael Brudno<sup>1</sup>, Jody R. Schwartz<sup>2</sup>, Alexander Poliakov<sup>2</sup>, Edward M. Rubin<sup>2</sup>, Kelly A. Frazer<sup>2</sup>, Lior S. Pachter<sup>3,\*</sup> and Inna Dubchak<sup>1,\*</sup>

<sup>1</sup>National Energy Research Scientific Computing Center, <sup>2</sup>Genome Sciences Department, Lawrence Berkeley National Laboratory, Berkeley, CA 94720, USA and <sup>3</sup>Department of Mathematics University of California at Berkeley, Berkeley, CA 94720, USA

## PipMaker vs. VISTA

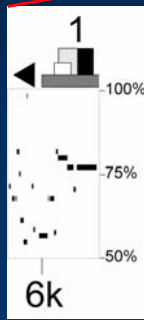
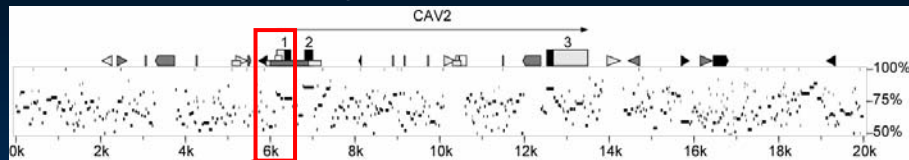
- Visualization
- Alignment Strategy
  - VISTA: *avid*
  - PipMaker: *blastz*
- East Coast – West Coast



## PipMaker

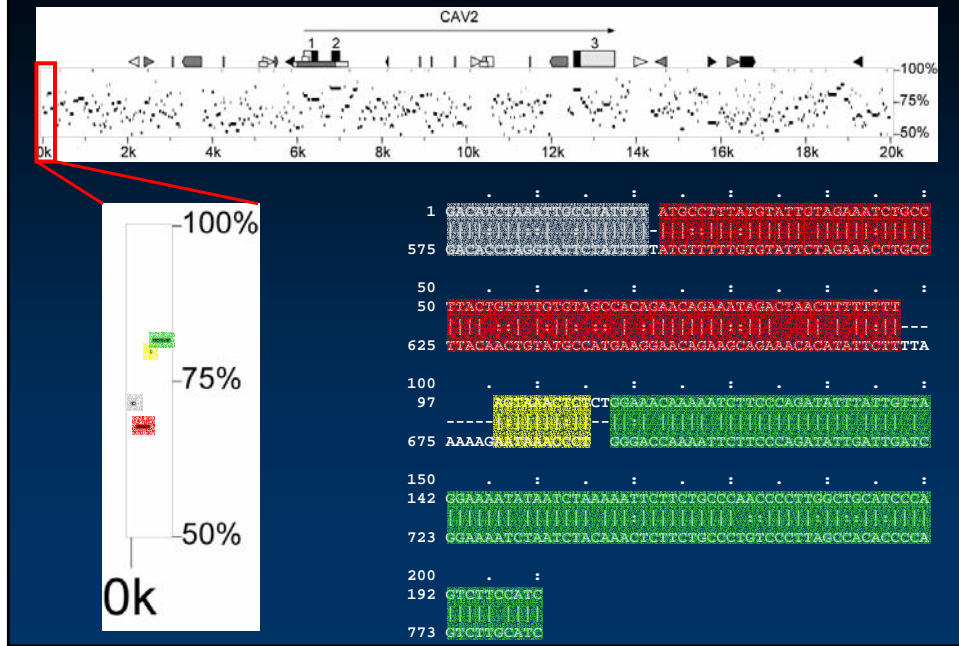
<http://bio.cse.psu.edu/pipmaker/>

- Percent Identity Plot

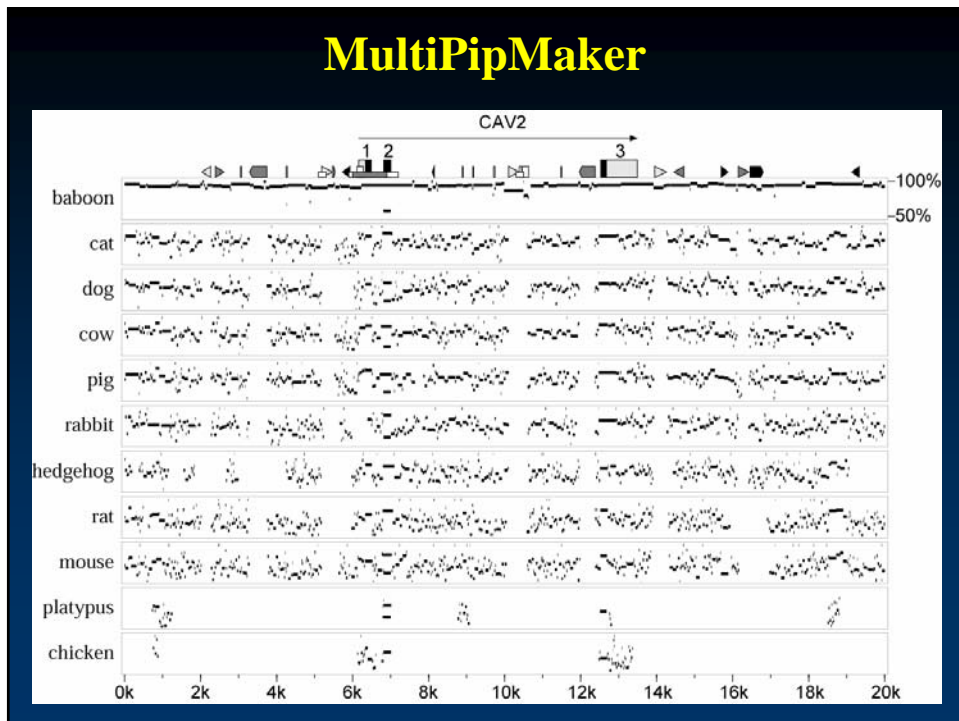


- X-axis is the reference sequence
- Horizontal lines represent gap-free alignments

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



## MultiPipMaker





<http://www-gsd.lbl.gov/vista/>

VISUALIZATION TOOLS FOR ALIGNMENTS

# VISTA

- Global Alignment (avid)
  - Bray et al. (2003) *Genome Res* 13:97-102
- Sliding Window Approach to Visualization
  - Plot Percent Identity within a Fixed Window Size, at Regular Intervals

```

GACATCTAAATTGCCTATTTT ATGCCTTTATGTATTGTAGAAATCTGCCTTACTGTTTTGTGTAGCCACAGAACAGAAATAGACTAACTTTTTTTT
||||:||||:|:| :|||||||:-|||:|:||||| |||||:||||||| :|:| |:|: :| :|||||||:|:|||| | | | | :||
GACACCTAGGTATTCTATTTTATGTTTTTGTGTTTCTAGAAACCTGCCTTACAACCTGTATGCCATGAAGGAACAGAAGCCAGAAACACATATTCTT
  
```

72% 68% 80% 88% 76% 52% 56% 64%

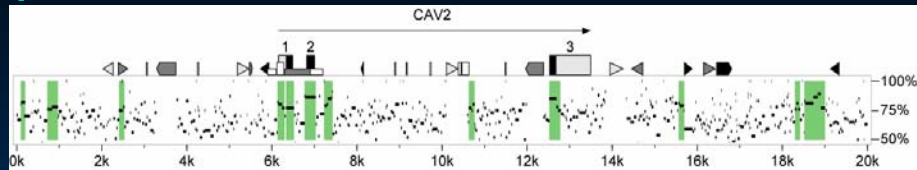
## VISTA

CAV2

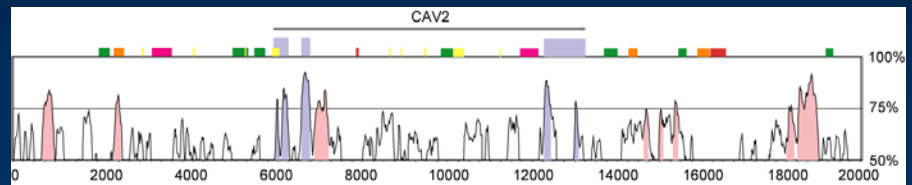
- Percent Identity is plotted from:
  - 100 base windows
  - Moved every 15 bases
- Colored regions meet certain alignment criteria
  - >100 bp >75% Identity

## What's Your Preference?

### PipMaker



### VISTA



## East & West Coast Unite

<http://zpicture.dcode.org/>

### Resource

## zPicture: Dynamic Alignment and Visualization Tool for Analyzing Conservation Profiles

Ivan Ovcharenko,<sup>1,2</sup> Gabriela G. Loots,<sup>2</sup> Ross C. Hardison,<sup>3</sup> Webb Miller,<sup>4,5</sup> and Lisa Stubbs<sup>2,6</sup>

<sup>1</sup>Energy, Environment, Biology and Institutional Computing, Lawrence Livermore National Laboratory, Livermore, California 94550, USA; <sup>2</sup>Genome Biology Division, Lawrence Livermore National Laboratory, Livermore, California 94550, USA;

<sup>3</sup>Department of Biochemistry and Molecular Biology, The Pennsylvania State University, University Park, Pennsylvania 16802, USA; <sup>4</sup>Department of Computer Science and Engineering, The Pennsylvania State University, University Park, Pennsylvania 16802, USA; <sup>5</sup>Department of Biology, The Pennsylvania State University, University Park, Pennsylvania 16802, USA

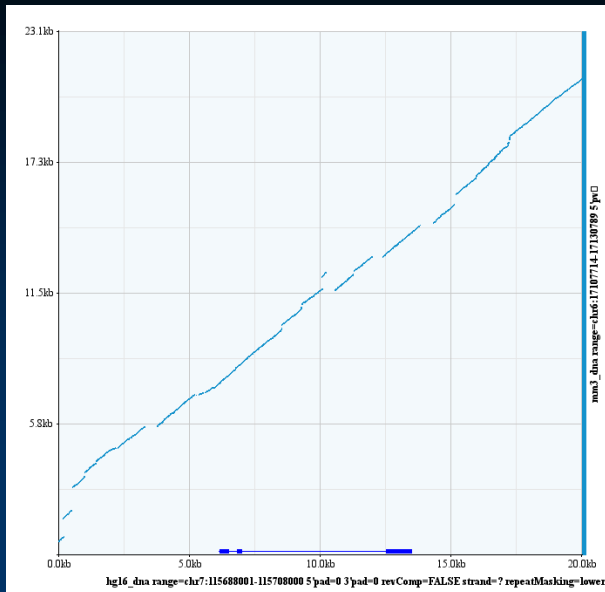
<sup>6</sup>Department of Biology, The Pennsylvania State University, University Park, Pennsylvania 16802, USA

*Genome Research*, 2004, 14(3):472–7

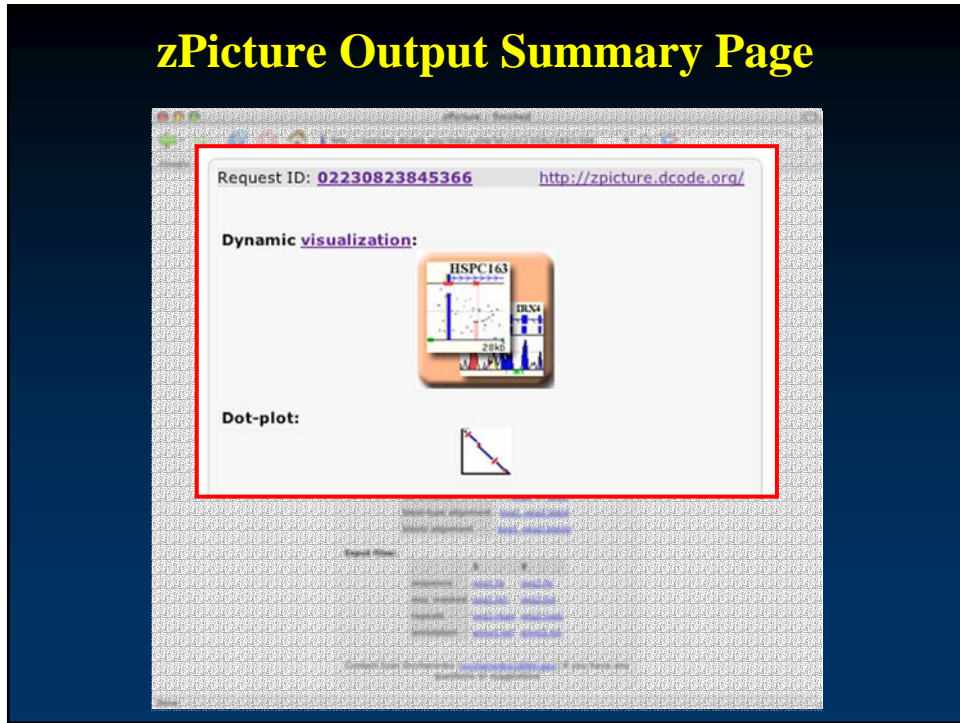
# zPicture Output Summary Page



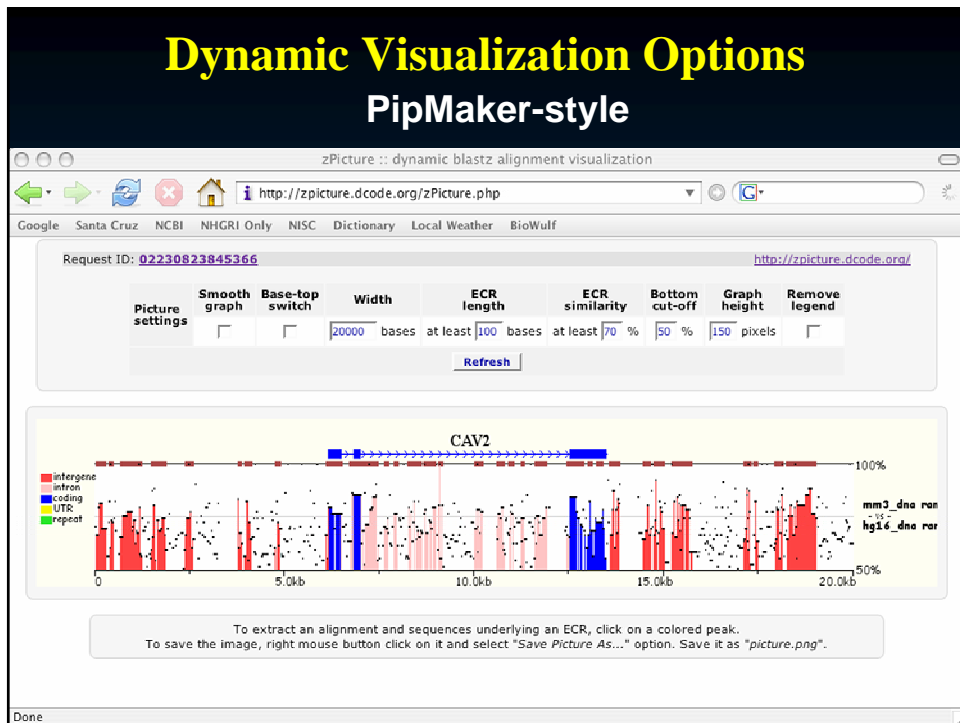
# zPicture: Dot Plot View



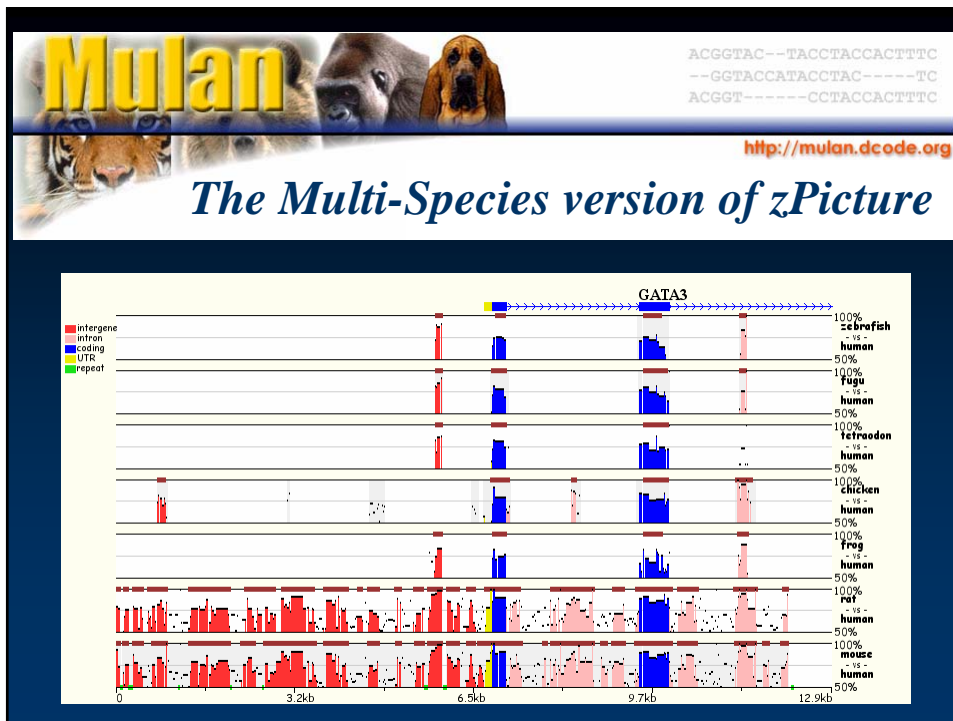
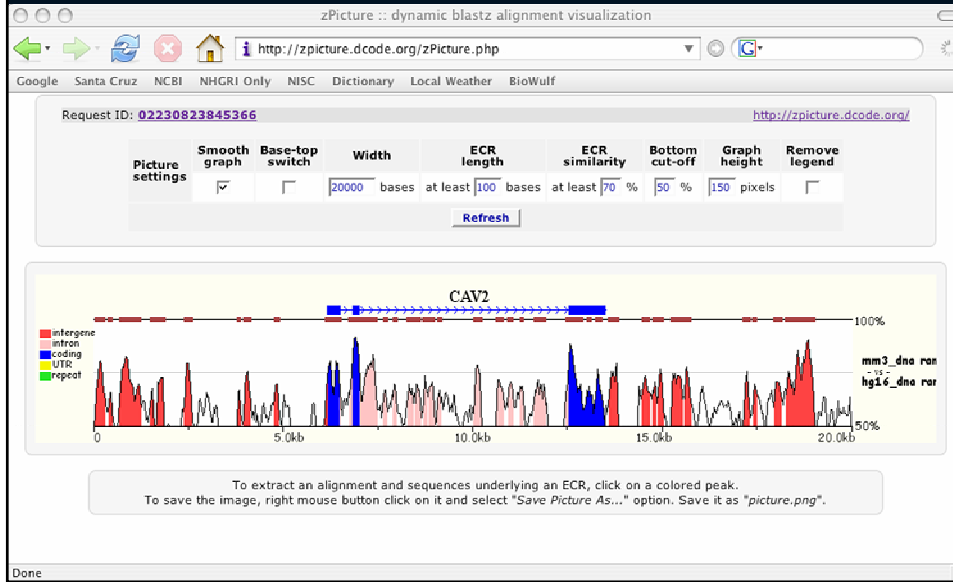
## zPicture Output Summary Page



## Dynamic Visualization Options PipMaker-style



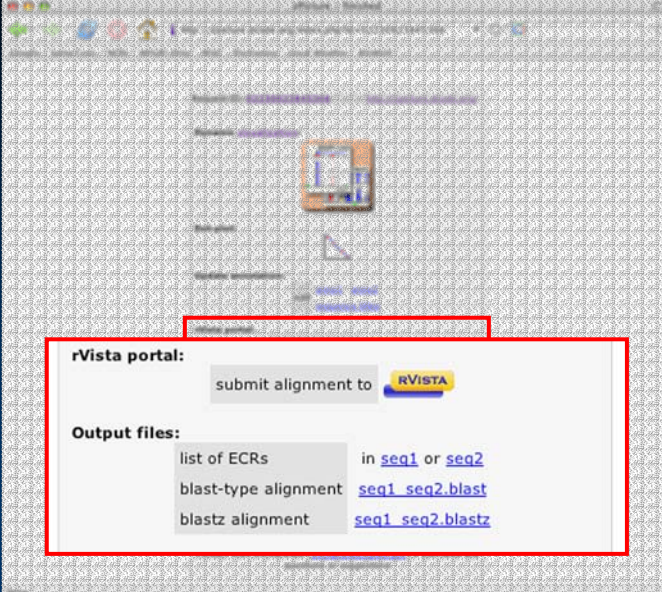
# Dynamic Visualization VISTA-style



## Outline

- Fundamental concepts of comparative genomics
- Alignment and visualization tools
  - Pair-wise and multi-species methods
  - Combining with transcription factor binding site data
- Motif Identification
- Comparative genomics resources available at UC Santa Cruz -- <http://genome.ucsc.edu>
  - Genome-wide sequence availability
  - Gene prediction and identification
  - Finding orthologous sequences in other species
  - Identifying conserved sequences
- Insights from vertebrate genome sequence comparisons
- Multi-species sequence analysis

## zPicture Output Summary Page



The screenshot shows a web browser window displaying the zPicture Output Summary Page. A red box highlights the 'rVista portal:' section, which includes a 'submit alignment to' button with the 'RVISTA' logo. Below this, the 'Output files:' section lists three types of alignments and their corresponding file names:

Alignment Type	File Name
list of ECRs	in <a href="#">seq1</a> or <a href="#">seq2</a>
blast-type alignment	<a href="#">seq1_seq2.blast</a>
blastz alignment	<a href="#">seq1_seq2.blastz</a>

## Are there any transcription factor binding sites in my alignment?



### Pre-Computed Sequence Alignments

```

human  GGAGAGTGTGACAGATGTTATCATTTGCTCCATTTGTGACGAGCGTAGGACCATGCTTC
chimp  GGAGAGTGTGACAGATGTTATCATTTGCTCCATTTGTGACGAGCGTAGGACCATGCTTC
dog    GGAGAGTGTGACAGATGTTATCATTTGCTCCATTTGTGACGAGCGTAGGACCATGCTTC
mouse  GGAGAGTGTGACAGAGCTTTCATTTGCCCATTTGTGACAGTGTGGCCGCGAGCTTC
rat    GGAGAGTGTGACAGATGTTATCATTTGCCCATTTGTGACAGTGTGGCCGCGATCTTC
chicken GGAGAGCGCTGACAGATGTTTCATTTGCCCATTTTCATGAGCGTAGGACCATGCTTC
Fugu   GGAGAGCGCTGACAGATGTTTCATTTGCCCATTTTCATGAGCGTAGGACCATGCTTC
zebrafish GGAGTAGGCTTTTAGACATATCATCTCCCTTCTTCGAGAGCGTAGGACCATGCTTC
  
```



**TRANSFAC**



<http://www.gene-regulation.com/>

## TRANSFAC



<http://www.gene-regulation.com/>

- A database of:
  - Eukaryotic transcription factors
  - Their genomic binding sites
  - And DNA binding profiles
- Data are collected from published studies
  - Non curated
  - Redundant data

# JASPAR: An Alternative to TRANSFAC

*Nucleic Acids Research*, 2004, Vol. 32, Database issue D91-D94  
DOI: 10.1093/nar/gkh012

## JASPAR: an open-access database for eukaryotic transcription factor binding profiles

Albin Sandelin, Wynand Alkema, Pär Engström, Wyeth W. Wasserman<sup>1</sup> and Boris Lenhard\*

Center for Genomics and Bioinformatics, Karolinska Institutet, Berzelius väg 35, S-17177 Stockholm, Sweden and  
<sup>1</sup>Centre for Molecular Medicine and Therapeutics, Department of Medical Genetics, University of British Columbia, Vancouver, BC, Canada

- Differences from TRANSFAC:
  - Manually curated for “high quality” experiments
  - Non redundant collection

<http://jaspar.cgb.ki.se/>

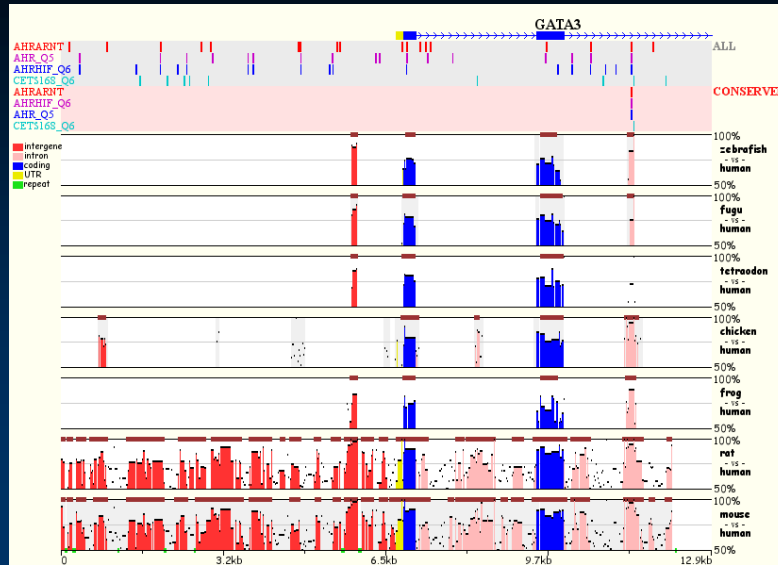


## TRANSFAC Data are inherently “noisy”

- Binding sites are very short  
6-10 bases in length
- Low complexity  
Only 4 “letters” in the DNA alphabet
- Frequently observe binding site by chance
- *Conservation can help reduce the noise*
- **rVISTA 2.0** for pair-wise alignments
- **multiTF** for multi-species alignments



## Example of multiTF Output



## Summary of Alignment Tools

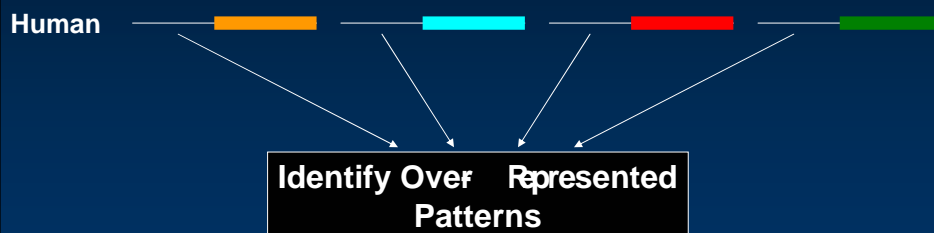
- PipMaker (blastz)
- VISTA (avid)
- zPicture and MULAN
- Lagan and mLagan (glocal alignments)  
– <http://lagan.stanford.edu/>
- rVISTA 2.0
- **Box 1** from:  
Ureta-Vidal, Ettwiller, and Birney (2003) Comparative Genomics: Genome-Wide Analysis in Metazoan Eukaryotes *Nature Reviews Genetics* 4: 251-262
- **Table 1** from:  
Miller, Makova, Nekrutenko, and Hardison (2004) Comparative Genomics *Annual Reviews in Human Genetics* 5:15-56

## Outline

- Fundamental concepts of comparative genomics
- Alignment and visualization tools
  - Pair-wise and multi-species methods
  - Combining with transcription factor binding site data
- Motif Identification
- Comparative genomics resources available at UC Santa Cruz -- <http://genome.ucsc.edu>
  - Genome-wide sequence availability
  - Gene prediction and identification Finding orthologous sequences in other species
  - Identifying conserved sequences
- Insights from vertebrate genome sequence comparisons
- Multi-species sequence analysis

## Motif Finding

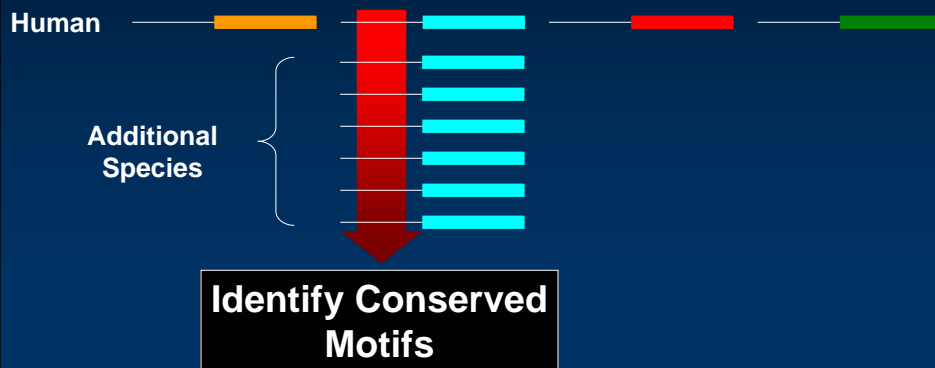
- Identify Transcription Factor Binding Sites
- What sequences should be searched?  
*Coordinately Regulated Genes*



## Phylogenetic Footprinting

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

### Orthologous Genes



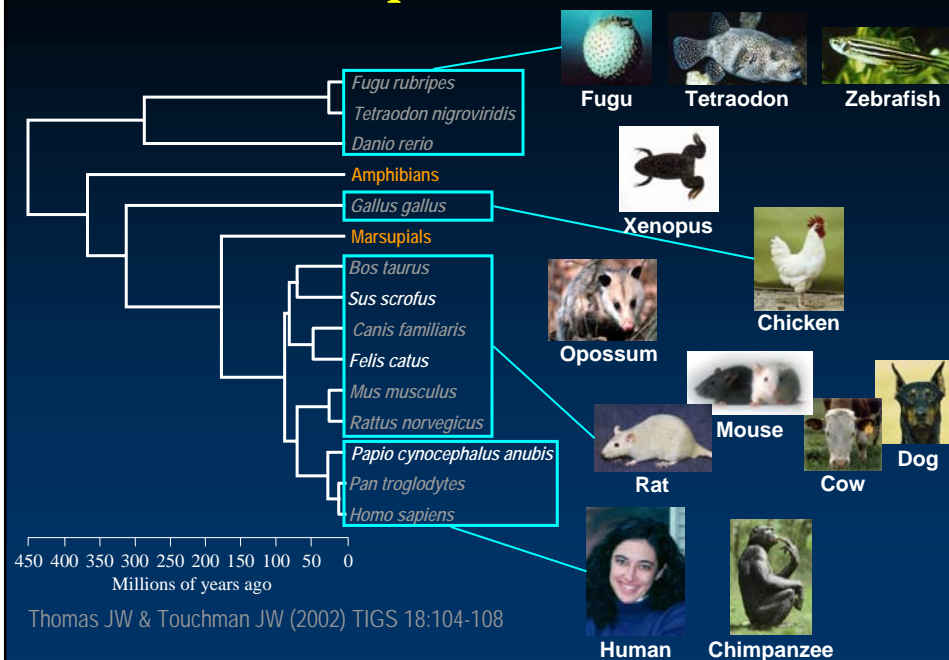
## Summary of Phylogenetic Footprinting Tools

- **FootPrinter** – <http://bio.cs.washington.edu/software.html>
  - Blanchette and Tompa (2003) *Nucleic Acids Research* **31**:3840–3842
- **phyloCon** – <http://oldural.wustl.edu/~twang/PhyloCon/>
  - Wang and Stormo (2003) *Bioinformatics* **19**:2369-80
- **phyME**
  - Sinha, Blanchette, and Tompa (2004) *BMC Bioinformatics* **28**:170
- **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

## Outline

- Fundamental concepts of comparative genomics
- Alignment and visualization tools
  - Pair-wise and multi-species methods
  - Combining with transcription factor binding site data
- Motif Identification
- Comparative genomics resources available at UC Santa Cruz -- <http://genome.ucsc.edu>
  - Genome-wide sequence availability
  - Gene prediction and identification Finding orthologous sequences in other species
  - Identifying conserved sequences
- Insights from vertebrate genome sequence comparisons
- Multi-species sequence analysis

## Genome-Wide Sequences

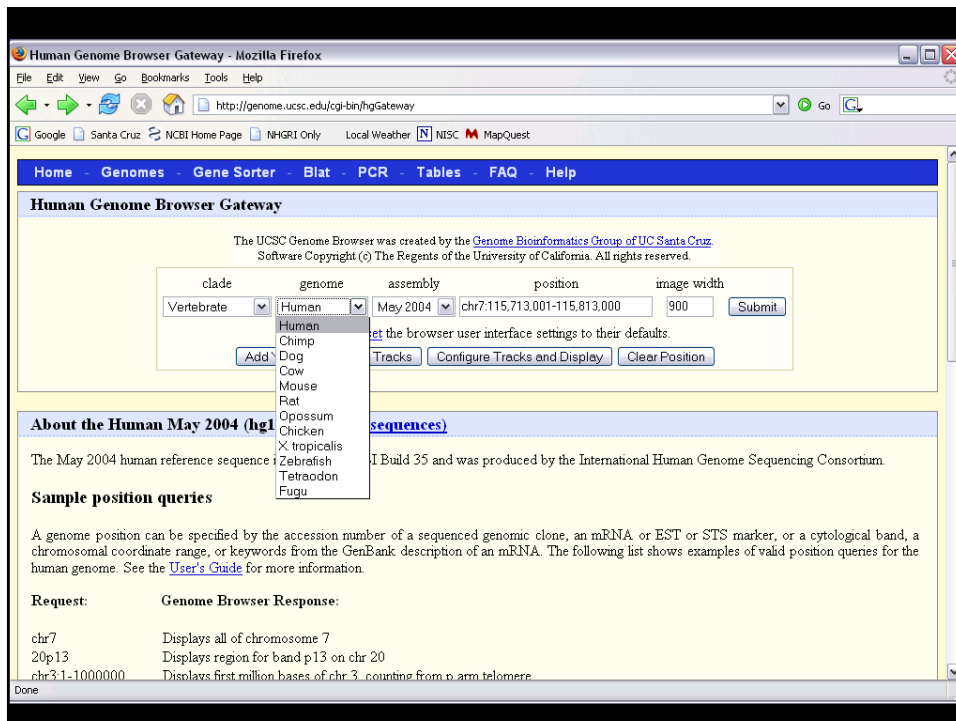


# Genome Browsers

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

 project **Ensembl**  
<http://www.ensembl.org>

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



Human Genome Browser Gateway - Mozilla Firefox

File Edit View Go Bookmarks Tools Help

<http://genome.ucsc.edu/cgi-bin/hgGateway>

Google Santa Cruz NCBI Home Page NHGRI Only Local Weather NISC MapQuest

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

**Human Genome Browser Gateway**

The UCSC Genome Browser was created by the [Genome Bioinformatics Group of UC Santa Cruz](#).  
Software Copyright (c) The Regents of the University of California. All rights reserved.

clade	genome	assembly	position	image width
Vertebrate	Human	May 2004	chr7:115,713,001-115,813,000	900

Submit

set the browser user interface settings to their defaults.

Tracks Configure Tracks and Display Clear Position

**About the Human May 2004 (hg1)**

The May 2004 human reference sequence is Build 35 and was produced by the International Human Genome Sequencing Consortium.

**Sample position queries**

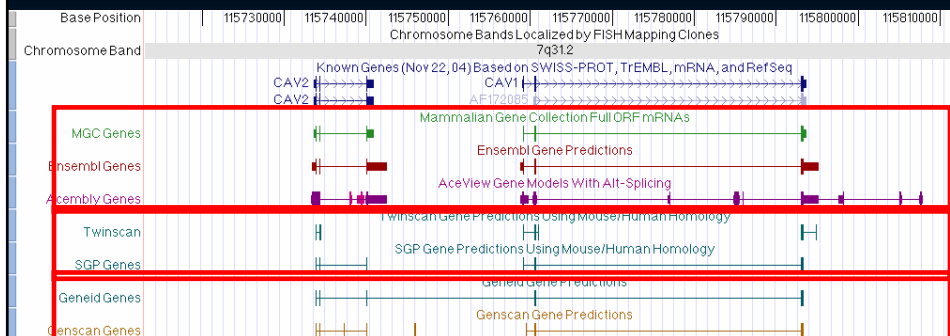
Request:	Genome Browser Response:
chr7	Displays all of chromosome 7
20p13	Displays region for band p13 on chr 20
chr3:1-1000000	Displays first million bases of chr 3, counting from p arm telomere

Done

## Outline

- Fundamental concepts of comparative genomics
- Alignment and visualization tools
  - Pair-wise and multi-species methods
  - Combining with transcription factor binding site data
- Motif Identification
- Comparative genomics resources available at UC Santa Cruz -- <http://genome.ucsc.edu>
  - Genome-wide sequence availability
  - Gene prediction and identification
  - Finding orthologous sequences in other species
  - Identifying conserved sequences
- Insights from vertebrate genome sequence comparisons
- Multi-species sequence analysis

## Approaches to Gene Prediction



- Evidence-Based
  - MGC
  - Acemby
  - Ensembl
- *Ab Initio*
  - Genscan
  - Geneid
- Dual-Genome
  - Twinscan
  - SGP

## Additional Gene Prediction Resources

- **Fugu BLAT Track at UCSC**
- **SLAM** – <http://baboon.math.berkeley.edu/~syntenic/slam.html>
  - Cawley et al. (2003) *Nucleic Acids Research* **31**:3507-3509
- **Exoniphy**
  - Siepel and Haussler. Computational identification of evolutionarily conserved exons. *Proc. 8th Annual Int'l Conf. on Research in Computational Biology*, pp. 177-186, 2004.
  - <http://www.soe.ucsc.edu/~acs/recomb2004.pdf>
  - Also see genome “test” browser for data
- **Box 1** from:
  - Ureta-Vidal et al. (2003) *Nature Reviews Genetics* **4**:251-262

## Outline

- Fundamental concepts of comparative genomics
- Alignment and visualization tools
  - Pair-wise and multi-species methods
  - Combining with transcription factor binding site data
- Motif Identification
- Comparative genomics resources available at UC Santa Cruz -- <http://genome.ucsc.edu>
  - Genome-wide sequence availability
  - Gene prediction and identification
  - Finding orthologous sequences in other species
  - Identifying conserved sequences
- Insights from vertebrate genome sequence comparisons
- Multi-species sequence analysis



## Chaining Alignments

- Chaining bridges the gulf between large syntenic blocks and base bybase alignments.

### The Challenge:

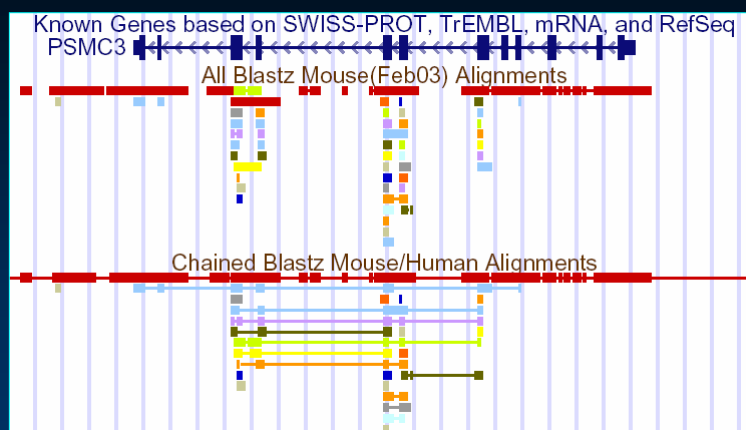
- Local alignments tend to break at transposon insertions, inversions, duplications, etc.
- Global alignments tend to force non homologous bases to align.

### The Solution:

- Chaining is a rigorous way of joining together local alignments into larger structures.

*Slide (though modified) Courtesy of Jim Kent*

## Chains join together related local alignments

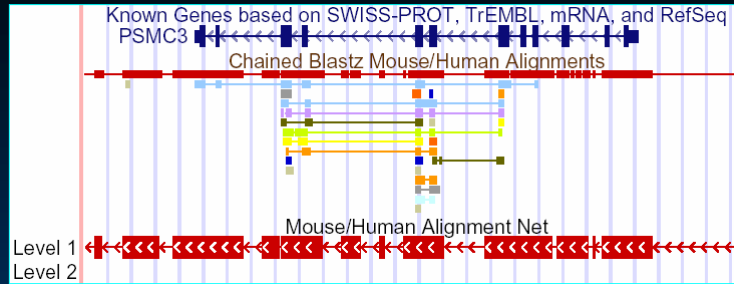


**Protease Regulatory Subunit 3**

*Slide Courtesy of Jim Kent*



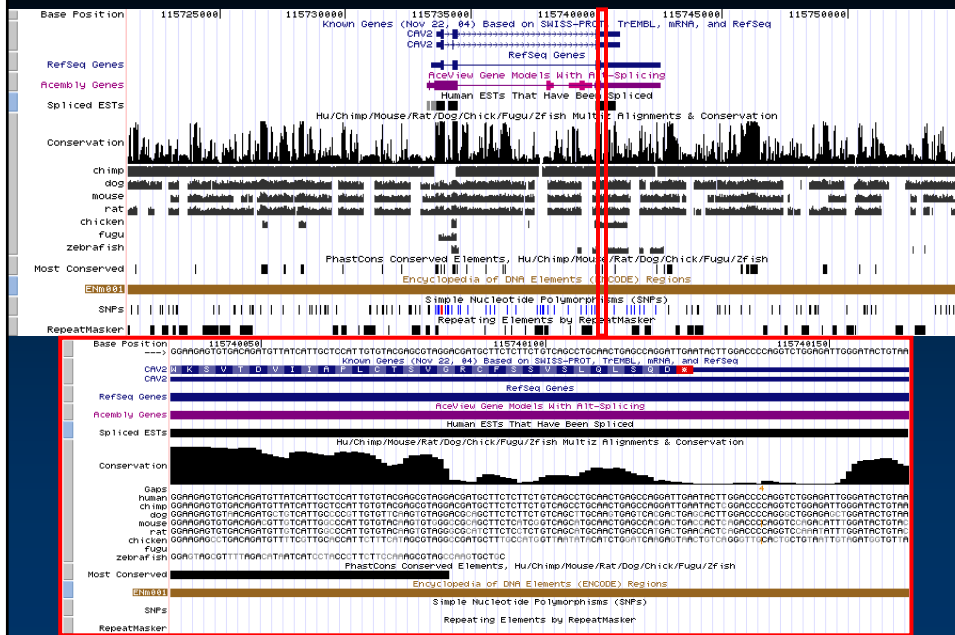
# Net Alignments: Focus on Orthology



- Frequently, there are numerous mouse alignments for any given human region, particularly for coding regions.
- Net finds best mouse match for each human region.

Slide (though modified) Courtesy of Jim Kent

# Genome-wide Multiple Sequence Alignments



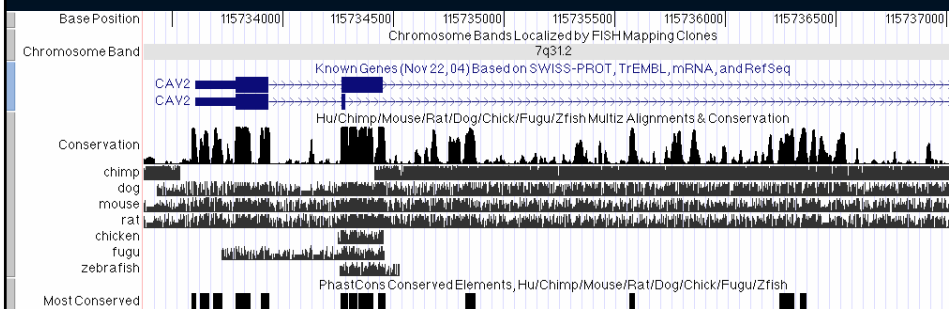
## Conservation Score at UCSC

- Displays evolutionary conservation based on a phylogenetic hidden Markov model

Probability that the given alignment column was generated by the “conserved” state

- “Most Conserved” track represents highly conserved regions
  - Tuned to cover ~4% of the genome

## “Most Conserved” Track at UCSC



*What if you want a different stringency than was used for the Most Conserved Track?*

Human chr7:115,733,372-115,737,074 - UCSC Genome Browser v100 - Mozilla Firefox

File Edit View Go Bookmarks Tools Help

http://genome.ucsc.edu/cgi-bin/hgTables

Google Santa Cruz NCBI Home Page NHGRI Only Local Weather NISC MapQuest

Home Genomes Blat PCR DNA **Tables** Gene Sorter Convert Ensembl NCBI PDF/PS Help

### UCSC Genome Browser on Human May 2004 Assembly

move <<< << < > >> >>> zoom in 1.5x 3x 10x base zoom out 1.5x 3x 10x

position chr7:115,733,372-115,737,074 jump clear size 3,703 bp configure

chr7 (q31.2) 115,734,000 115,734,500 115,735,000 115,735,500 115,736,000 115,736,500 115,737,000

Base Position  
Chromosome Band  
Known Genes (Nov 22, 04) Based on S155-FROT, TRENEL, mRNA, and RefSeq  
Conservation  
chr:imp  
dog  
mouse  
rat  
chicken  
fish  
zebrafish  
Most Conserved

move start < 2.0 > move end < 2.0 >

Click on a feature for details. Click on base position to zoom in around cursor. Click on left mini-buttons for track-specific options.

default tracks hide all configure refresh

Use drop down controls below and press refresh to alter tracks displayed. Tracks with lots of items will automatically be displayed in more compact modes.

#### Mapping and Sequencing Tracks

Base Position	Chromosome Band	STS Markers	RGD QTL	FISH Clones
dense	dense	hide	hide	hide
Recomb Rate	Map Contigs	Assembly	Gap	Coverage

http://genome.ucsc.edu/cgi-bin/hgTables?db=hg17&position=chr7:115733372-115737074&hgta\_regionType=range&hgid=39729767

## Using the Table Browser to get Highly Conserved Sequences

Output phastCons as Custom Track - Mozilla Firefox

File Edit View Go Bookmarks Tools Help

http://genome.ucsc.edu/cgi-bin/hgTables?hgid=39729767

Google Santa Cruz NCBI Home Page NHGRI Only Local Weather NISC MapQuest

Home Genomes Gene Sorter Blat PCR Tables **FAQ** Help

### Output phastCons as Custom Track

Custom track header:

name= tb\_phastCons

description= table browser query on phastCons

visibility= pack

url=

Select type of data output:

BED format (no data value information, only position)

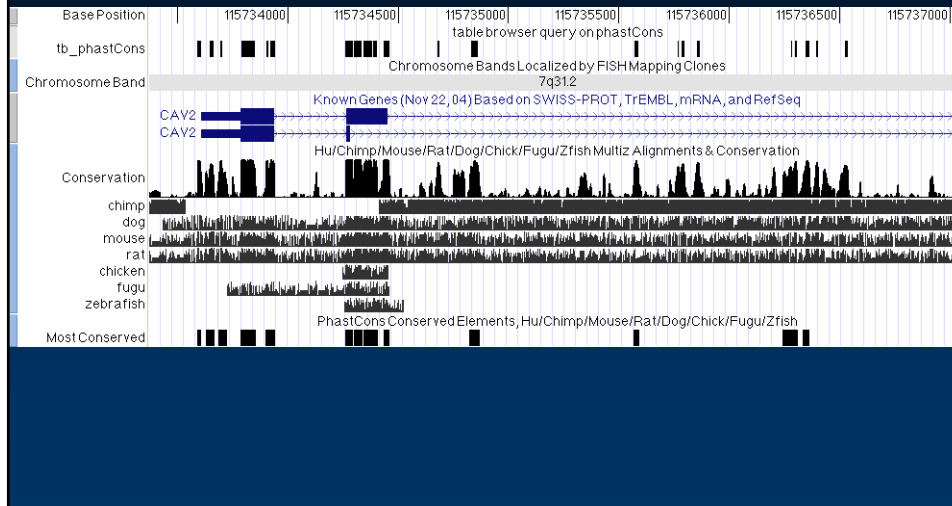
DATA VALUE format (position and real valued data)

Get Custom Track in Table Browser Get Custom Track in File

Get Custom Track in Genome Browser Cancel

Done

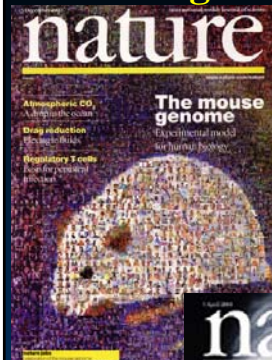
## Using the Table Browser to get Highly Conserved Sequences



## Outline

- **Fundamental concepts of comparative genomics**
- **Alignment and visualization tools**
  - Pair-wise and multi-species methods
  - Combining with transcription factor binding site data
- **Motif Identification**
- **Comparative genomics resources available at UC Santa Cruz -- <http://genome.ucsc.edu>**
  - Genome-wide sequence availability
  - Gene prediction and identification Finding orthologous sequences in other species
  - Identifying conserved sequences
- **Insights from vertebrate genome sequence comparisons**
- **Multi-species sequence analysis**

## Insights from Human-Rodent Sequence Comparisons



Nature 420:520, 2002



Nature 428:493, 2004

- 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
- See Jan 2003 & April 2004 issues of *Genome Research*

## 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

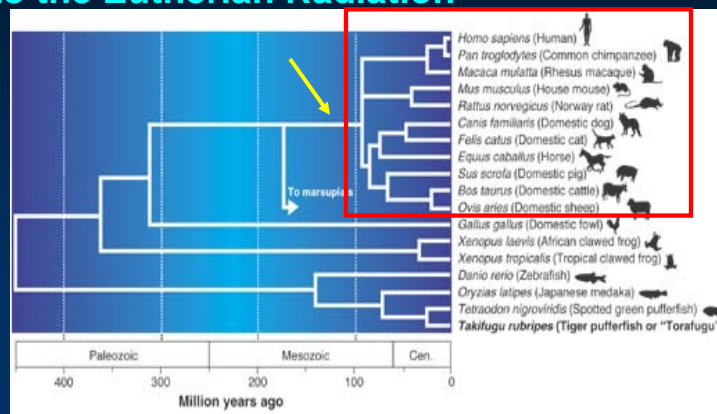
- Third position of codons which can be any base and code for the same amino acid

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

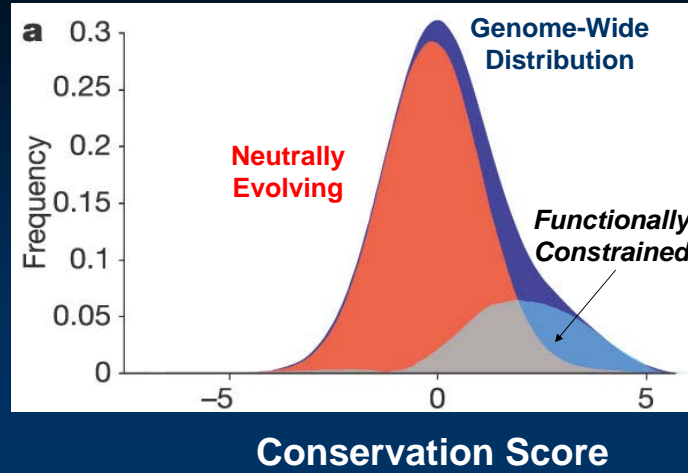
- Ancient Relics of Transposons Inserted Prior to the Eutherian Radiation



Adapted from Hedges & Kumar, *Science* 297:1283-5

# Determining the Fraction of Sequence Under Purifying Selection

Adapted From Figure 28, *Nature* 420:553



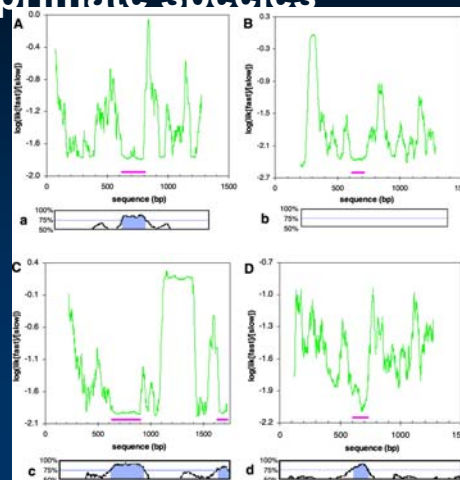
## Outline

- Fundamental concepts of comparative genomics
- Alignment and visualization tools
  - Pair-wise and multi-species methods
  - Combining with transcription factor binding site data
- Motif Identification
- Comparative genomics resources available at UC Santa Cruz -- <http://genome.ucsc.edu>
  - Genome-wide sequence availability
  - Gene prediction and identification Finding orthologous sequences in other species
  - Identifying conserved sequences
- Insights from vertebrate genome sequence comparisons
- Multi species sequence analysis

## Phylogenetic Shadowing

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

- Identifying sequence **differences** between multiple primate species





# Multi-Species Comparative Sequence Analysis

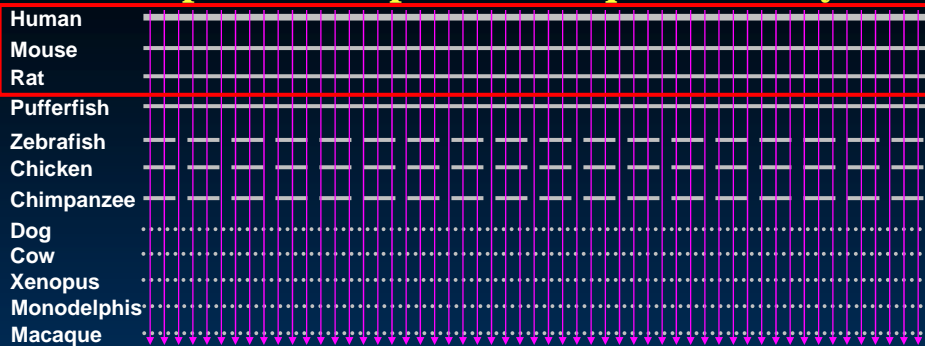
## Comparative analyses of multi-species sequences from targeted genomic regions

J. W. Thomas<sup>1,4</sup>, J. W. Touchman<sup>1,2,4</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>1</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>3</sup>, S. Schwartz<sup>2</sup>, L. Elnitski<sup>3</sup>, J. R. Idol<sup>3</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. Haghghi<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. Tiongson<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>1</sup>, D. Haussler<sup>10</sup>, P. Green<sup>10</sup>, W. Miller<sup>9</sup> & E. D. Green<sup>1,2</sup>

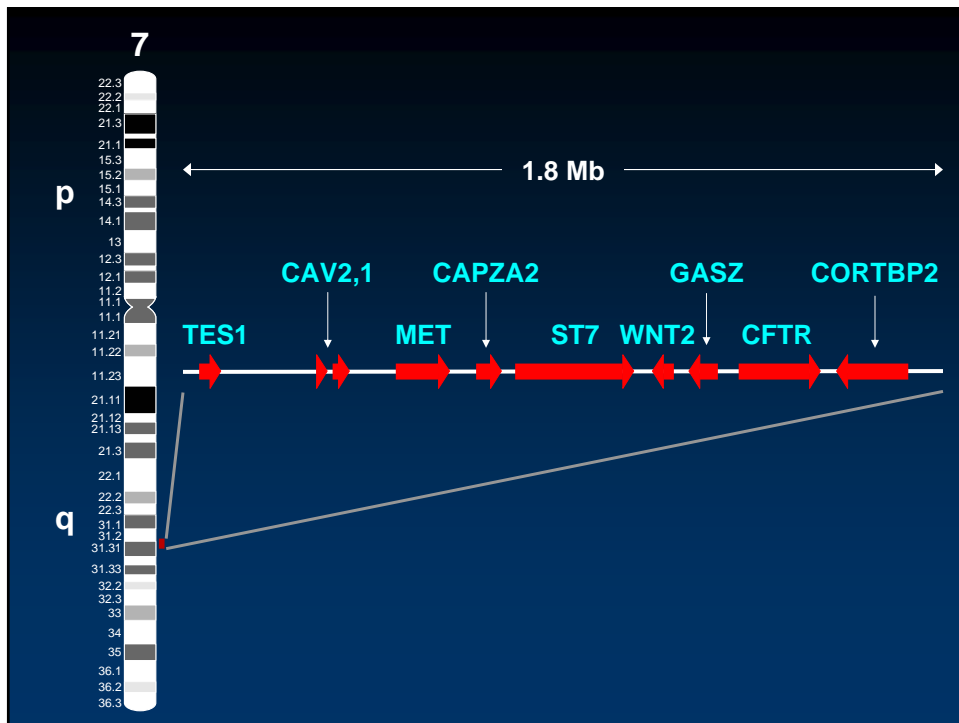
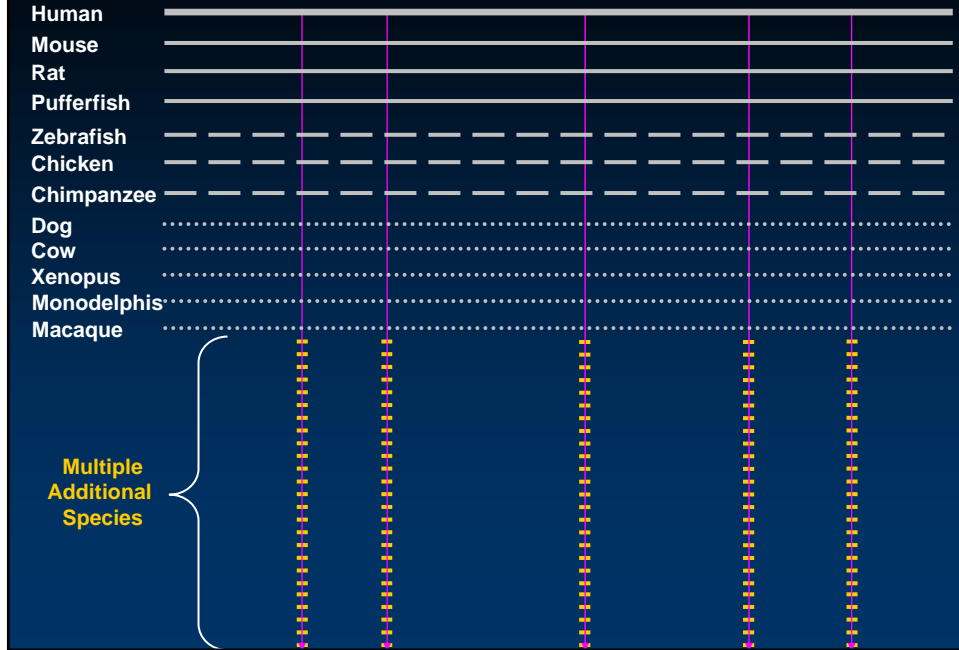


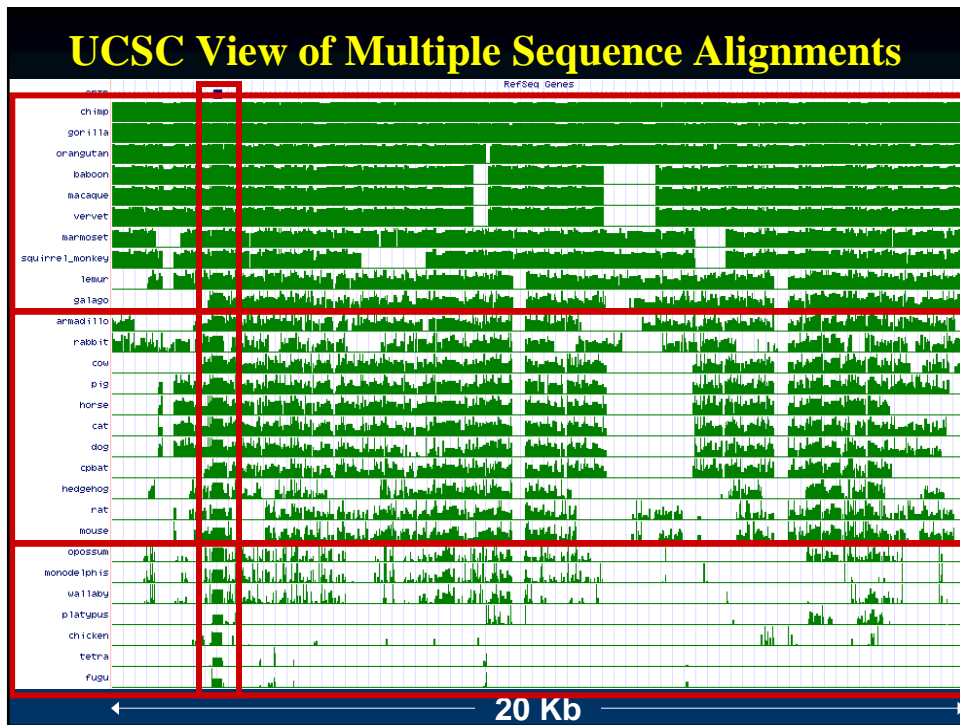
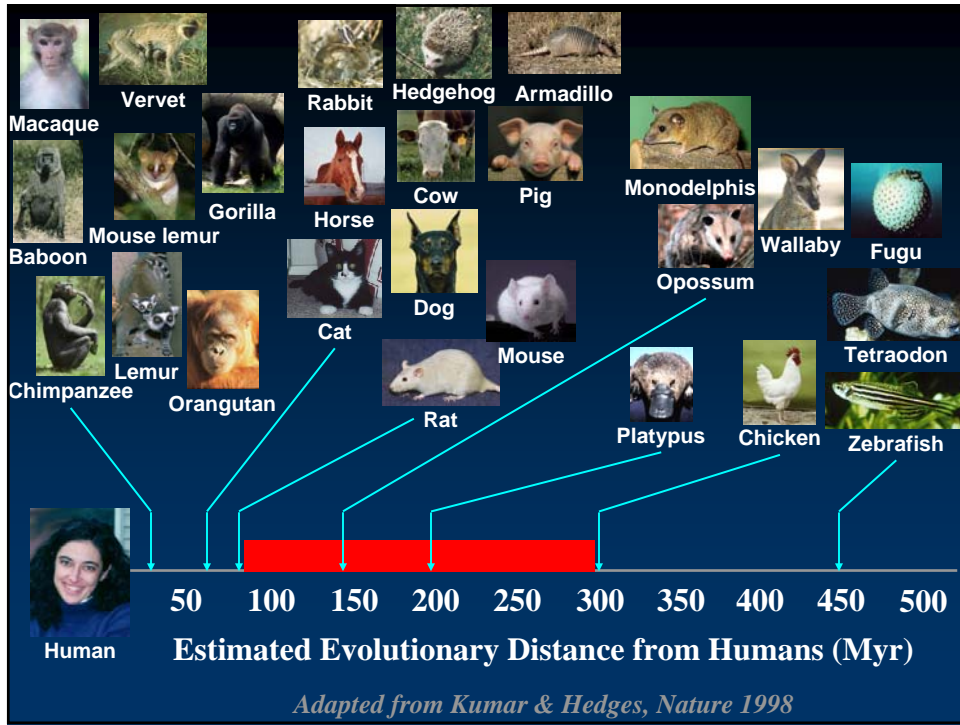
*Nature 424:788, 2003*

# Multi-Species Comparative Sequence Analysis



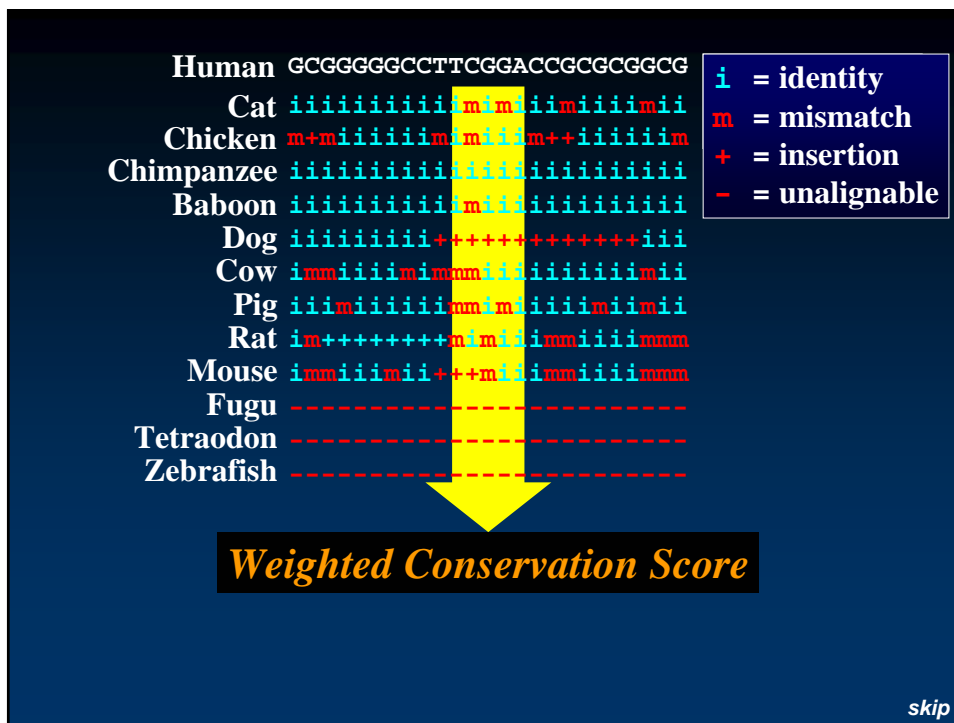
## Multi-Species Comparative Sequence Analysis



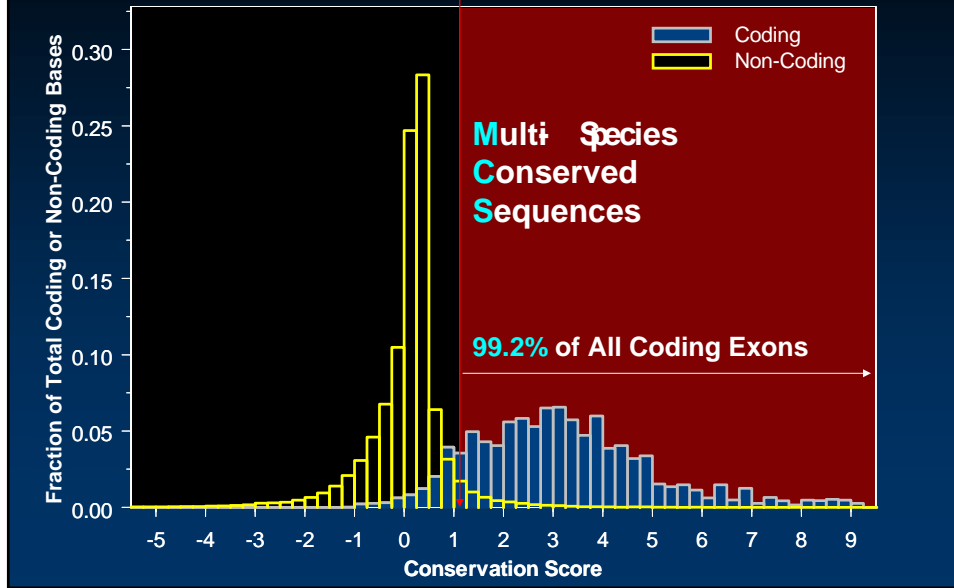


## Multi-Species Weighted Conservation Score

- 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



# Multi-Species Conservation Score Distribution



## Multi-Species Conservation Score

Article

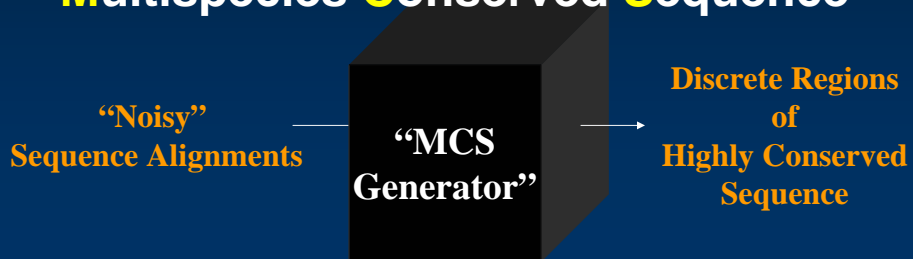
**Identification and Characterization of Multi-Species Conserved Sequences**

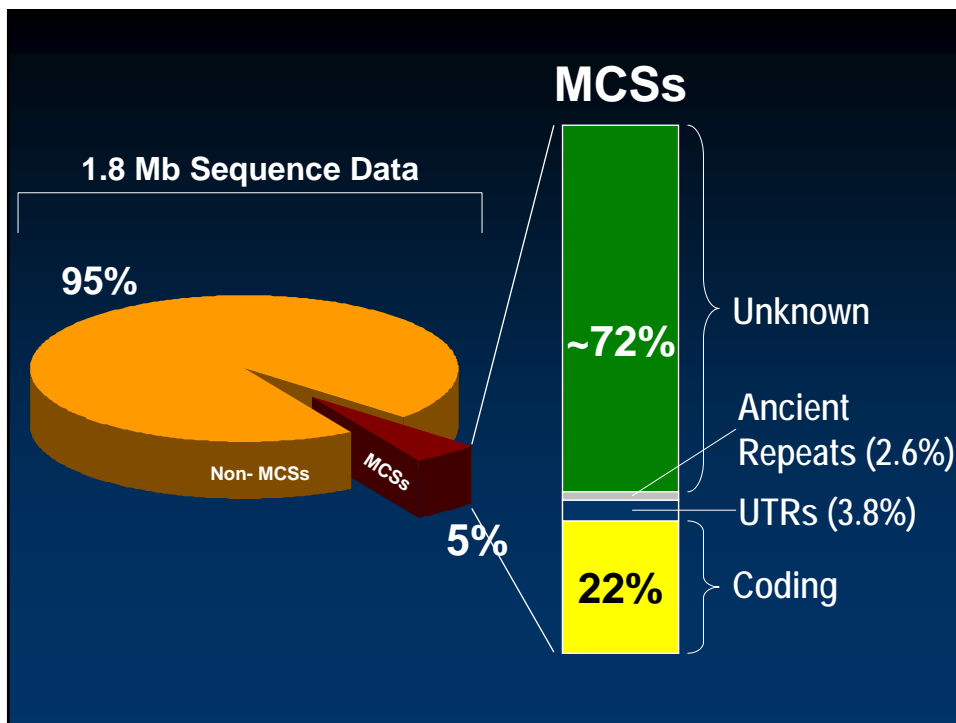
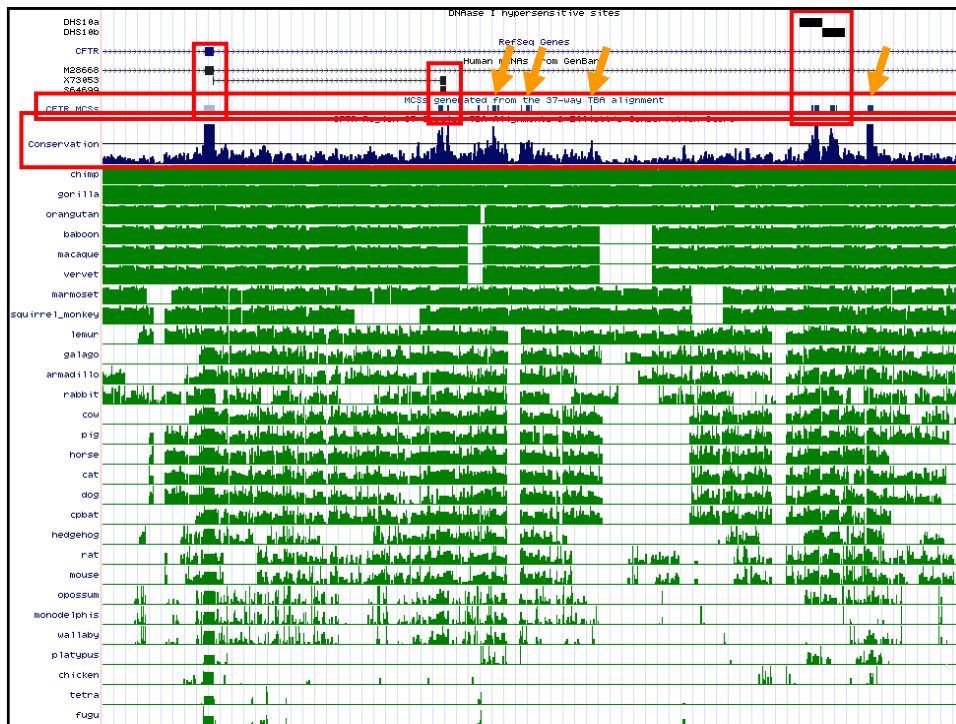
Elliott H. Margulies,<sup>1</sup> Mathieu Blanchette,<sup>3</sup> NISC Comparative Sequencing Program,<sup>1,2</sup> David Haussler,<sup>3,4,5</sup> and Eric D. Green<sup>1,2,5</sup>

*Genome Research* (2003) 13:2507-2518

**MCS**

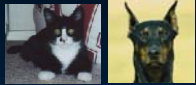
**Multispecies Conserved Sequence**





# Lineage-Specificity of MCSs in Mammals

## Carnivores



Cat Dog

## Artiodactyls



Cow Pig

## Rodents



Mouse Rat

## Monotreme



Platypus

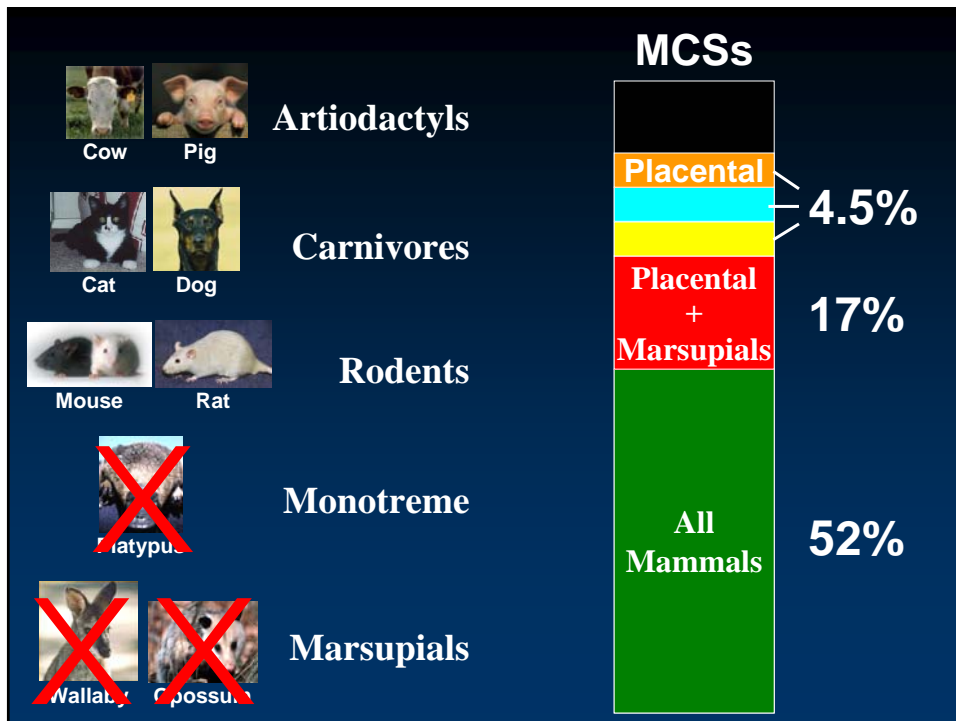
## Marsupials



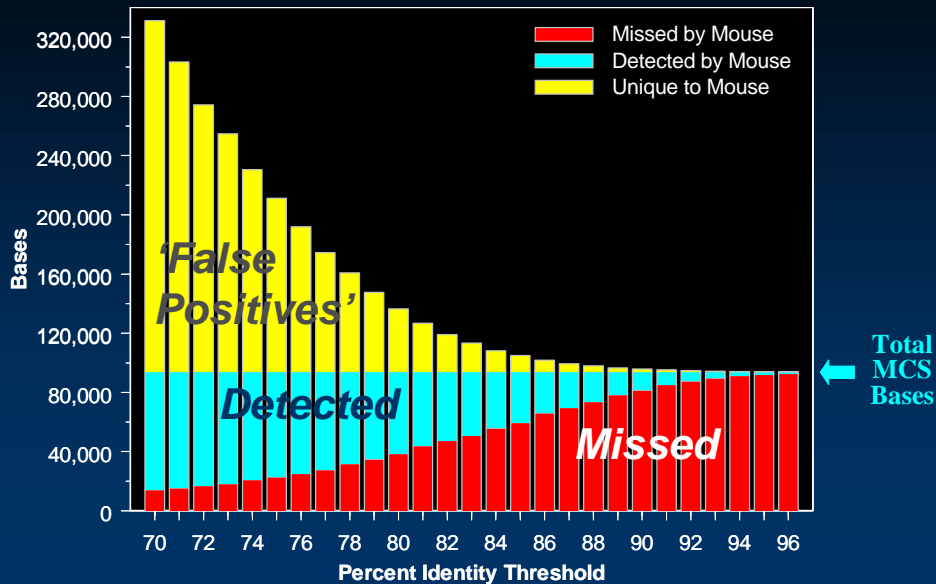
Wallaby Opossum

For each MCS:  
Catalog All the Species in Which it is Present

skip



## MCS Overlap with Mouse Alignments

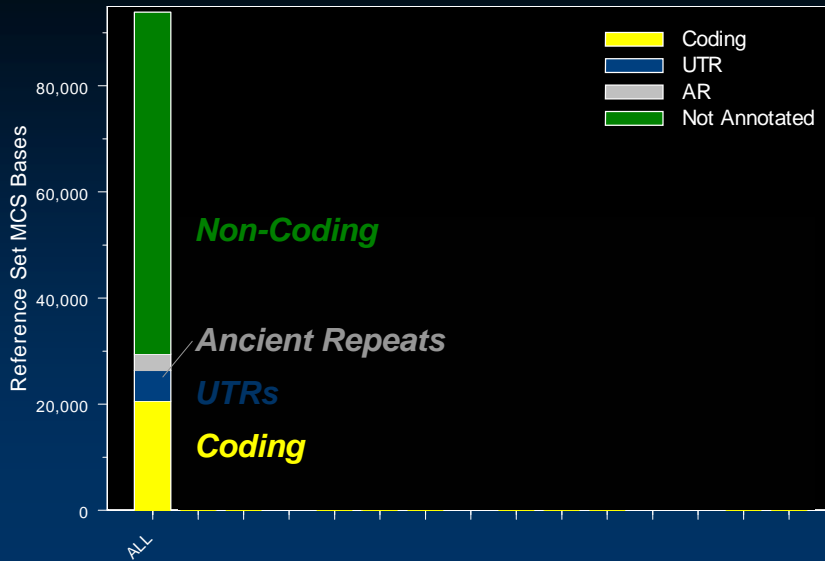


## Detection of MCSs with Different Species

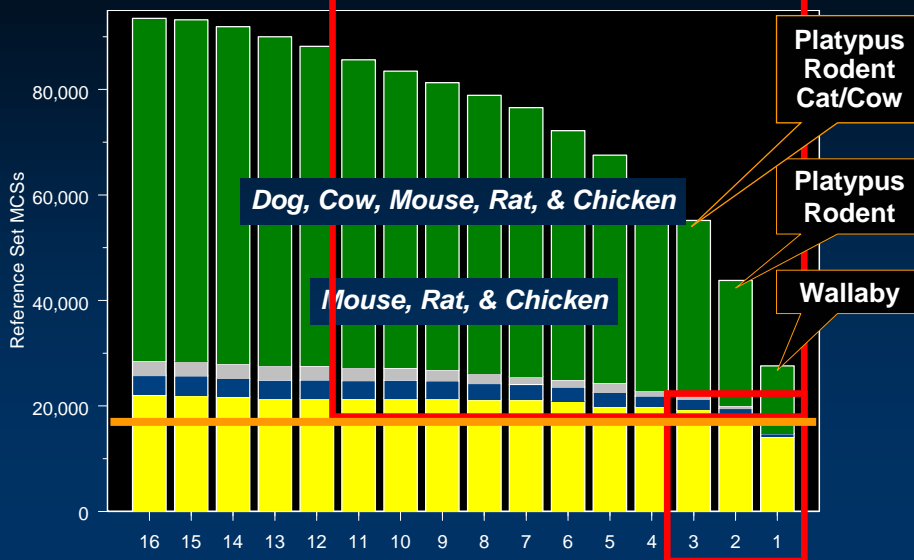
- Investigating the Relative Contribution of Different Species' Sequences to MCS Detection using More Quantitative Approaches
- Re-Compute Conservation Score for All\* Possible Subsets of Species
- Compare to a 'Reference Set' of MCSs
  - Generated with All Species
  - Surrogates for Conserved *Functional* Elements



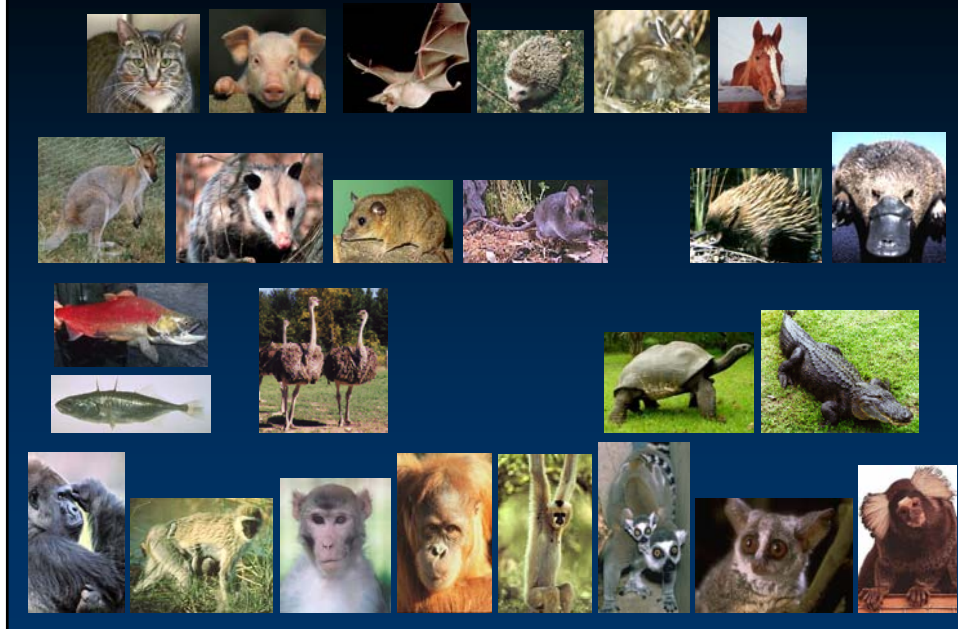
## Single Species Performance



## Best Performing Subsets



## More Species is Better



## MCS Detection and Sequence Quality

- To date, MCS detection has been with reasonably high-quality sequence
- What quality of sequence is desired for MCS detection — especially provided a set of high-quality reference sequences?

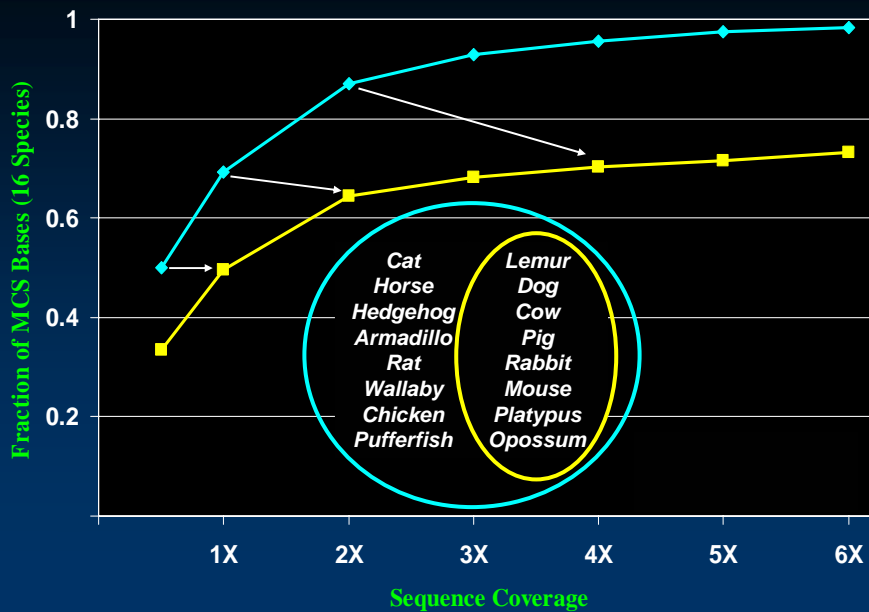
## MCS Detection and Sequence Quality

- What Tradeoffs are encountered between **sequence coverage** vs. **number of species**?

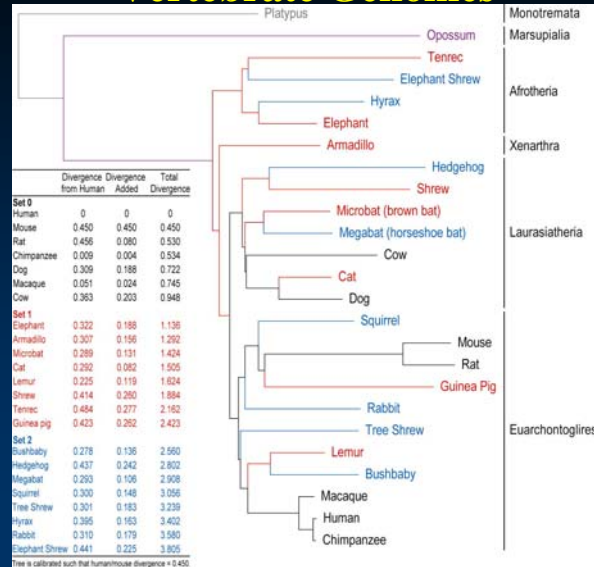
*Our Data Set Provides a Unique Opportunity to Explore these Questions*

- 1) Re-create 0.5X, 1X, 2X... Read-Coverage Datasets
- 2) Analyze for MCSs
- 3) Compare to “Finished” MCSs

## Sequence Coverage vs. MCS Detection



# Low-Redundancy Sequencing of Multiple Vertebrate Genomes



Margulies et al., (2005) *PNAS*, 102:3354-3359