

# **Bovine Annotation in Ensembl**

**Tim Hubbard**

**Wellcome Trust Sanger Institute**

**Bovine Genome Project Workshop**

**29-31st March 2005, Houston, Texas**

## Ensembl Genome Browser

### Search Ensembl

Search all species for  with

### About Ensembl



Ensembl is a joint project between [EMBL - EBI](#) and the [Sanger Institute](#) to develop a software system which produces and maintains automatic annotation on metazoan genomes. Ensembl is primarily funded by the [Wellcome Trust](#).

This site provides free access to all the data and software from the Ensembl project. Click on the species buttons to the right to browse the data.

Access to all the data produced by the project, and to the software used to analyse and present it, is provided free and without constraints. Some data and software may be subject to third-party constraints [\[details\]](#).

For all enquiries, please contact the Ensembl [HelpDesk](#) ([helpdesk@ensembl.org](mailto:helpdesk@ensembl.org)).

### Help and documentation

- Take the [Ensembl tour](#), go through a step-by-step [worked example](#), or read [these papers](#).
- For help on any web page click:
- There is also an [index](#) of help pages, and a set of guided [How do I...?](#) trails.

Display your own data in Ensembl

Questions or suggestions? Try the

Documentation (includes tutorial on direct data access & instructions for installing Ensembl on your own site)

Try the site map as a good starting point for exploring what Ensembl has to offer



### Species - Ensembl v29

<input type="button" value="Human"/>	<a href="#">NCBI 35</a>	<a href="#">Mar 05</a>
<input type="button" value="Mouse"/>	<a href="#">NCBI m33</a>	<a href="#">Feb 05</a>
<input type="button" value="Zebrafish"/>	<a href="#">WTSI Zv4</a>	<a href="#">Sep 04</a>
<input type="button" value="Rat"/> <small>pre!</small>	<a href="#">RGSC 3.1</a>	<a href="#">Jul 04</a>
<input type="button" value="Chicken"/>	<a href="#">WASHUC1</a>	<a href="#">Jul 04</a>
<input type="button" value="Mosquito"/>	<a href="#">MOZ 2</a>	<a href="#">Feb 05</a>
<input type="button" value="Fugu"/>	<a href="#">Fugu v2.0</a>	<a href="#">May 04</a>
<input type="button" value="Fruitfly"/>	<a href="#">BDGP 3.2.1</a>	<a href="#">Feb 05</a>
<input type="button" value="Chimp"/>	<a href="#">CHIMP1</a>	<a href="#">Mar 05</a>
<input type="button" value="Honeybee"/>	<a href="#">Amel1.1</a>	<a href="#">Sep 04</a>
<input type="button" value="Tetraodon"/>	<a href="#">TETRAODON7</a>	<a href="#">Sep 04</a>
<input type="button" value="Dog"/>	<a href="#">BROADD1</a>	<a href="#">Mar 05</a>
<input type="button" value="C. elegans"/>	<a href="#">WS 130</a>	<a href="#">Dec 04</a>
<input type="button" value="X. tropicalis"/>	<a href="#">JGI3</a>	<a href="#">Feb 05</a>
<input type="button" value="S. cerevisiae"/>	<a href="#">S228C</a>	<a href="#">Mar 05</a>
<input type="button" value="Cow"/> <small>pre!</small>	<a href="#">Btau_1.0</a>	
<input type="button" value="Opossum"/> <small>pre!</small>	<a href="#">BROAD0.5</a>	

### Data

Sequence similarity searches

Batch data/sequence retrieval

Ensembl Archive sites

Vertebrate Genome Annotation (VEGA)

Access to whole genome shotgun data (includes additional species)

Download Ensembl data via FTP

## Cow Genome Browser

### Ensembl Entry Points

Show scaffold

From

To

Lookup

BLAST your sequence

BLAST

For fast identity search try

SSAHA

### Cow Genome Project



Btau\_1.0 is a preliminary 3x assembly of the draft genome sequence of cow (*Bos taurus*), Hereford breed, using whole genome shotgun (WGS) reads from small insert clones. The project coordination and genome

sequencing and assembly is provided by the Human Genome Sequencing Center at Baylor College of Medicine.

The N50 size is the length such that 50% of the assembled genome lies in blocks of the N50 size or longer. The N50 of the contigs is 4.2 kb. The N50 of the scaffolds is 13.5 kb. The total length of all contigs is 2.26 Gb. When the gaps between contigs in scaffolds are included, the total span of the assembly is 2.34 Gb.

As this is a pre-release, the database does not contain any genes. Subsequent annotations including the ensembl genebuild are ongoing and will be added as soon as they are completed. Future assemblies will include WGS sequences with a larger insert sizes, BAC end sequences, BAC sequences, and marker information for more contiguous assembly, better scaffolding, and chromosome assignment.

Please refer to more [details](#), [conditions of use](#) and [credits](#).

### Example Data Points

This release of cow data is unassembled, so there are no chromosomes available to browse. Use the BLAST, SSAHA in the Entry Points section above to locate data.

A few example data points :

- ▶ Scaffold: [Scaffold 90001](#)
- ▶ Scaffold: [Scaffold 42](#)

The assembly comprises 795212 contigs (which are named by Genbank accession) organised into 449727 scaffolds.

### Documentation & Help

About Ensembl

Home

For context-sensitive help on any web page click

Help

Questions or suggestions? Try

Help Desk

## Pre-Ensembl

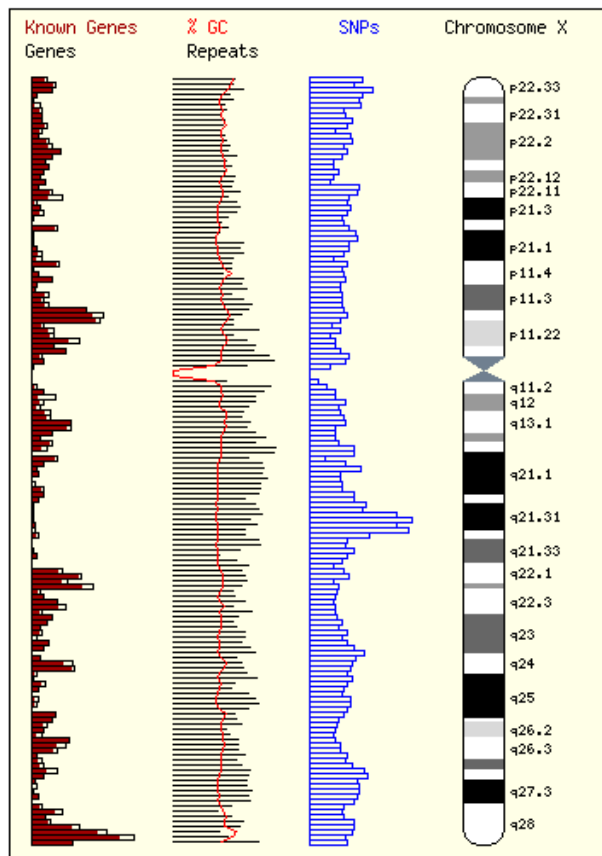
- **Provides rapid access to new genome assemblies before a full gene build is complete**
  - Shows raw alignments to mRNA and Proteins
  - Provides BLAST and SSAHA search services

Find All

Lookup [e.g. 6, 22]

Help

# MapView



### Chromosome X

Length: 154,824,264 bps  
 Gene Count: 931  
 Known Gene Count: 766  
 PseudoGene Count: 380  
 SNP Count: 320997

### Change Chromosome

Chromosome:

### Jump to Contigview

Click anywhere on the chromosome ideogram or one of the feature distribution plots to jump to a contig-level view of features at that point. Alternatively, you can jump to contigview between any two markers on this chromosome:

Between:    
 and:

[Display contig-level view between any two features.](#)

### Synteny

View Human Chr X vs

### OMIM Diseases

[Browse OMIM Diseases](#) on this chromosome.

### Map your data

[Map your own data](#) using KaryoView.

Display your own features





# Ensembl ContigView

[Home](#) > [Human](#) > [What's New](#) > [Text Search](#) > [BLAST Search](#) > [EnMap](#) > [Export Data](#) > [Download](#) > [Disease Browser](#) > [Docs](#)

Find    [e.g. [AC067852](#), [AP000869](#)]

### Chromosome 13

Chr 13

#### Overview

Rat synteny  
 Mouse synteny  
 Chromosome band

31.48 Mb 31.54 Mb 31.60 Mb 31.70 Mb 31.84 Mb 31.90 Mb 32.00 Mb 32.14 Mb 32.20 Mb 32.29 Mb

DNW(contigs)  
 Markers  
 Ensembl Genes  
 Gene legend

Legend:  
 CURATED KNOWN GENE (blue), CURATED NOVEL GENE (green), ENSEMBL PREDICTED GENES (KNOWN) (red), ENSEMBL PREDICTED GENES (NOVEL) (purple), CURATED PUTATIVE (yellow), CURATED PSEUDOGENE (orange), ENSEMBL PREDICTED PSEUDOGENE (pink), CURATED POLYCOMPLEX (cyan), ENSEMBL PSEUDOGENES (grey)

#### Detailed View

Jump to Chromosome: 13 bp  to

Features ▾ DAS Sources ▾ Repeats ▾ Decorations ▾ Export ▾ Jump to ▾ Image size ▾ Help ▾

Length  
 Human proteins  
 Proteins  
 Ensembl trans.  
 Vega trans.  
 Amino acids  
 Sequence  
 DNW(contigs)  
 Sequence  
 Amino acids  
 SNPs  
 Gene legend

Note: Amino acids only displayed for less than 4.5 kb. Sequence only displayed for less than 4.2 kb.

There are currently 58 tracks switched off, use the menu above the page to turn these on.

#### Basepair View

Length  
 Ensembl trans.  
 Vega trans.  
 Amino acids  
 Sequence  
 DNW(contigs)  
 Sequence  
 Amino acids  
 Restr. Enzymes  
 Gene legend

Legend:  
 CURATED KNOWN GENE (blue), CURATED NOVEL GENE (green), ENSEMBL PREDICTED GENES (KNOWN) (red), ENSEMBL PREDICTED GENES (NOVEL) (purple), CURATED PUTATIVE (yellow), CURATED PSEUDOGENE (orange), ENSEMBL PREDICTED PSEUDOGENE (pink), CURATED POLYCOMPLEX (cyan), ENSEMBL PSEUDOGENES (grey)

new
SETUP
CONFIG
RESULTS
DISPLAY

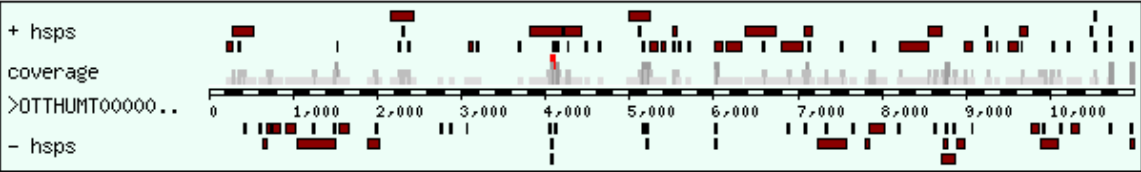
**Displaying OTTHUMT00000046000 sequence alignments vs Bos\_taurus LATESTGP database**

Showing top  alignments of 194, sorted by  refresh

[Alignment Locations vs. Karyotype \(click arrow to hide\)](#)

Karyotype unavailable for Bos\_taurus

[Alignment Locations vs. Query \(click arrow to hide\)](#)



Key (%ID): 0 - 20 20 - 40 40 - 60 60 - 80 80 - 100

[Alignment Summary \(click arrow to hide\)](#)

Select rows to include in table, and type of sort (Use the 'ctrl' key to select multiples) refresh

Query	Subject	Scaffold	Contig	Stats	Sort By						
Name Start End Ori	_off_ Name Start End	Name Start End Ori	_off_ Name Start End	E-val P-val %ID Length	>Score <E-val >E-val <P-val						
Links	Query	Scaffold	Stats								
	Start	End	Ori	Name	Start	End	Ori	Score	E-val	%ID	Length
<a href="#">[A]</a> <a href="#">[S]</a> <a href="#">[G]</a> <a href="#">[C]</a>	8202	8558	+	<a href="#">SCAFFOLD256786</a>	2755	3111	+	202	4.3e-146	89.11	358
<a href="#">[A]</a> <a href="#">[S]</a> <a href="#">[G]</a> <a href="#">[C]</a>	1049	1510	-	<a href="#">SCAFFOLD106828</a>	6949	7413	-	163	6.0e-214	83.73	467
<a href="#">[A]</a> <a href="#">[S]</a> <a href="#">[G]</a> <a href="#">[C]</a>	7842	8033	-	<a href="#">SCAFFOLD256786</a>	767	958	-	152	6.7e-85	94.79	192
<a href="#">[A]</a> <a href="#">[S]</a> <a href="#">[G]</a> <a href="#">[C]</a>	6381	6744	+	<a href="#">SCAFFOLD78058</a>	1429	1798	+	149	0.	84.88	377
<a href="#">[A]</a> <a href="#">[S]</a> <a href="#">[G]</a> <a href="#">[C]</a>	3823	4197	+	<a href="#">SCAFFOLD388566</a>	98	475	+	135	1.9e-99	83.81	383
<a href="#">[A]</a> <a href="#">[S]</a> <a href="#">[G]</a> <a href="#">[C]</a>	6796	7051	+	<a href="#">SCAFFOLD78058</a>	1850	2102	+	125	0.	87.16	257
<a href="#">[A]</a> <a href="#">[S]</a> <a href="#">[G]</a> <a href="#">[C]</a>	292	531	+	<a href="#">SCAFFOLD207125</a>	2446	2682	+	124	4.6e-61	87.92	240
<a href="#">[A]</a> <a href="#">[S]</a> <a href="#">[G]</a> <a href="#">[C]</a>	9877	10101	-	<a href="#">SCAFFOLD96316</a>	9	236	+	120	7.4e-75	87.93	232
<a href="#">[A]</a> <a href="#">[S]</a> <a href="#">[G]</a> <a href="#">[C]</a>	7238	7566	-	<a href="#">SCAFFOLD1352</a>	3494	3826	-	106	3.7e-54	82.84	338

**Summary**

- ▶ **setup**
  - Bos\_taurus
  - Genomic sequence
  - BLASTN
  - Low sensitivity
- ▶ **configure**
  - -E: 10
  - -B: 100
  - -filter: dust
  - -RepeatMasker
  - -W: 15
  - -M: 1
  - -N: -3
  - -Q: 3
  - -R: 3
- ▶ **results**
- ▶ **display**

ⓘ Not yet initialised

# Blast hit shown on scaffold

## [-] Scaffold SCAFFOLD106828

Scaffold SCAFFOLD106828

## [-] Overview



## [-] Detailed view

Jump to region: SCAFFOLD bp -17818 to 32180 Refresh

Window + Zoom - Window

Features ▾ Compara ▾ DAS Sources ▾ Repeats ▾ Decorations ▾ Export ▾ Jump to ▾ Image size ▾ Help ▾

Length 8.25 Kb

EMBL mRNAs

Proteins

Blast hits

Preliminary data

DNA(contigs) < ARFC01754195 < ARFC01150063

Repeats

Length 8.25 Kb

0 1,000 2,000 3,000 4,000 5,000 6,000 7,000 8,000

There are currently 12 tracks switched off, use the menus above the image to turn these on.

Ensembl Bos\_taurus SCAFFOLD106828:1-8253 Wed Mar 30 12:06:46 2005





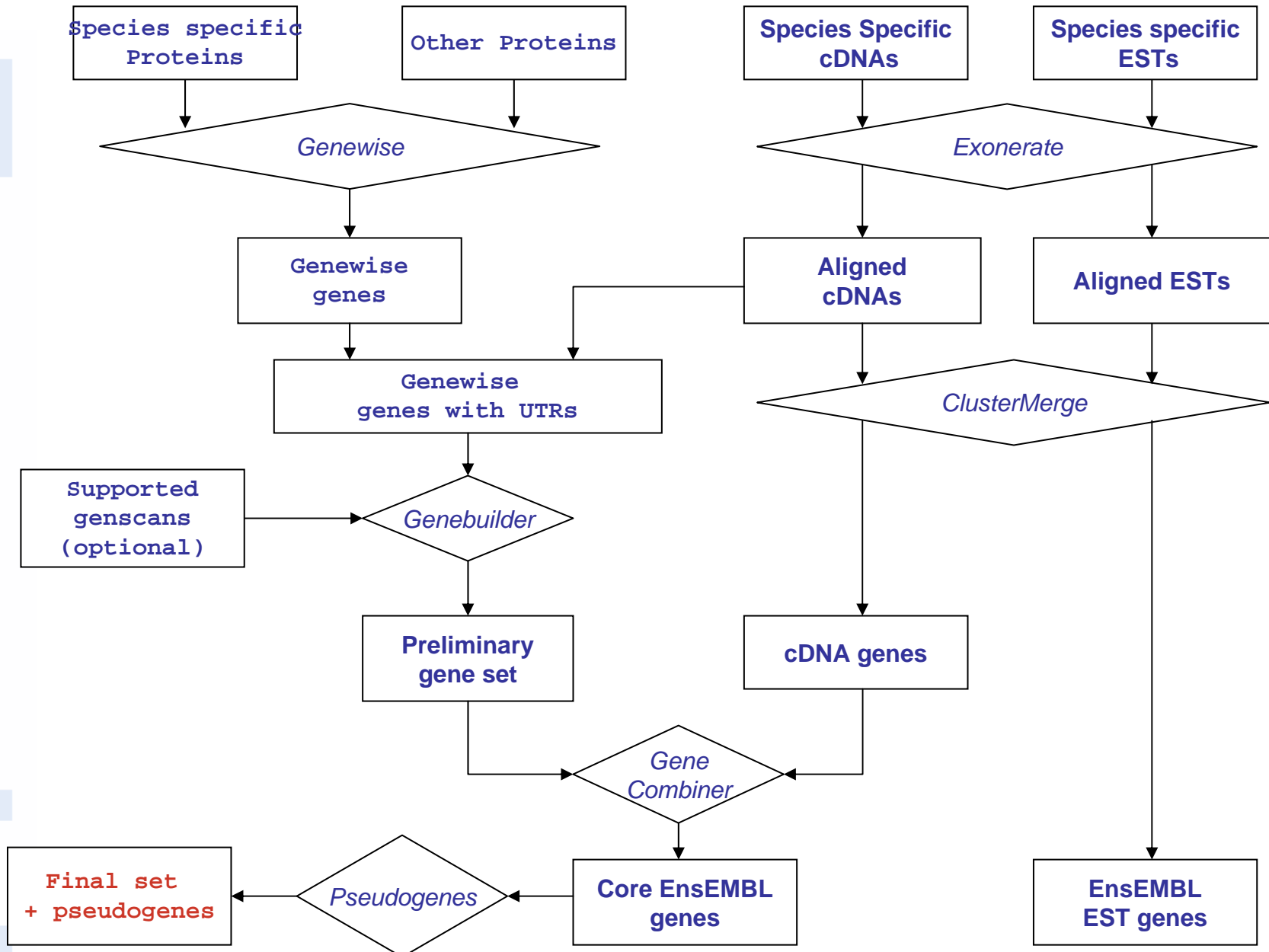
# Builds released since last SAB

Species	# Coding Genes	# Transcript	# Builds released since Nov 2003
Human	22218	33869	3 full, 4 patch, 1 pre
Mouse	24461	31535	2 full, 3 patch, 2 pre
Rat	22159	28545	1 full in progress, 2 pre
Dog	18201	30308	1 full, 1 pre
Chicken	17709	28416	2 full (1 release), 1 patch (ncRNA), 1 pre
Fugu	20796	33003	1 full (1 <sup>st</sup> in house)
Honey Bee	9671	16948	1 full (+ 1 in progress), 2 extended pre
Mosquito	14364	15802	1 patch (vectorbase team)
X. tropicalis	24405	52786	1 full, 1 pre (zfish group)
Zebrafish	23524	32062	2 full, 2 pre (zfish group)
Cow	-	-	1 pre, low coverage build development
Chimp	22475	43000	3 transfers from human (API team), 1 pre
Opossum	18936	32270	1 full in progress, 1 extended pre
C. intestinalis	-	-	1 full in progress
Tetraodon	28005	28005	1 import from Genoscope
D. melanogaster	13792	19178	1 import from Flybase
C. elegans	19765	24278	2 imports from Wormbase
S. cerevisiae	6680	6680	1 import from SGD
C. briggsae			1 full, REMOVED

# Stages in Genebuild Pipeline

- **the ‘raw compute’ (eg. blasts, repeat masking)**
  - All dependencies added at the start and run fully automated
- **A set of somewhat species specific steps requiring experimentation and assessment which we call the ‘genebuild’**
  - Stages added individually (genewise steps) or in small groups (cDNA build, targetted build) to the pipeline
  - Allows us to modify the procedure, but gives us the pipeline control benefits (automatic job submission, retries, flagging failures, batching jobs)
- **A set of post processing steps**
  - Pseudogene labelling
    - Pipelined during the last year. Can be customised.
  - Protein annotation
    - Fully automated like the ‘raw compute’ stage
  - Cross reference generation
    - Recently automated by the API team (separate system)

# Gene build Summary

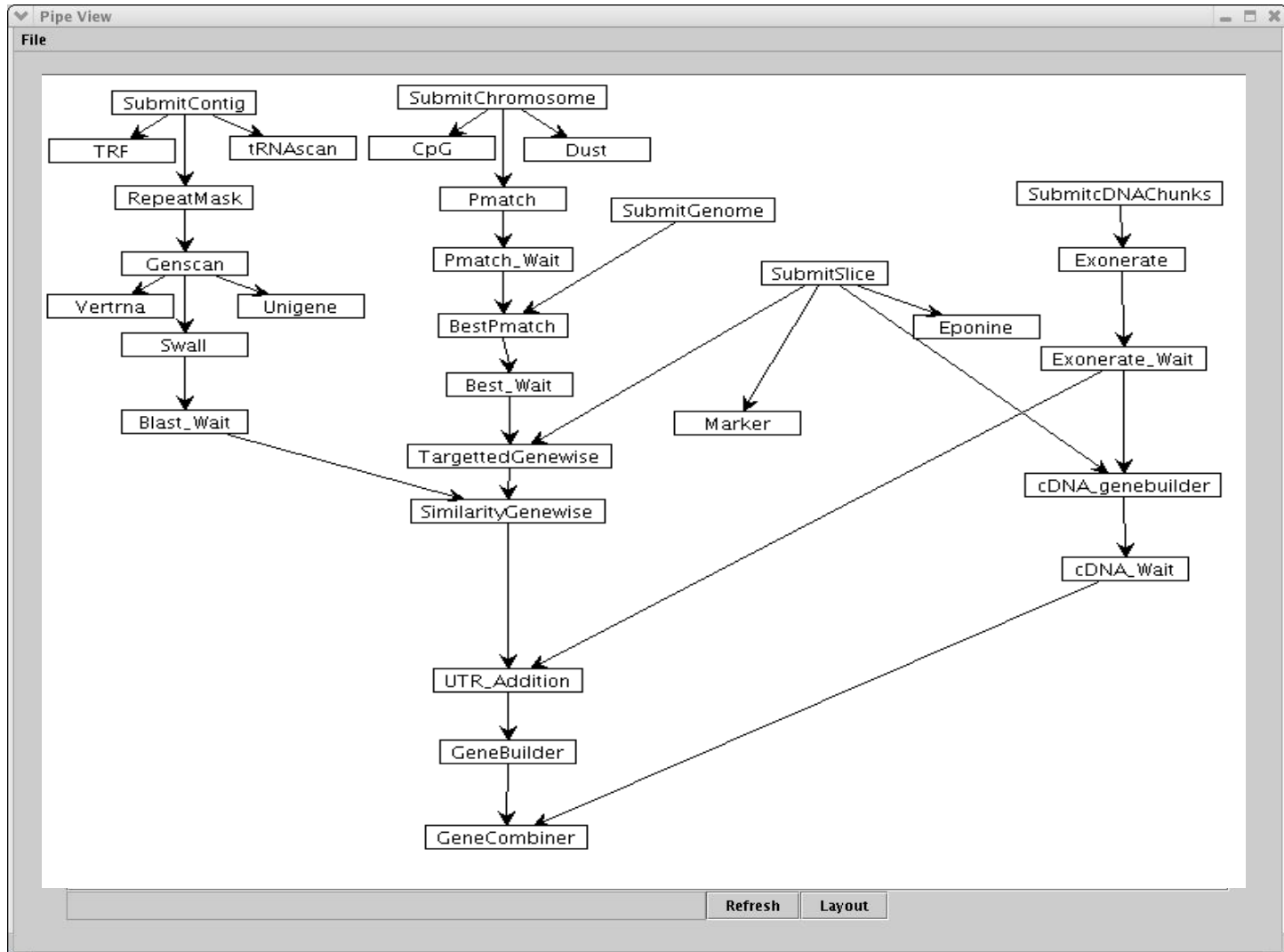




Ensembl



# A Genebuild Pipeline



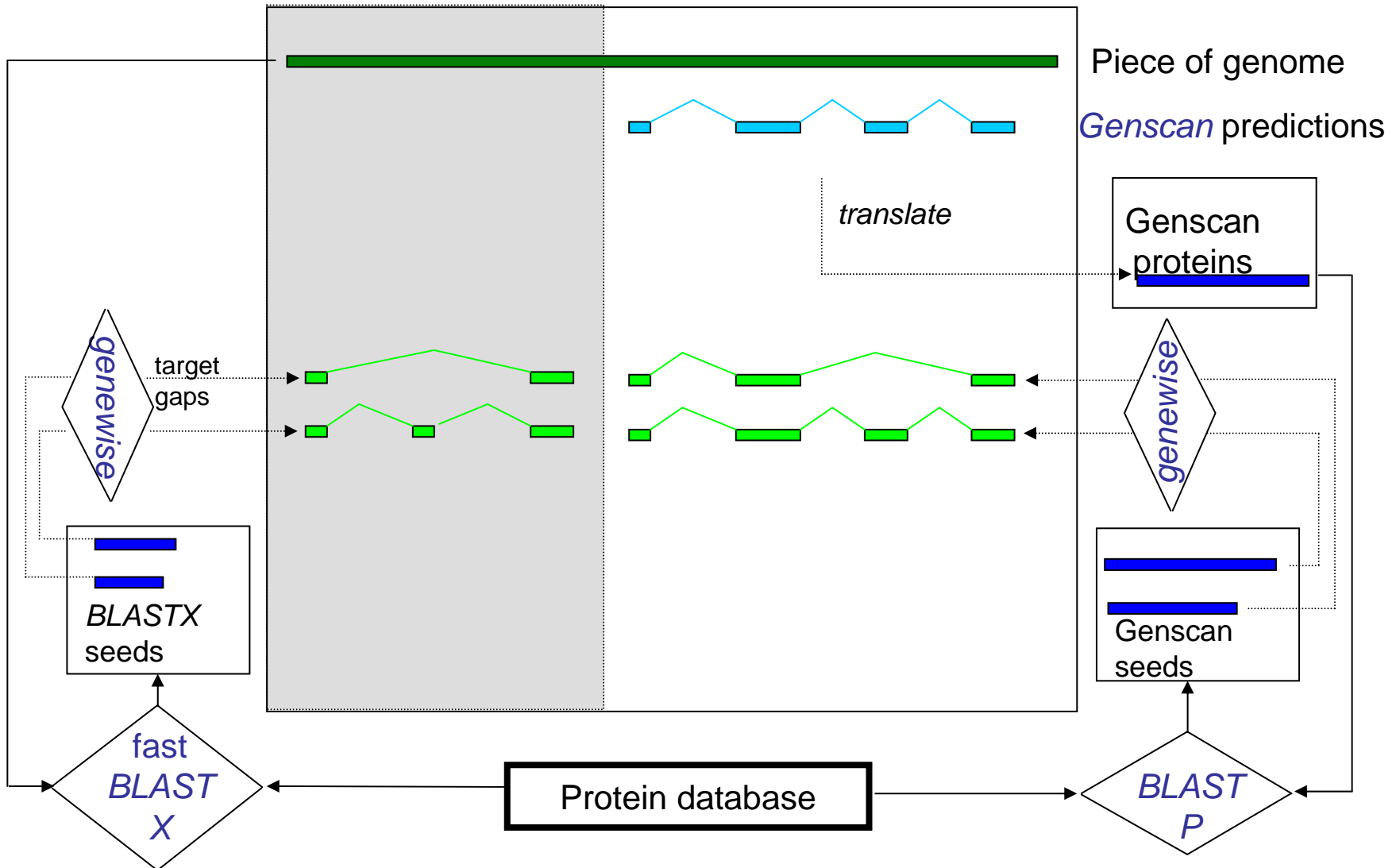


# Each organism has different characteristics and resources which affect genebuild design

- Evolutionary distance to closest well characterised species
- Number of species specific protein and cDNA sequences
- Number and quality of ESTs
- Genome size
- Coverage
- Availability of manual annotation data



# Fugu Genebuild



# Low Coverage Genomes

- **8 mammalian genomes sequenced to 2x are expected over the next 2-3 years.**
- **Ensembl is aiming to provide gene sets for these based on alignments to human, building predictions on scaffolds which align to genomic locations of human genes**
- **Test case**
  - Cow preliminary 3x assembly 449727 scaffolds, 795212 contigs
  - Good test case because 6x assembly will be available soon

# Cow 3x gene build

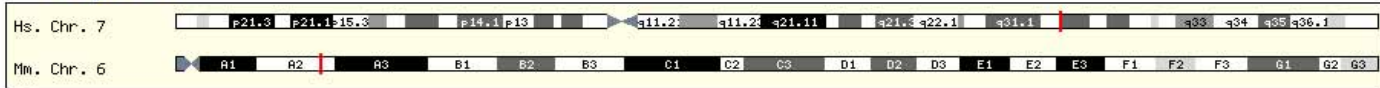
- **Define order of cow scaffolds from genome alignment to human (virtual assembly)**
- **Project human genes across alignment onto virtual cow assembly**
- **Identify Cow specific genes using standard pipeline**
  - coverage depends on coverage of Cow cDNAs)
  - Initial results: 80/2500 cDNA alignments gave matches that did not overlap with human projections. Some appear to be junk; some likely to be real.
- **Many issues resolving Cow-Human genome matches, assembly issues etc.**
- **Will investigate projection from Dog as Dog gene build/assembly improves**



# MultiContigView

Home Human What's New TextSearch BlastSearch MartSearch Export Data Download Disease Browser Docs Archive sites  
Find [All] [ ] [Lookup] [e.g. cancer, AC104620.5.1.155643, RH9632] [Help]

Top level

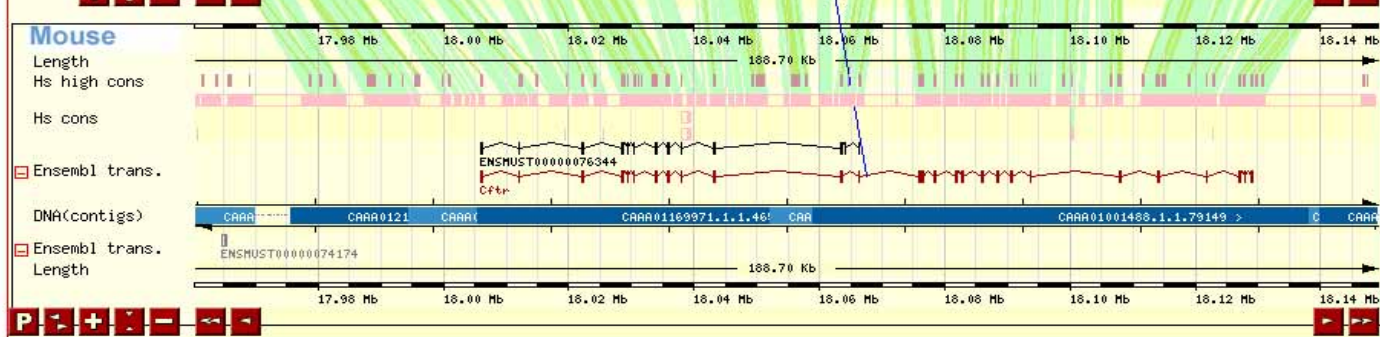
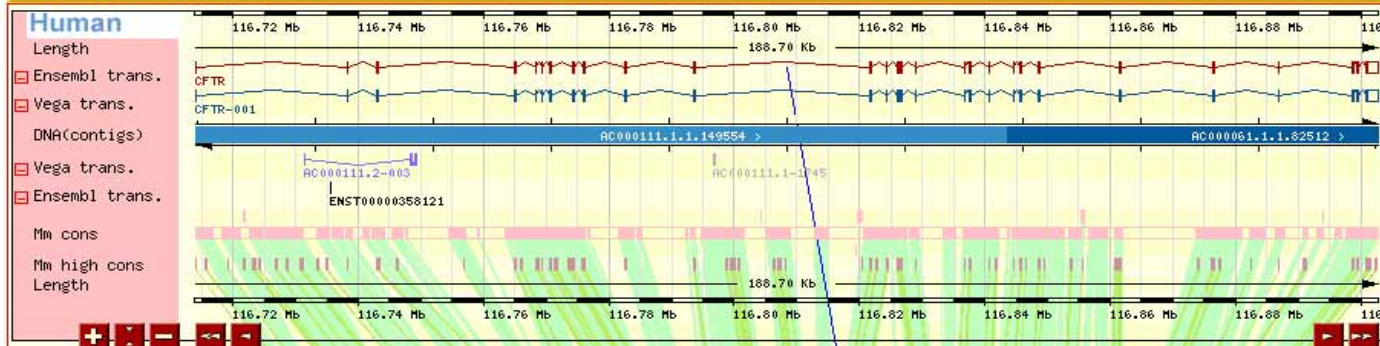


Overview

Detailed view

Jump to region: 7 bp [116713968] to [116902666] [Refresh]  
[<< 2 Mb] [1 Mb] [Window] [Zoom] [Window] [1 Mb] [2 Mb >>]

Features [v] Compara [v] Repeats [v] Decorations [v] Jump to [v] Export [v] Image size [v] [Help]



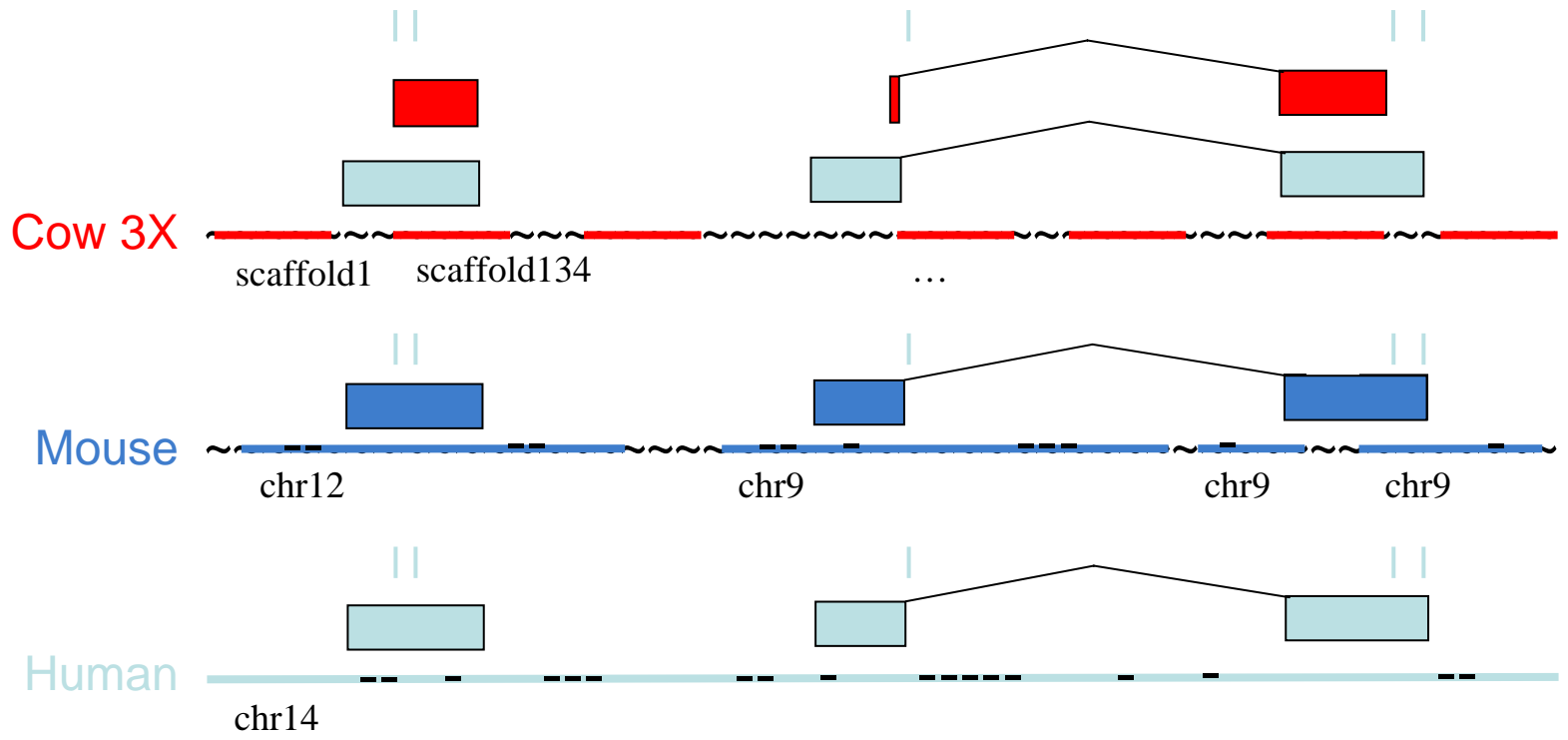
Ensembl



# AlignSlice API

- Using whole genome pairwise/multiple alignment data, to generate a reference coordinate system common to the aligned species in the genomic region of interest.
- Being able to project a transcript from one species to another through the alignment data
- Give gene context information across species, and more generally give annotation context information.
- Needed as a significant number of genomes are going to be 2X/3X where no sensible gene building is likely to give good quality gene set.

# AlignSliceView mock-up



# Acknowledgements

## **EBI Ensembl**

### **Comparison**

*Abel Ureta-Vidal*

*Jessica Severin*

*Cara Woodwark*

*Javier Herrero*

### **Mart**

*Arek Kasprzyk*

*Craig Melsopp*

*Glenn Proctor*

*Damian Smedley*

### **Core API**

*Arne Stabenau*

*Ian Longden*

### **Helpdesk**

*Xose Fernandez-S.*

*Damien Keefe*

*Guy Slater*

*Yuan Chen*

*Darin London*

*Ewan Birney*

## **Sanger Ensembl**

### **Genebuild**

*Steve Searle*

*Val Curwen*

*Dan Andrews*

*Laura Clarke*

*Kevin Howe*

*Vivek Iyer*

*Felix Kokocinski*

*Jan Vogel*

*Simon White*

### **Web**

*James Stalker*

*James Smith*

*Fiona Cunningham*

*Paul Bevan*

## **Zebrafish**

*Kerstin Jekosch*

*Mario Caccaro*

## **Vega**

*Patrick Meidl*

*Steve Trevanion*

## **Systems**

*Guy Coates*

*Tim Cutts*

*Mark Rae*

*Simon Kelley*

*Tim Hubbard*

*Tony Cox*

*Richard Durbin*

NCBI



# NCBI and the Bovine Genome

- Bovine Genome specific public resources at NCBI
- Bovine assembly and annotation activities
- Ensuring the Bovine Genome information is included in a broad information context

# Bovine Specific Resources

Genome Project - Microsoft Internet Explorer

File Edit View Favorites Tools Help

Address <http://www.ncbi.nlm.nih.gov/genome/sts/sts.cgi?uid=278327>

NCBI UniSTS Integrating Markers and Maps

PubMed Entrez BLAST OMIM Taxonomy Structure

Search UniSTS for [ ] Go

Entrez UniSTS

Help  
Query tips  
Submit  
Submit map  
FTP site  
Statistics

UniSTS:278327 [Links](#)

**EDG1**

*Bos taurus* chromosome 3  
*Homo sapiens* locus EDG1  
*Pan troglodytes* chromosome 1, locus LOC457063

Found by e-PCR in sequences from *Homo sapiens*, *Pan troglodytes* and *Sus scrofa*.

**Primer Information** ?

Forward primer: **TGGCCCTCTCAGACCTGTTG**  
Reverse primer: **TGGCGAGGAGACTGAACACG**  
PCR product size: not available  
GenBank Accession: **M31210**

**Bos taurus**

Name: EDG1

**Mapping Information** ?

EDG1 ILTX-2004 Map: Chr 3|LG1 [Map Viewer](#)  
Position: 138.21 (cR5000)

**Homo sapiens**

**Cross References** ?

Gene GeneID: 1901  
Symbol: EDG1  
Description: endothelial differentiation, sphingolipid G-protein-coupled receptor, 1  
Position: 1p21

Map Viewer - Microsoft Internet Explorer

File Edit View Favorites Tools Help

Address <http://www.ncbi.nlm.nih.gov/mapview/maps.cgi?TAXID=9913&CHR=3&BEG=0.00&END=312.6>

NCBI NCBI Map Viewer

PubMed Entrez BLAST OMIM Taxonomy Structure

Search [ ] Find Find in This View Advanced Search

**Bos taurus (cow) Build 1.1** [BLAST The Cow Genome](#)

Chromosome: [1](#) [2](#) [3](#) [4](#) [5](#) [6](#) [7](#) [8](#) [9](#) [10](#) [11](#) [12](#) [13](#) [14](#) [15](#) [16](#) [17](#) [18](#) [19](#) [20](#) [21](#) [22](#) [23](#) [24](#) [25](#) [26](#) [27](#) [28](#) [29](#) [X](#) [Y](#) [MT](#)

Master Map: ILTX [Summary of Maps](#) [Maps & Options](#)

Region Displayed: 0.00-312.81 cR

MARC	ILTX	Marker	cR	LG	LinkOut	MARC	X	Polymorphism
BMS871		UAP1	1.6	LG1	<a href="#">STS</a> <a href="#">acc</a>			N
		MUC1	55.6	LG1	<a href="#">STS</a> <a href="#">acc</a>			N
		S100A6	57	LG1	<a href="#">STS</a> <a href="#">acc</a>			N
		TUFT1	63	LG1	<a href="#">STS</a> <a href="#">acc</a>			N
		RME23	75.1	LG1	<a href="#">STS</a> <a href="#">acc</a>			Y
		PSMD4	90	LG1	<a href="#">STS</a> <a href="#">acc</a>			N
		NRAS	96.8	LG1	<a href="#">STS</a> <a href="#">acc</a>			N
		OVGP1	124.9	LG1	<a href="#">STS</a> <a href="#">acc</a>			N
		FLJ10330	128.4	LG1	<a href="#">STS</a> <a href="#">acc</a>			N
		BE217498	138.2	LG1	<a href="#">STS</a> <a href="#">acc</a>			N
		EDG1	138.2	LG1	<a href="#">STS</a> <a href="#">acc</a>			N
		CDC7L1	142.1	LG1	<a href="#">STS</a> <a href="#">acc</a>			N
		TGFBR3	142.1	LG1	<a href="#">STS</a> <a href="#">acc</a>			N
		ILSTS044	8.3	LG2	<a href="#">STS</a> <a href="#">acc</a>			N
		AW267062	23.3	LG2	<a href="#">STS</a> <a href="#">acc</a>			N
		NEDD5	63.5	LG2	<a href="#">STS</a> <a href="#">acc</a>			N
		MUF1	88.3	LG2	<a href="#">STS</a> <a href="#">acc</a>			N
		ACF7	93.6	LG2	<a href="#">STS</a> <a href="#">acc</a>			N
		MIP-T3	121.5	LG2	<a href="#">STS</a> <a href="#">acc</a>			N
		SFPQ	125.3	LG2	<a href="#">STS</a> <a href="#">acc</a>			N

Map Viewer Home  
Map Viewer Help  
Cow Maps Help  
Data As Table View  
[Maps & Options](#)  
Compress Map  
Region Show: [ ] [ ] Go  
out zoom in  
You are here: MARC  
default master

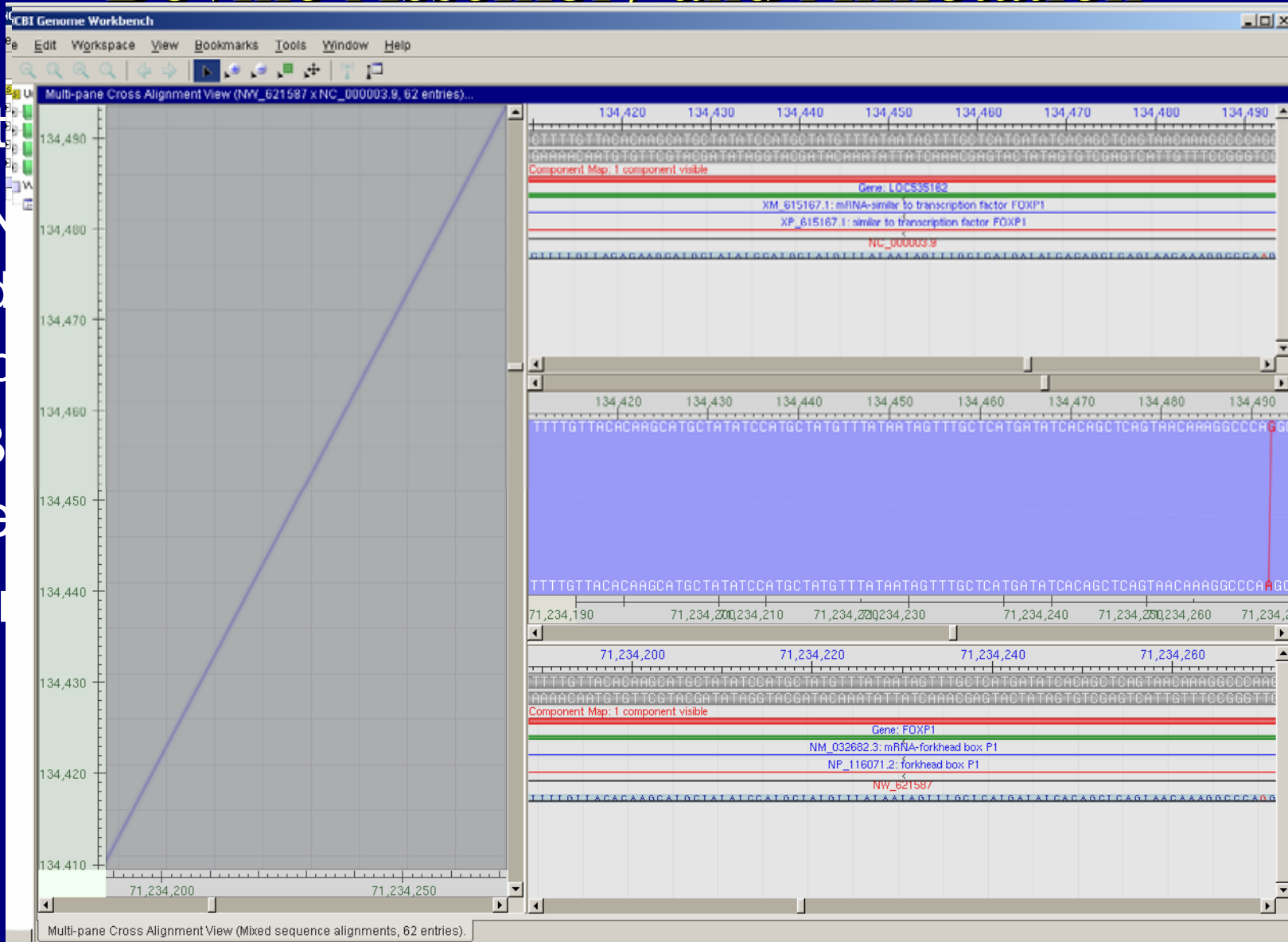


# Bovine Specific Resources

The image shows a screenshot of a web browser displaying two pages from the NCBI Entrez database. The main page is the Entrez Gene entry for Dopamine beta-hydroxylase (DBH) in *Bos taurus* (GeneID: 280758). The page includes a search bar, navigation tabs (All Databases, PubMed, Nucleotide, Protein, Genome, Structure, OMIM, PMC, Journals, Books), and a sidebar with various links like 'Entrez SITE MAP', 'Gene About Search', 'My NCBI', 'FAQ', 'FTP site', 'Related sites', 'Feedback', and 'Subscriptions'. The main content area provides details such as 'Genomic context: chromo:', 'Gene type: protein coding', 'Gene description: dopamin', 'RefSeq status: Provisional', 'Organism: *Bos taurus*', 'Lineage: Eukaryota; Meta', 'Eutheria; Laurasiatheria;', 'Gene aliases: DBH', 'Bibliography: Gen', 'PubMed links', 'GeneRIFs: 1. dioxygen and substrate ac', '2. The present study present', 'beta-hydroxylase (DBH; do', 'independently by pH and by', 'General gene informati', 'Markers (Sequence Tagg', 'DBH (e-PCR)', 'General protein inform', 'Name: dopamine beta-hy', 'NCBI Reference Sequ', 'mRNA Sequence', and 'Source Sequence'. A secondary window shows the PubMed abstract for a paper by Evans JP, Ahn K, and Klinman JP (J Biol Chem. 2003 Dec 12;278(50):49691-8), titled 'Evidence that dioxygen and substrate activation are tightly coupled in dopamine beta-monoxygenase. Implications for the reactive oxygen species.' The abstract text discusses the mechanism of copper mono-oxygenase, dopamine beta-mono-oxygenase (DbetaM), and the reoxidation of the enzyme-bound copper sites in the presence of O2.

# Bovine Assembly and Annotation

- Contigs
- mRNA
  - Ad
  - inc
- NCBI gene
- Bovine
- Trans



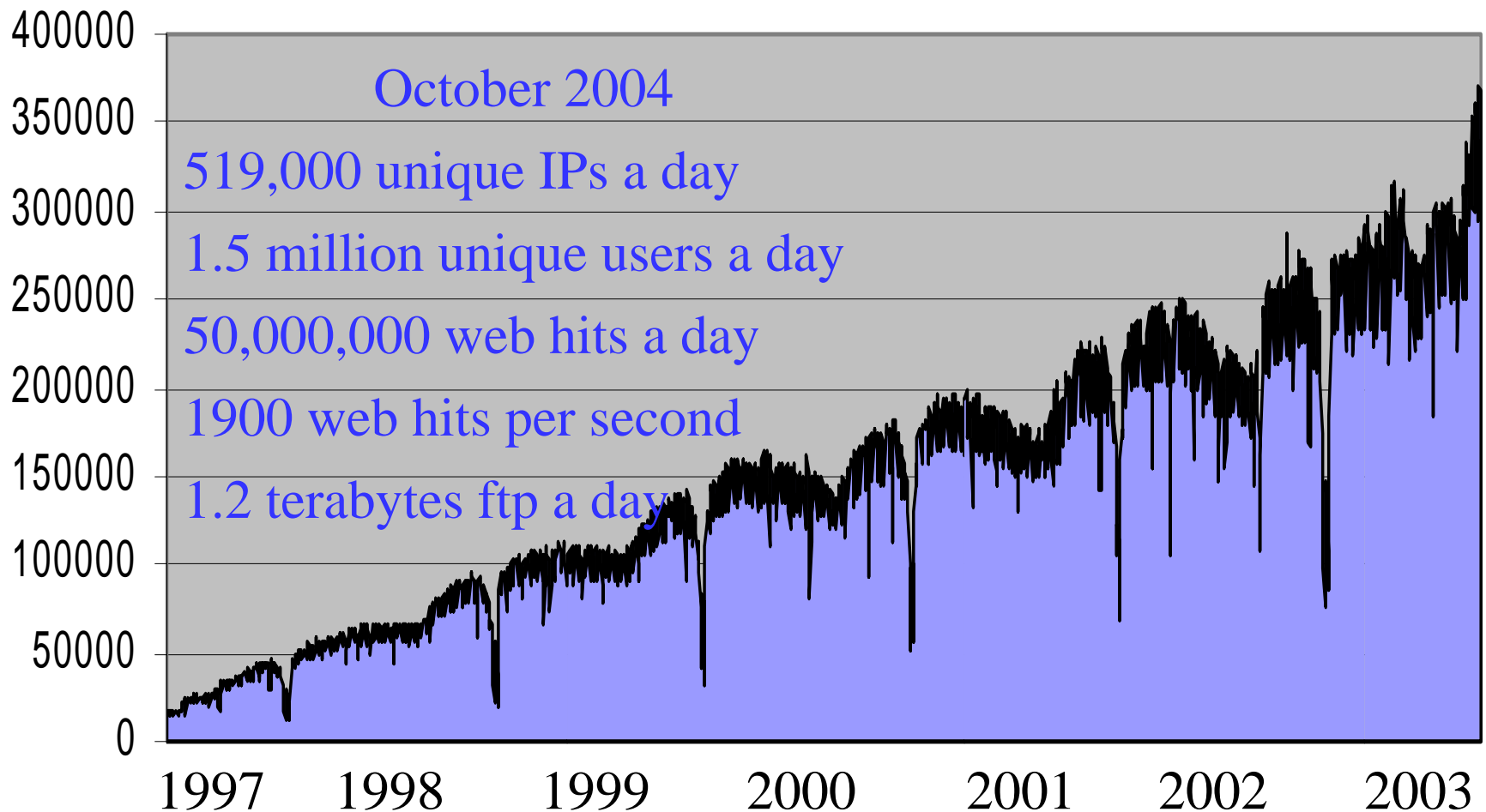
Multi-pane Cross Alignment View (Mixed sequence alignments, 62 entries).



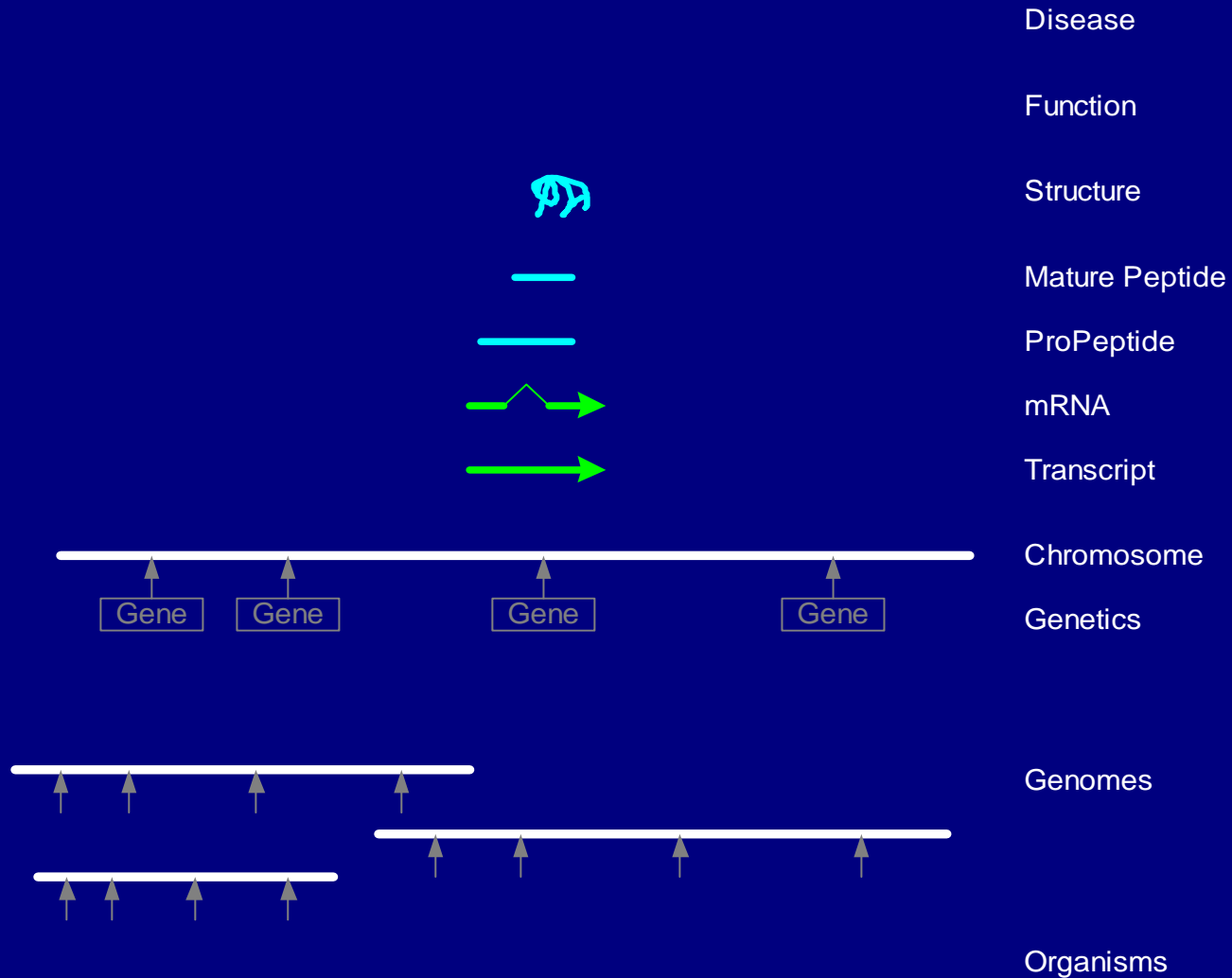
# Bovine Assembly and Annotation

- Contamination Screening – 167 contigs in 3x
- mRNA Based Super-scaffolds
  - Added 6297 aligned human cDNAs, 169 bovine
  - increased max contig length 2X
- NCBI genome pipeline placed 2340 known bovine genes (created 35,483 models)
- Bovine MGC picking, sequence evaluation
- Transcript anchored genomic alignment
- Traces, Assembly Database, SNPs

# Bovine in a Broader Context

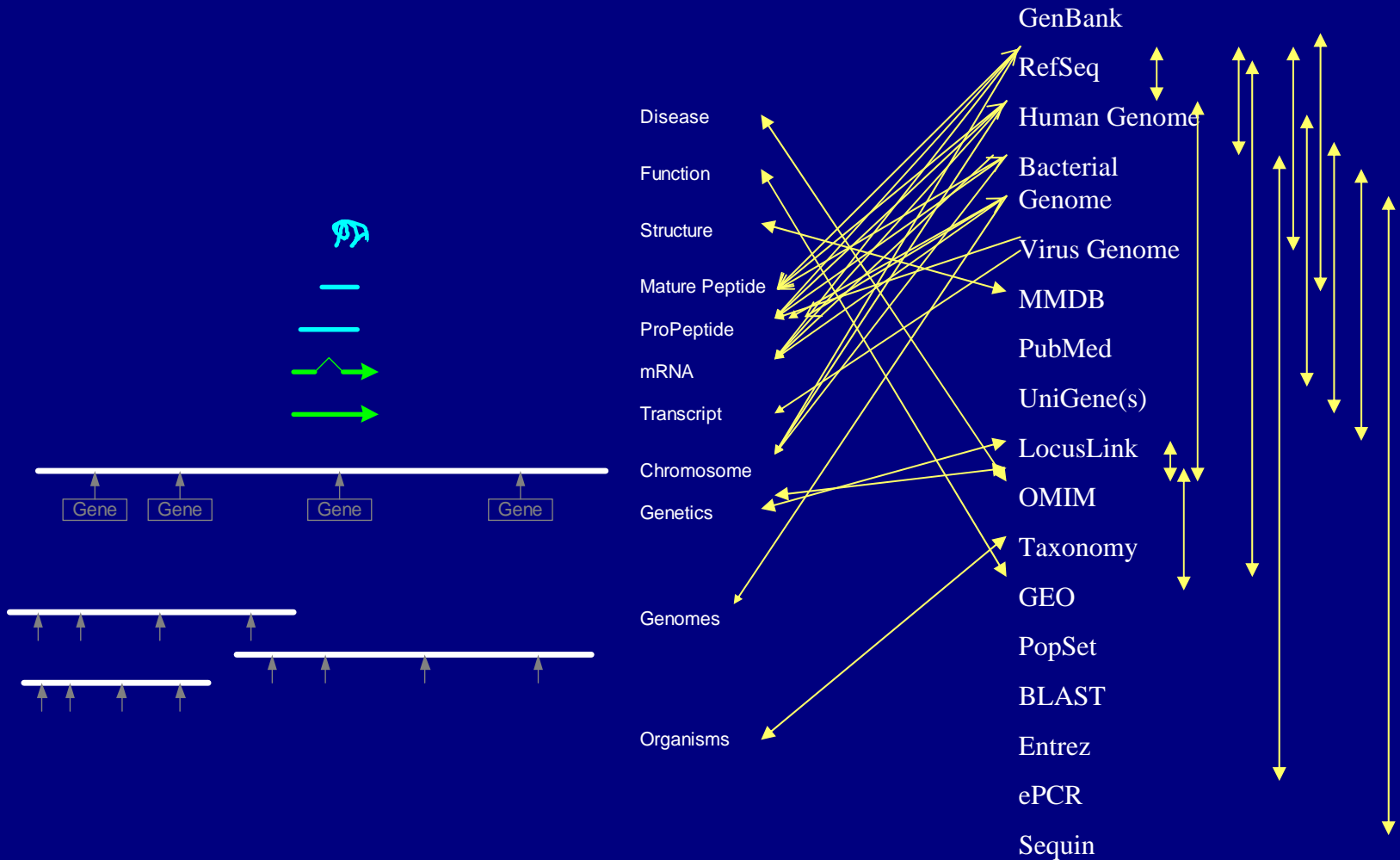


# Bovine in a Broader Context

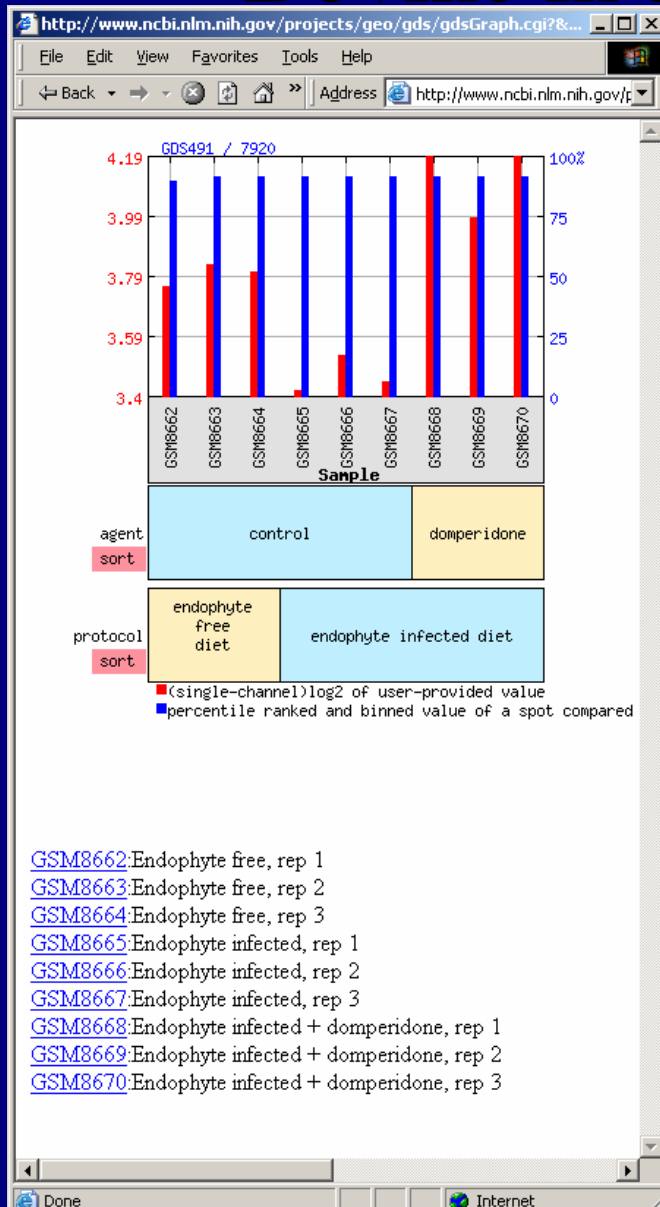




# Bovine in a Broader Context



# Bovine in a Broader Context



Entrez Protein - Microsoft Internet Explorer

Search Protein for [ ] Go Clear

Display Summary Show 20 Sort by Send to

All: 211 bacteria: 1 RefSeq: 46

Items 1 - 20 of 46 Page 1 of 3 Next

- [NP\\_037290](#) Reports BLink, Domains, Links  
 dopamine beta hydroxylase [Rattus norvegicus]  
 g|25742780|ref|NP\_037290.1|[25742780]
- [NP\\_620392](#) Reports BLink, Domains, Links  
 dopamine beta hydroxylase [Mus musculus]  
 g|20336728|ref|NP\_620392.1|[20336728]
- [NP\\_851338](#) Reports BLink, Domains, Links  
 dopamine beta-hydroxylase (dopamine beta-monoxygenase) [Bos taurus]  
 g|30794286|ref|NP\_851338.1|[30794286]
- [NP\\_000778](#) Reports BLink, Domains, Links  
 dopamine beta-hydroxylase precursor [Homo sapiens]  
 g|18426906|ref|NP\_000778.2|[18426906]
- [XP\\_520341](#) Reports BLink, Domains, Links  
 PREDICTED: similar to dopamine beta-hydroxylase precursor; dopamine beta-monoxygenase [Pan troglodytes]

# Bovine in a Broader Context

The image shows a screenshot of the NCBI Gene database interface. The main window displays the entry for **Gad1** (glutamate decarboxylase 1) in *Rattus norvegicus*. The GeneID is 24379, and the Locus tag is RGD:2652, RATMAP:33958. The official symbol is Gad1, and the name is glutamate decarboxylase 1 provided by Rat Genome Database. The RefSeq status is Provisional. The gene is located on chromosome 3, Maps: 3q21. The gene type is protein coding. The gene description is glutamate decarboxylase 1. The RefSeq status is Provisional. The organism is *Rattus norvegicus*. The lineage is Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi; Mammalia; Eutheria; Rodentia; Sciurognathi; Muridae; Murinae; Rattus. The gene aliases are Gad67. The summary states: SUMMARY: plays a role in the biosynthesis of gamma-aminobutyric acid [RGD].

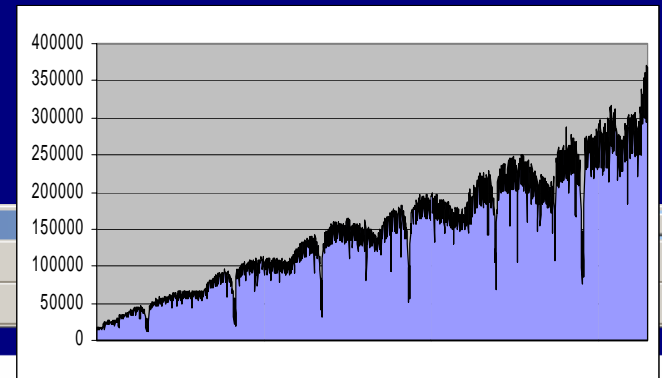
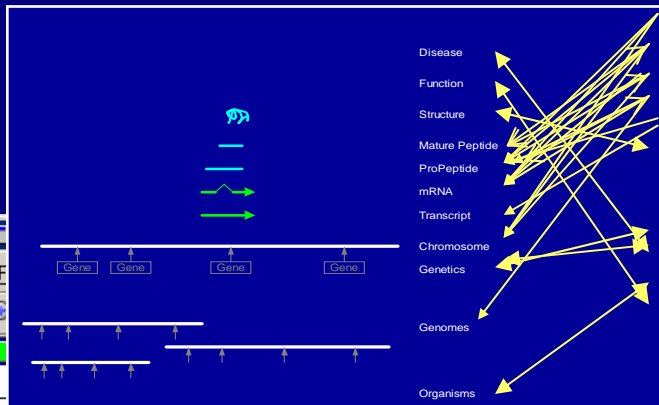
A callout box on the right side of the image highlights the identifier **LOC529488**. A red arrow points from this callout box to a red bar in the RefSeq track of the gene model, which is located between coordinates 25,995,417 and 25,995,710 on chromosome 3. This bar is labeled with the RefSeq ID **NC\_047655**.

The screenshot also shows a list of PubMed links for the gene, including:

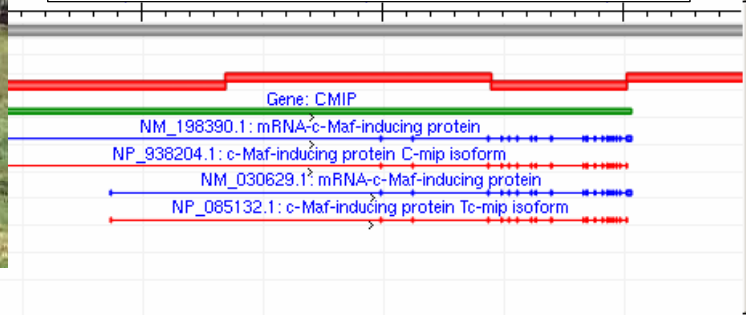
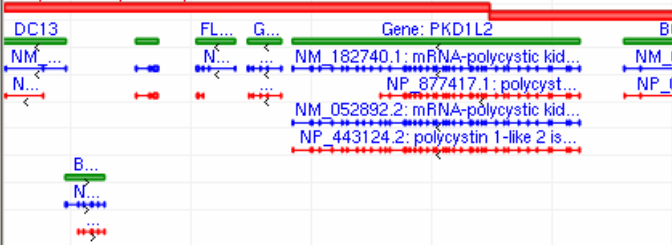
- 1. palmitoylation of GAD65 regulates the trafficking of the protein from Golgi membranes to an endosome in axons that is dependent on Rab5a and is required for the targeting of several synaptic vesicle protein clusters
- 2. The results revealed a significant decline in GAD-IR cells between middle and old age in CA1 but not in CA3.
- 3. With age, overall expression of GAD(67) mRNA decreases in the area surrounding the organum semicirculare and in the anteroventral periventricular nucleus. Young rats display a diurnal rhythm in GAD-IR cells.
- 4. Many areas of the brainstem and cerebellum involve double-labeled neurons with GLYT2 and GAD67 mRNAs and suggests that the corelease of glycine and GABA from single neurons is more widespread than has been reported.
- 5. The effects of a subthalamic nucleus lesion on 6-OHDA- or repeated D2 antagonist-induced changes in globus pallidus GAD(67) mRNA expression in parvalbumin[PV]+ and PV- neurons was examined in rats
- 6. The mRNA expression of GAD65 is up-regulated in the vestibular nuclei bilaterally 50 h after labyrinthectomy. In the flocculus, GAD65 mRNA expression is bilaterally up-regulated 50 h post-operatively.
- 7. acute furosemide treatment reduced the expression of GAD67 mRNA, the active holoenzyme predictive of GABA synthesis, within the lamina terminalis
- 8. in 6-hydroxydopamine-lesioned rats, GAD67 mRNA levels in striatonigral and striatopallidal pathways were selectively modified, and the modification correlated to dopamine agonist priming
- 9. The GAD67 and GAD65 is expressed in the same class of dopaminergic neurons in the GAD67-deficient mouse

# NCBI and the Bovine Genome

- Bovine Genome specific public resources at NCBI
- Bovine assembly and annotation activities
- Ensuring the Bovine Genome information is included in a broad information context



Component Map: 7 components visible



NCBI

WWW  
110101

# Bovine Genome in Context





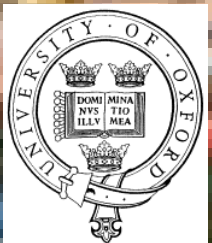
# ANALYSES FOR THE CATTLE GENOME PROJECT

**Chris Ponting**

**Leo Goodstadt**

**Caleb Webber**

**Andreas Heger**



**MRC**

Functional  
Genetics  
Unit

2001

2002

2003

2004

2005



*EVOLUTION*

Domain architectures

Protein Evolution  
Gene Duplication

Genome Remodelling

Gene Loss  
Gene Gain

*VARIATION*

SNPs

$K_S$  elevation

Allelic nulls

*DISEASE*

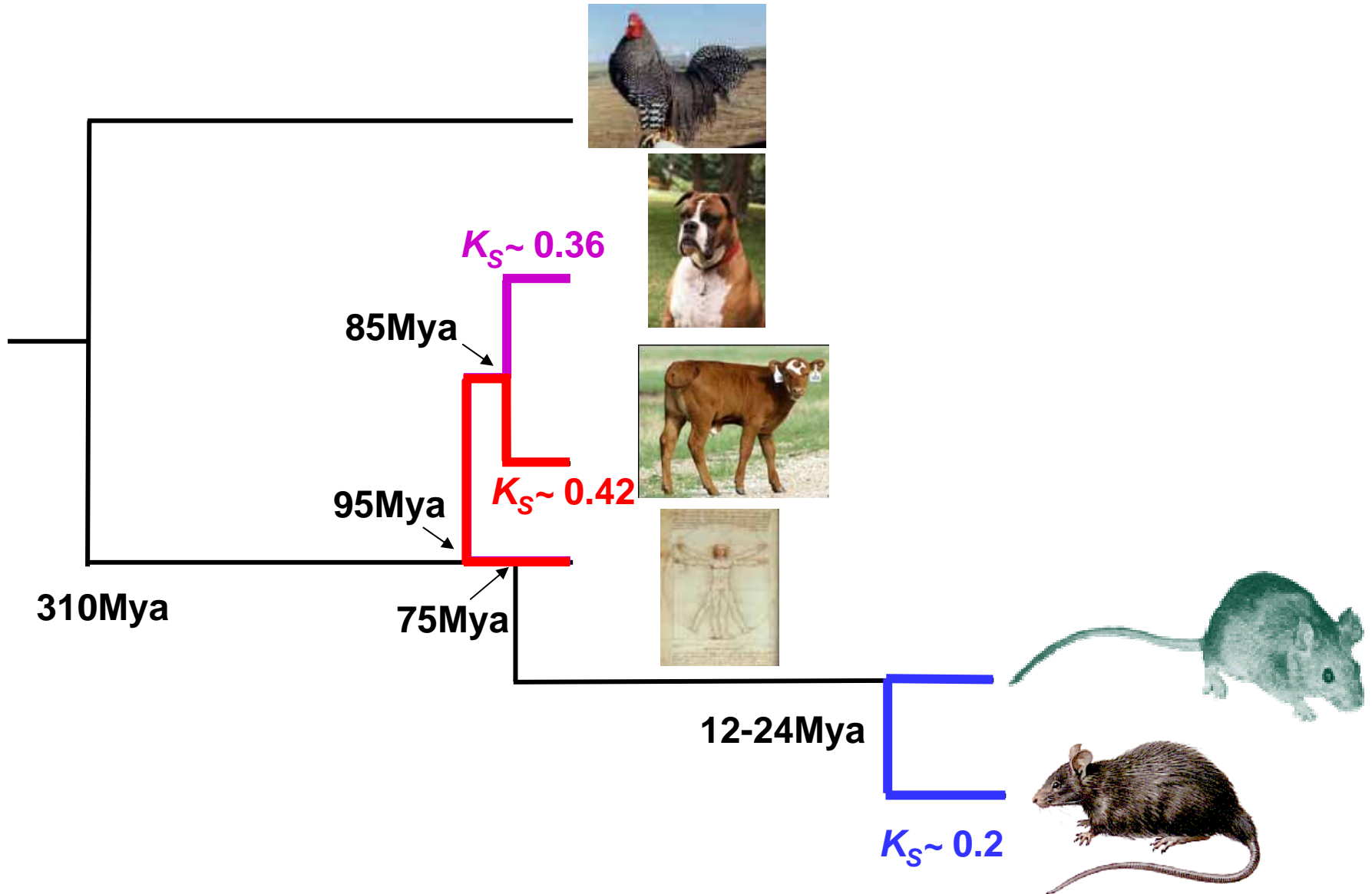
Paralogues

Conserved sites

$K_S$  elevation

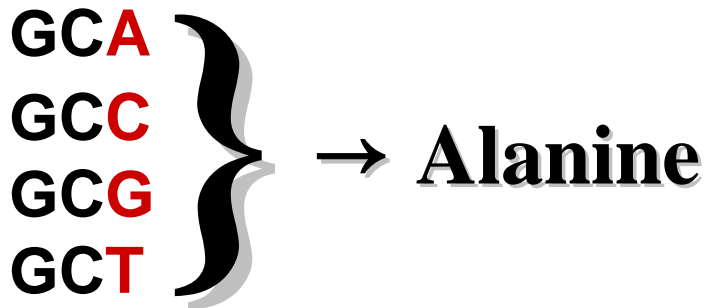


# Amniote Phylogeny



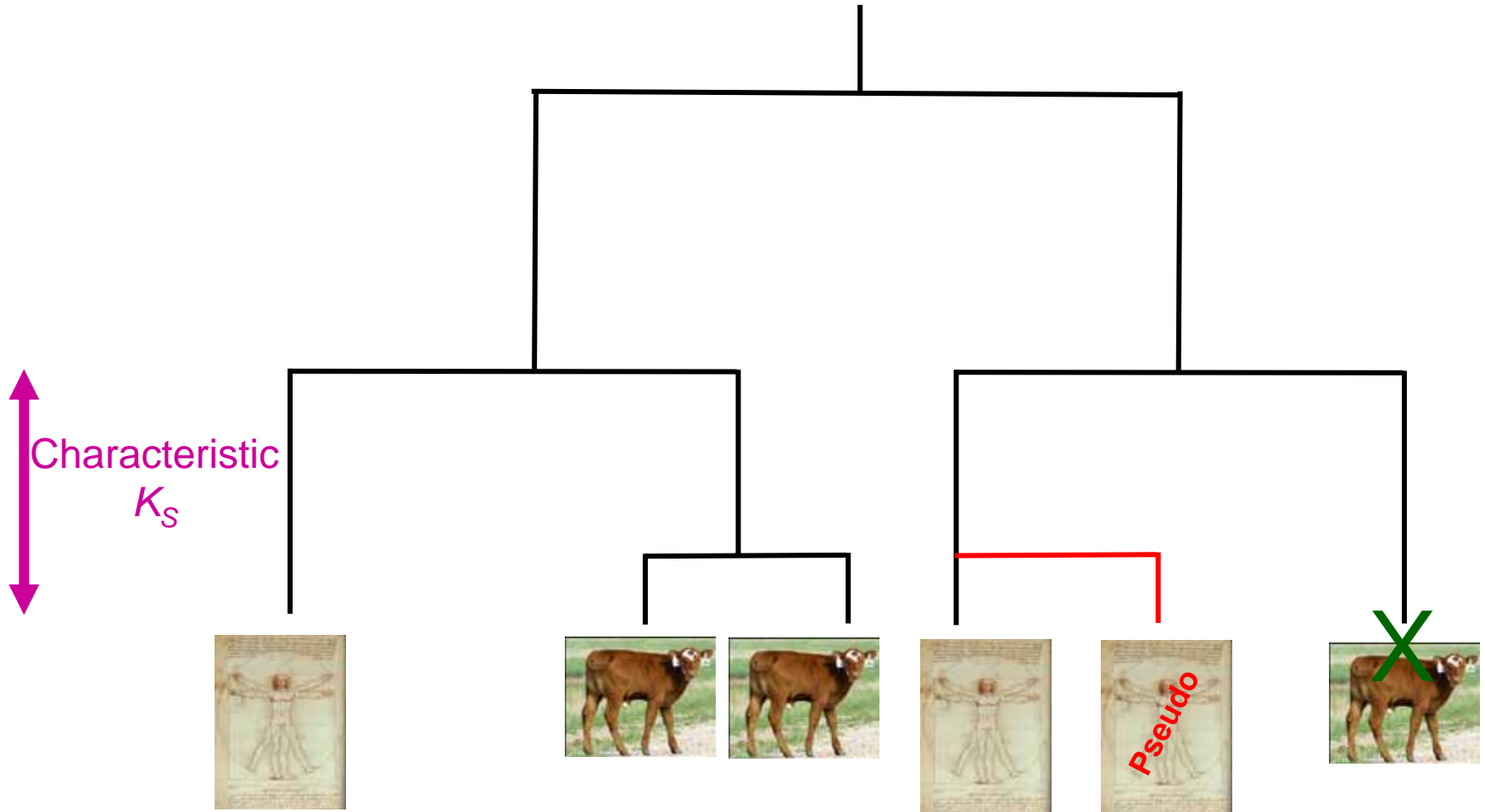
# $K_S$ : fraction of synonymous changes at synonymous sites

- Redundant genetic code, e.g.



- Third base of a codon “wobbles” without changing the translated amino acid
- $K_S$  measures neutral mutation rate in coding regions without selection

# Gene Phylogeny



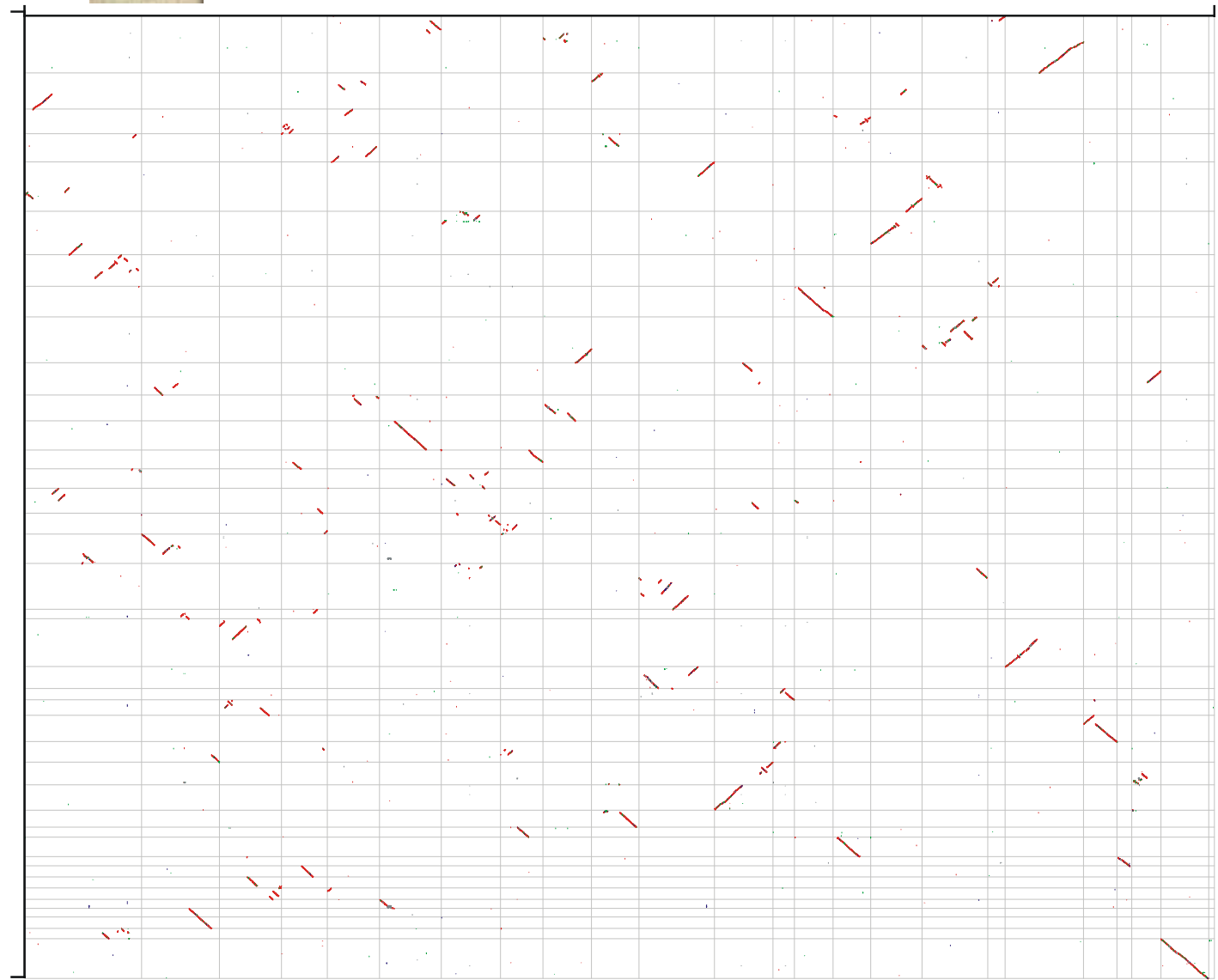
# $K_s$ Trees: predict...

- Lineage-specific paralogues
- Orthologues
- Retrotransposed pseudogenes
- Genes absent from assembly
- Conserved Synteny



# Ortholog gene positions in the human genome

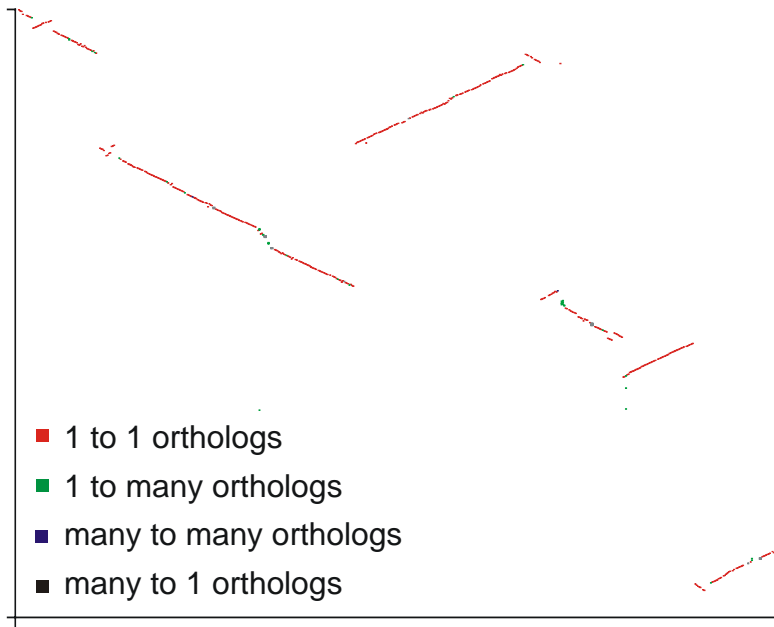
Ortholog gene positions in the dog genome



13,326: 73.2%



Ortholog gene positions in HSA17

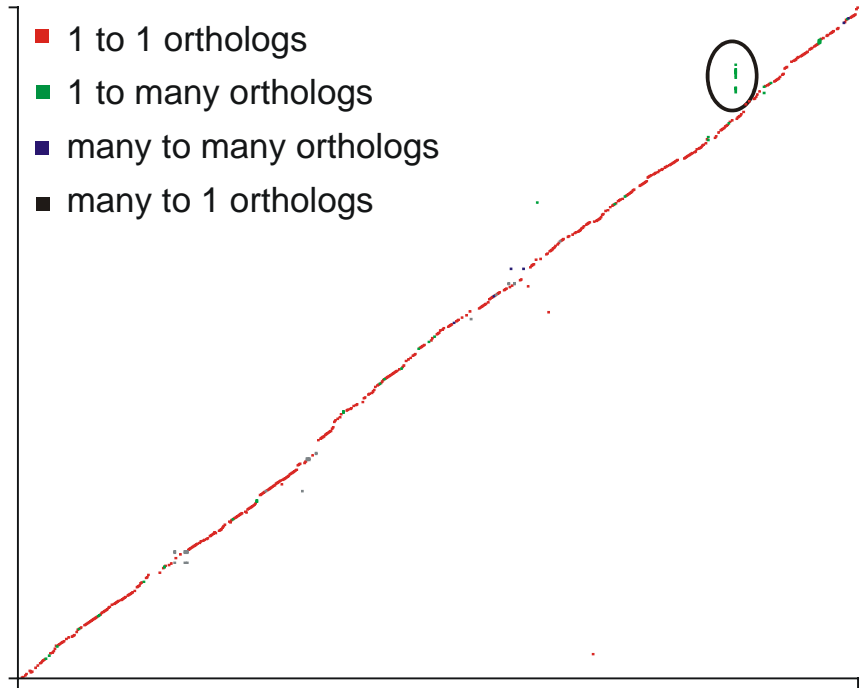


- 1 to 1 orthologs
- 1 to many orthologs
- many to many orthologs
- many to 1 orthologs

Ortholog gene positions in CFA9



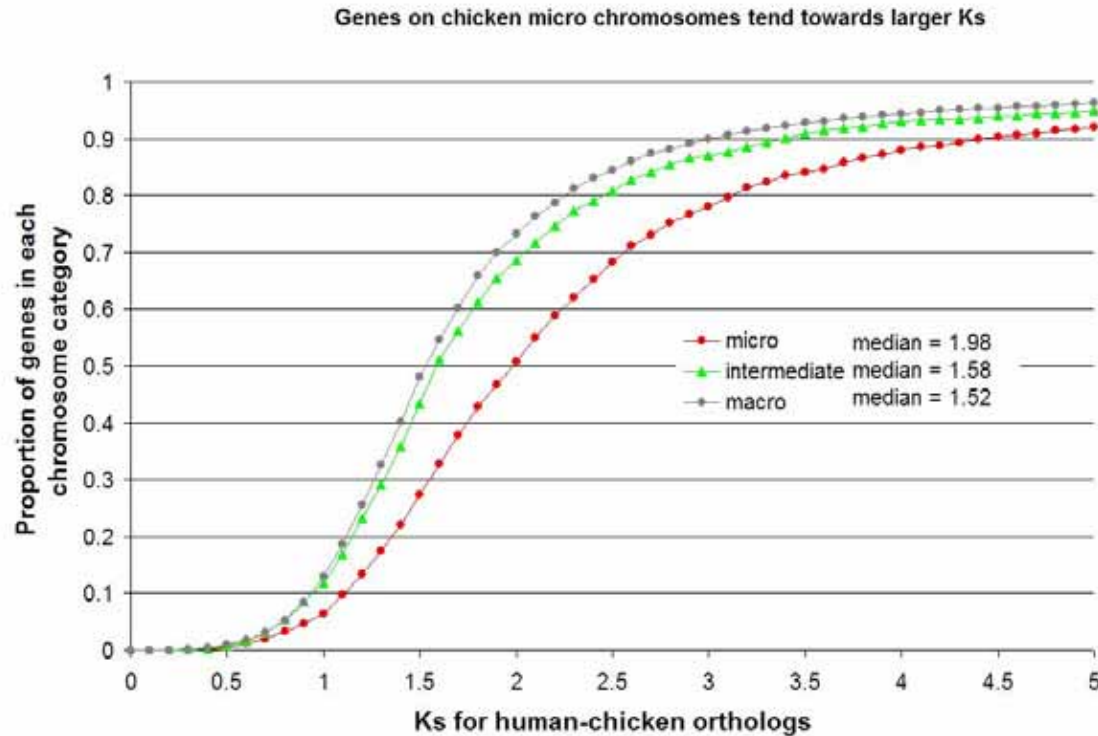
Ortholog gene positions in HSAX



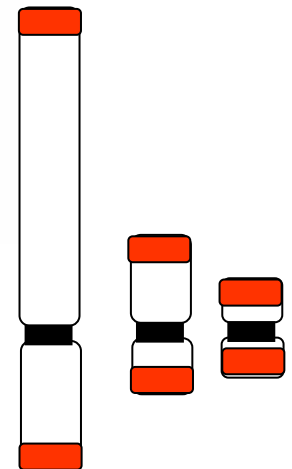
- 1 to 1 orthologs
- 1 to many orthologs
- many to many orthologs
- many to 1 orthologs

Ortholog gene positions in CFA9

# Variation in $K_S$



10Mb subtelomeric genes  
indistinguishable from microchromosomal  
genes, with respect to  $K_S$

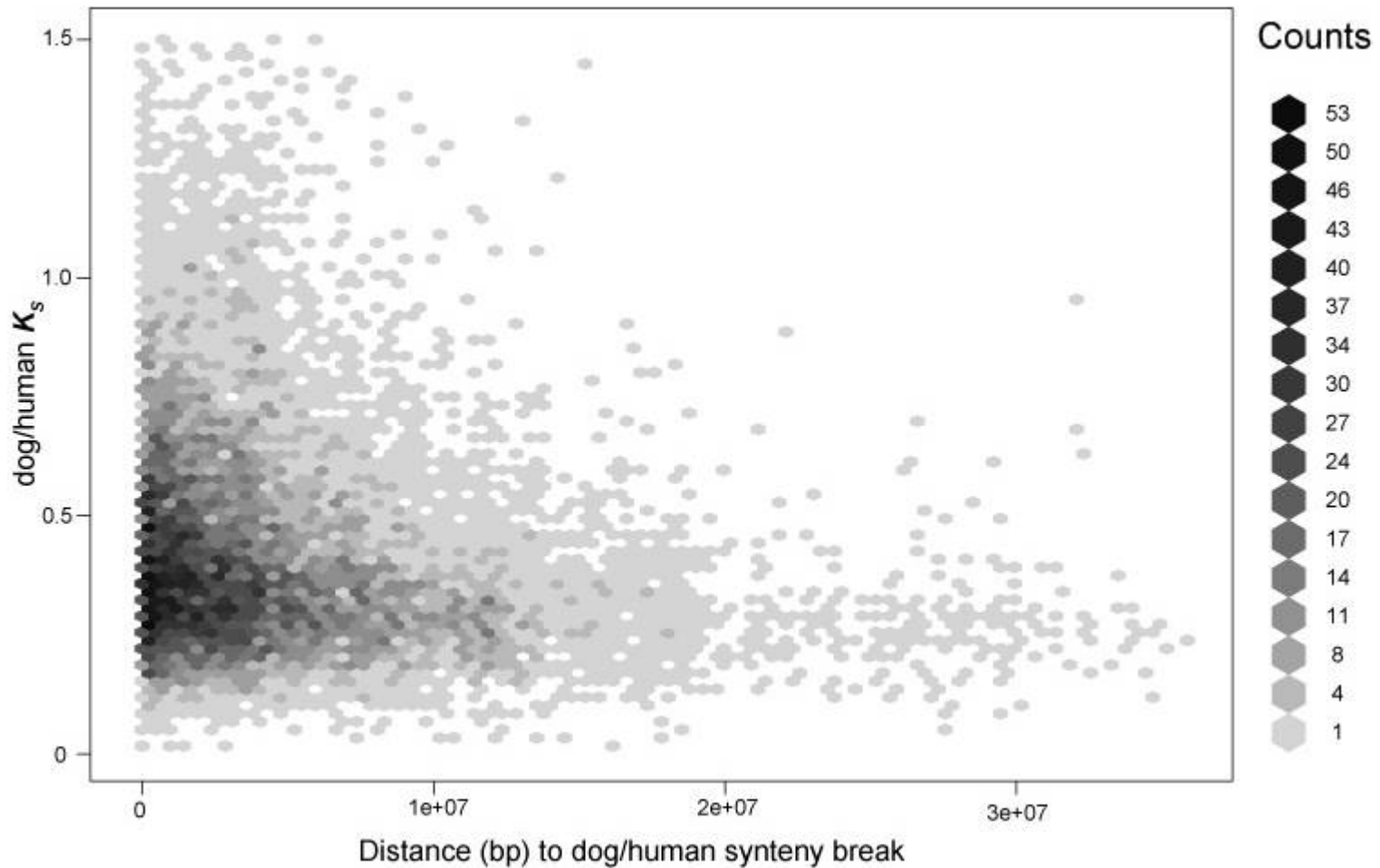


# $K_S$ variation

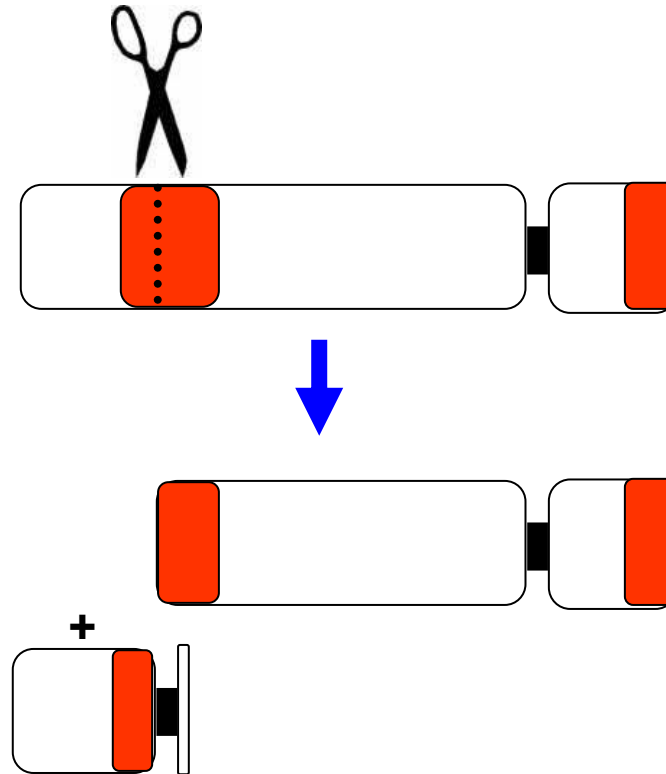
<b>Species</b>	<b>Subtelomeric / Subtelomeric</b>	<b>Subtelomeric / Interstitial</b>	<b>Interstitial / Subtelomeric</b>	<b>Interstitial / Interstitial</b>
Chicken / Human	3.15	1.95	2.40	1.53
Dog / Human	0.69	0.35	0.67	0.31
Dog / Chicken	2.15	2.08	1.98	1.64



# $K_S$ peaks towards synteny breaks

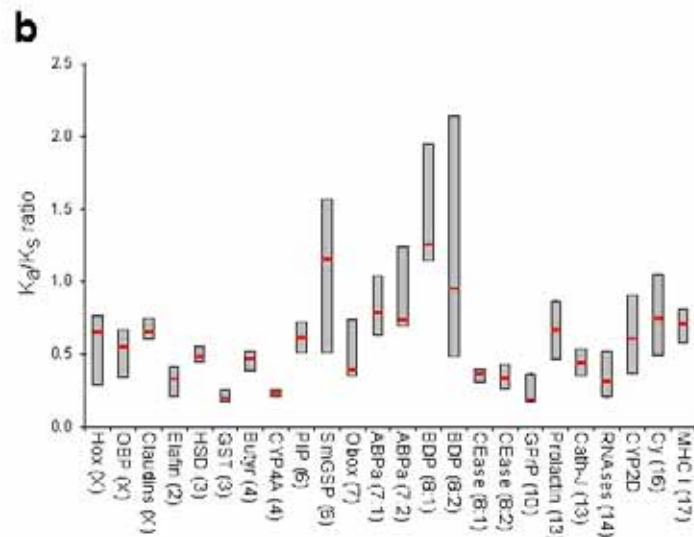
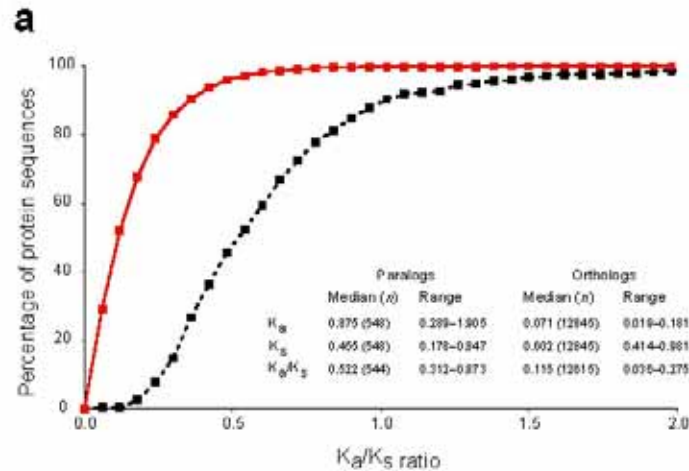


# Fragile breakage model

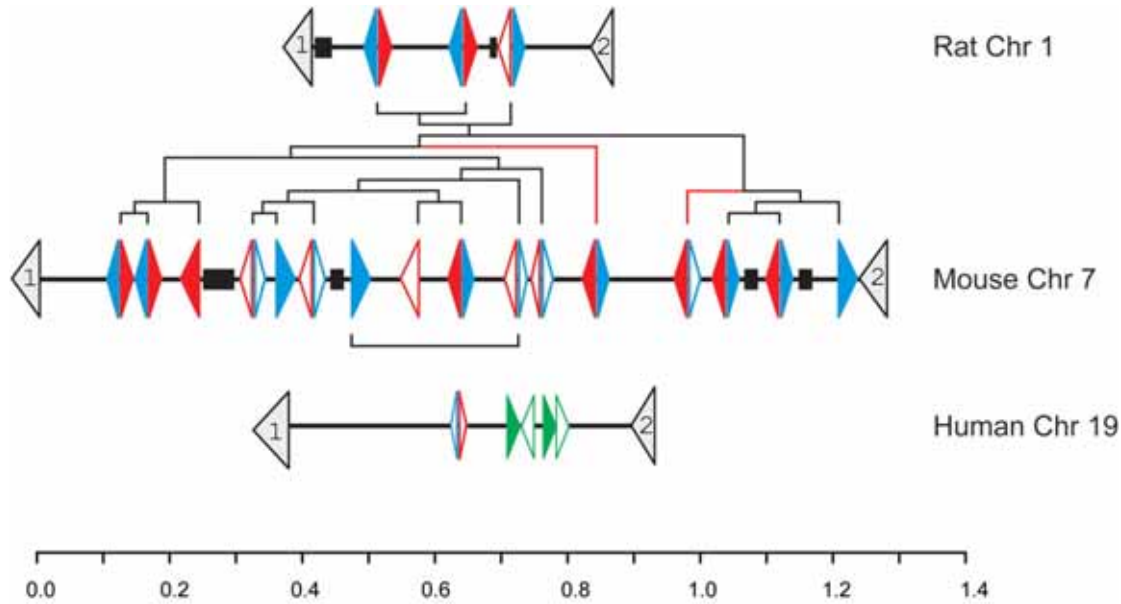


It still remains unclear *WHY* the human and cattle karyotypes are relatively stable, whereas those of dogs, mice and rats are not.

# Rapid Evolution in Gene Clusters

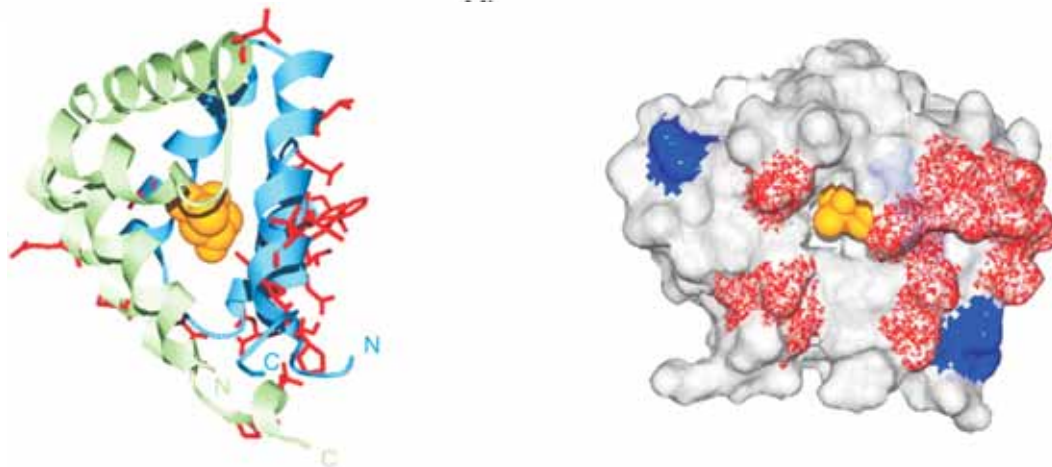


# Gene Duplication leading to Genome remodelling



Cattle: 7  $\alpha$  genes, and 11  $\beta$  genes (3 of which are pseudogenes)

*Directed Finishing?*



# To Do List?

Chris Ponting  
Leo Goodstadt  
Caleb Webber  
Andreas Heger

- Genes unique to cattle
- Complex orthology & paralogy relationships.
- QC on genes and assembly (pseudogenes and absent genes)
- Gene Reprediction
- Evolution of the bovine karyotype
- Low & high evolutionary rates:  $K_A/K_S$  &  $K_S$ .
- Intronic rates
- Adaptive evolution within genes

# Issues

- Gene Build: 6X coverage or more?
- Timing: ASAP
- Lineage-specific gene & lineage-specific biology: interactions with community.

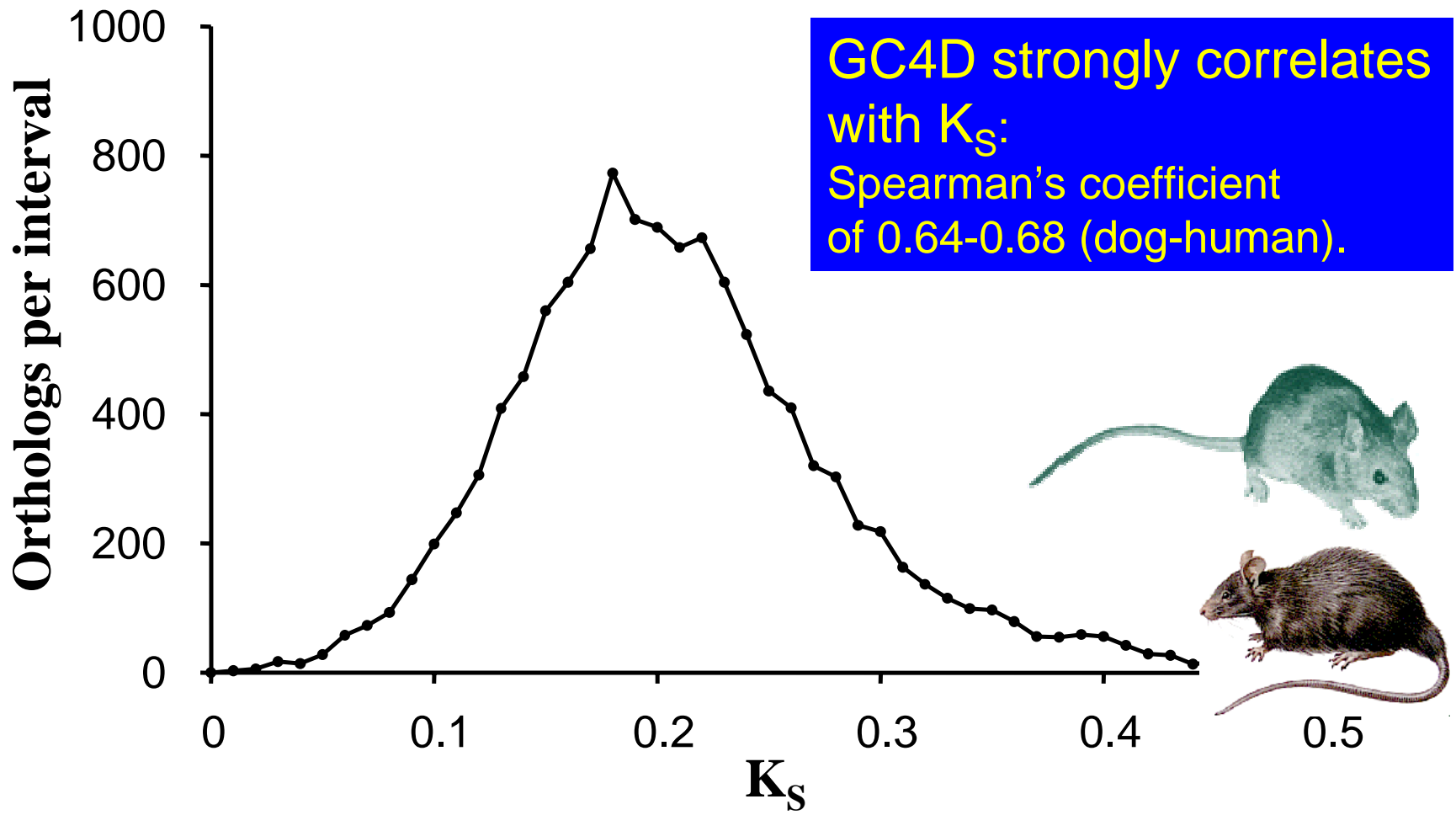




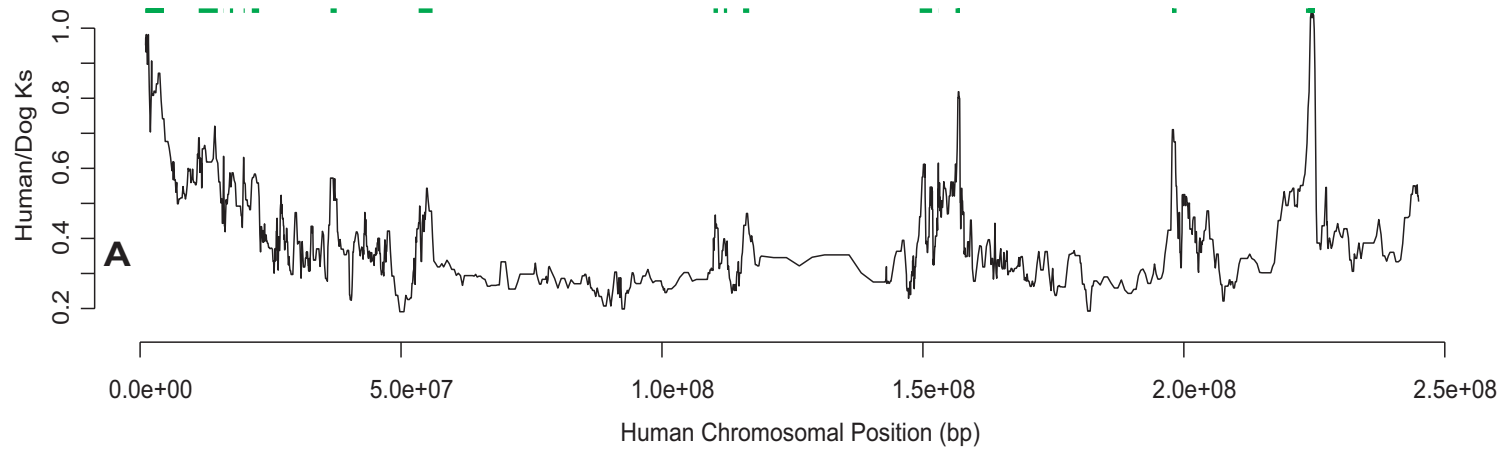
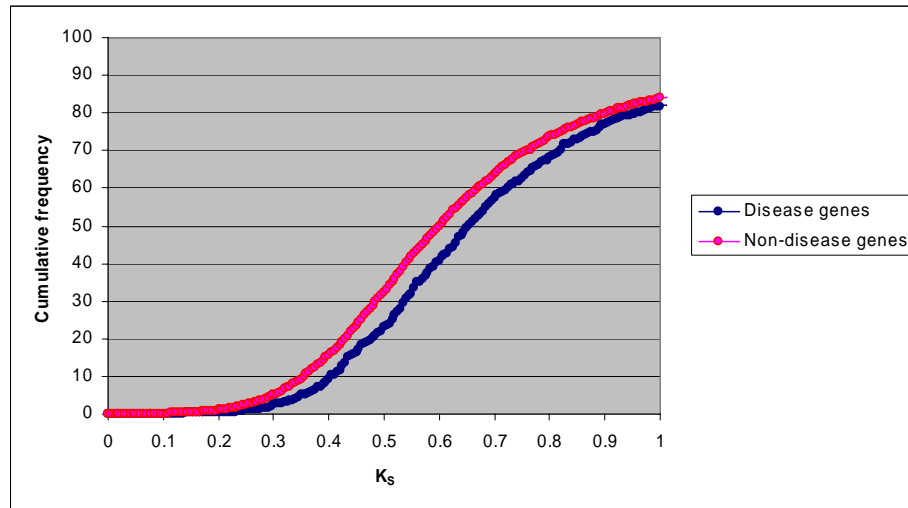




# Variations in $K_S$ values



# $K_S$ variation



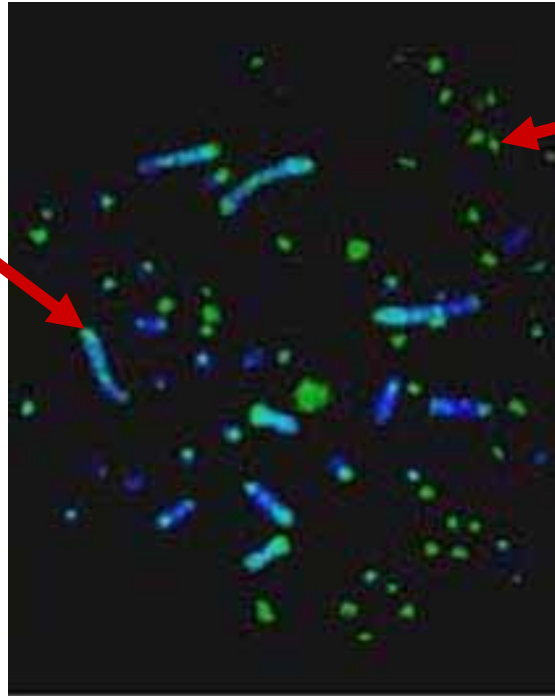
# Chicken:

## Large and Small Chromosomes

Large (macro-) chromosomes:

more DNA but lower gene density;

lower mutation rate; lower G+C



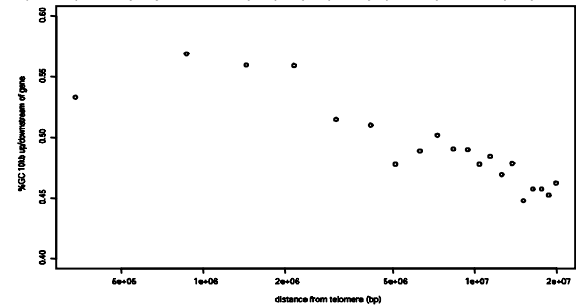
Small (micro-) chromosomes:

25% of the DNA but half the genes;

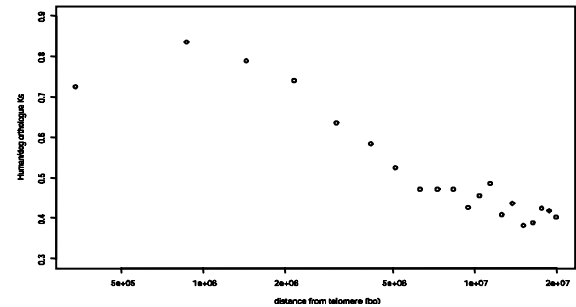
higher mutation rate; higher G+C

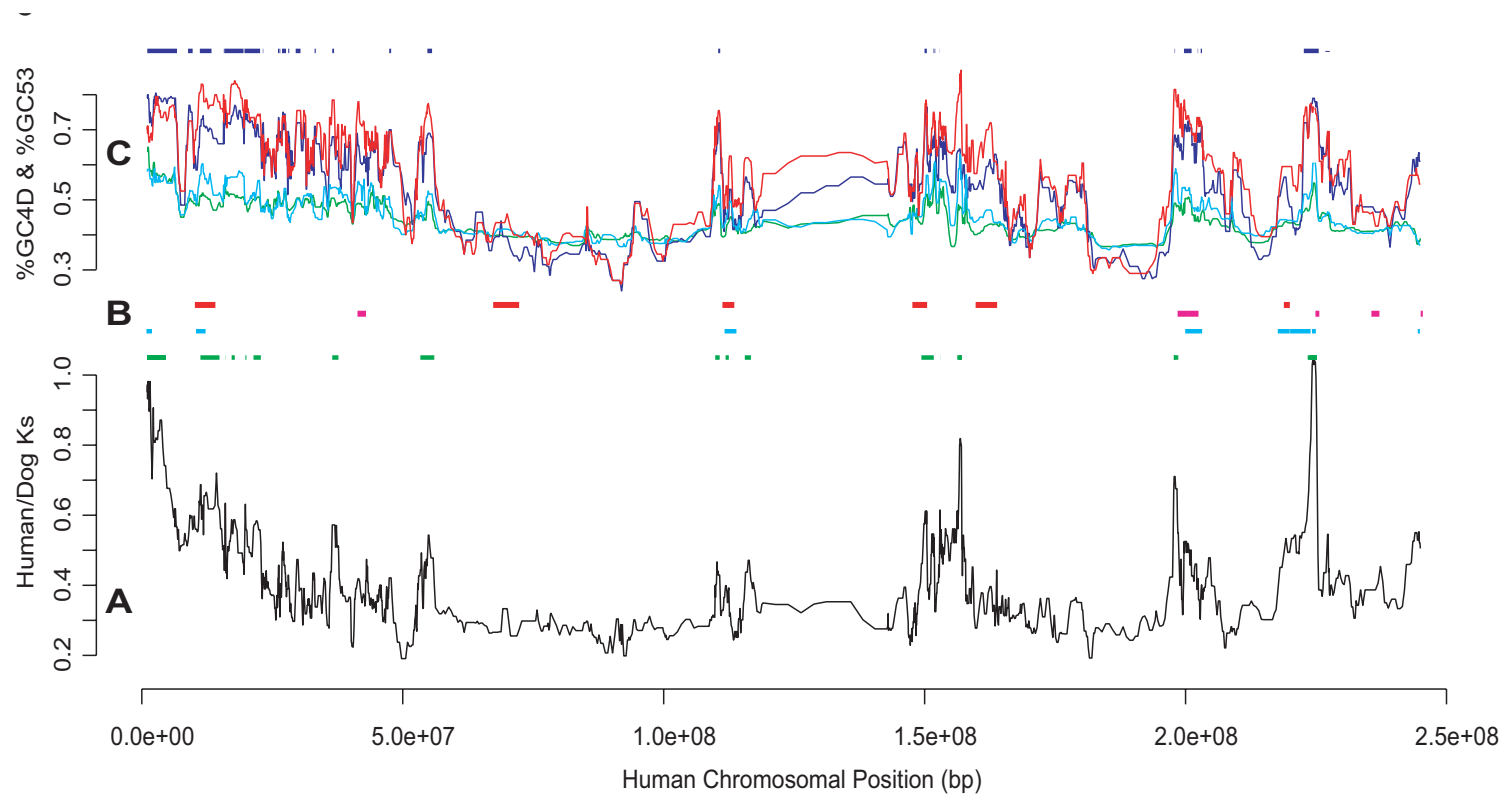
but dog telomeres  
are recent (derived)...?

Human GC4D vs distance from telomere



Hs-Cf Ks vs distance from human telomere





Breakpoints are not random:  
they occur preferentially in  
high G+C / high  $K_S$  regions

- $K_S$  hot spots that occur in ancestral interstitial regions have preferentially relocated to dog subtelomeres ( $p$ -value  $< 10^{-4}$ ).
- $K_S$  is significantly elevated towards dog-human synteny breakpoints ( $p$ -value  $< 2.2 \times 10^{-16}$ )

# Pathogenic Breakage

Occurs most frequently in:

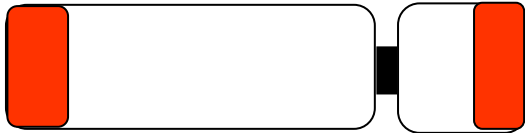
- subtelomeric regions
- high G+C regions.

(Yu et al., 1978; Stoll, 1980; Aula and von Koskull, 1976;  
Nakagome and Chiyo, 1976; Abeyasinghe et al., 2003)

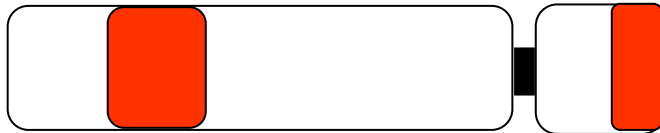




# Two models



**FAST** increase in G+C and in neutral rates in the last 60MY, since the dog karyotype shattered (DERIVED)



stable G+C and neutral rates in the last 60MY, since the dog karyotype shattered (ANCESTRAL)



Correlation coefficients (Spearman's  $\rho$ ) between 1:1 orthologues' GC4D fractions.

	Human	Dog	Mouse	Chicken
Human	-			
Dog	<b>0.945</b>	-		
Mouse	<b>0.836</b>	<b>0.818</b>	-	
Chicken	<b>0.539</b>	<b>0.595</b>	<b>0.539</b>	-
Rat	<b>0.825</b>	<b>0.810</b>	<b>0.937</b>	<b>0.524</b>

# Discussion

- Status and plans for gene prediction for bovine at Ensembl Sanger/EBI
- NCBI plans for the bovine sequence
- Using comparative genomics on protein families
- what the international bovine community wants and needs
- potential scientific partners