



# Examining reference gene sets

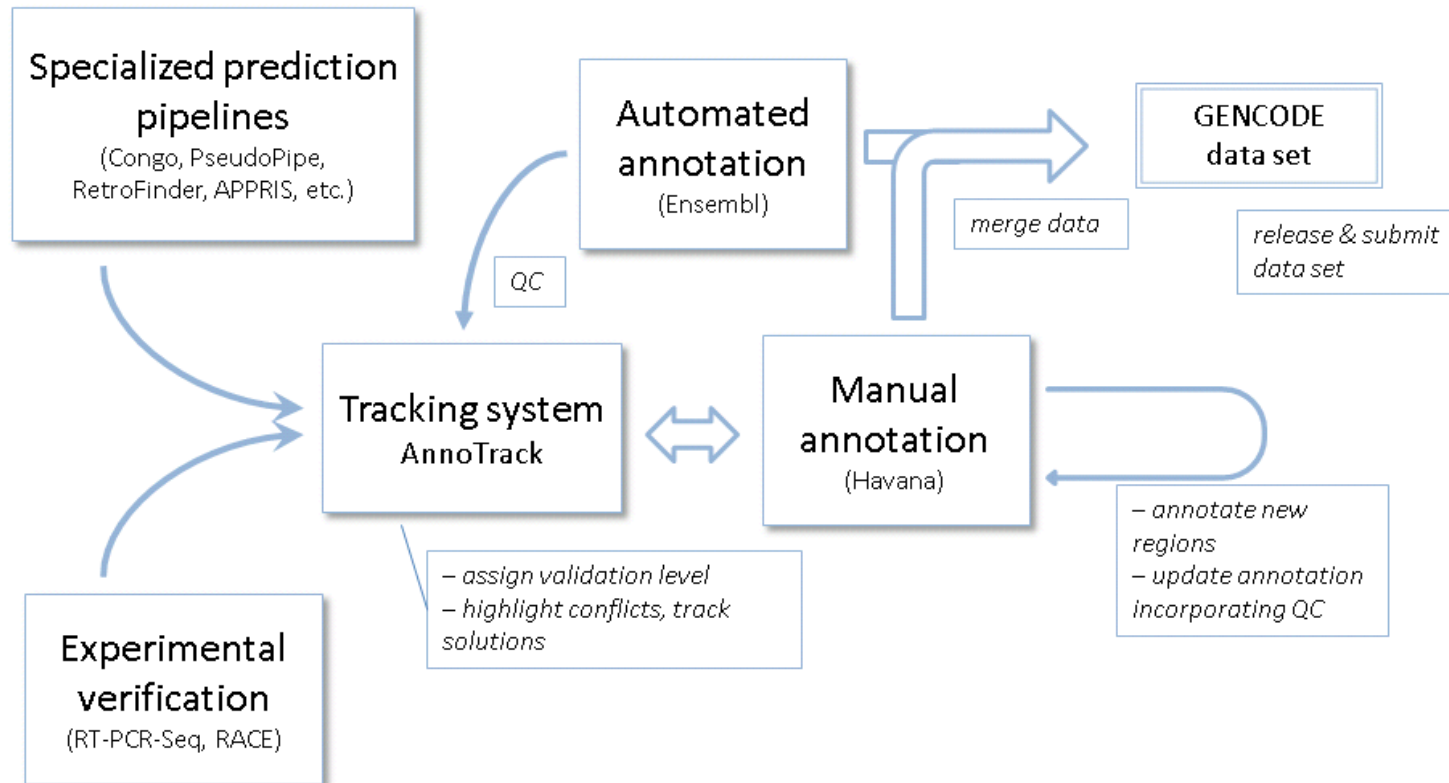
Jen Harrow  
WTSI



# What is GENCODE ?

- Human Reference gene set generated for the Encode Project
- Merge of automatic (Ensembl) and manual annotation (Havana) -release 4 May
- Experimental verification by RTPCR and RACE of transcripts(Lausanne)
- Manually annotate all loci biotypes (protein-coding, pseudogenes, “processed-transcripts/non coding RNAs”), small RNAs from Rfam and miRbase mapped to genome by ensembl
- Visualise in UCSC and Ensembl, collaborate with WashU, UCSC, Broad, Yale, CRG

# GENCODE pipeline



# gencodegenes.org



Project
Data
Participants
Publications
RGASP 1/2
RGASP 3
Contact us

## Statistics about the current Gencode freeze (version 11)

\*The statistics derive from the [gtf files](#), which include only the main chromosomes of the human reference genome.

### Version 11 (October 2011 freeze, GRCh37)

#### General stats

<b>Total No of Genes</b>	53639	<b>Total No of Transcripts</b>	180272
<b>Protein-coding genes</b>	20107	<b>Protein-coding transcripts</b>	81040
<b>Long non-coding RNA genes</b>	11600	- full length protein-coding:	60661
<b>Small non-coding RNA genes</b>	8801	- partial length protein-coding:	20379
<b>Pseudogenes</b>	12761	<b>Nonsense mediated decay transcripts</b>	10525
- processed pseudogenes:	9387	<b>Long non-coding RNA loci transcripts</b>	18566
- unprocessed pseudogenes:	2446		
- unitary pseudogenes:	156		
- polymorphic pseudogenes:	27		
- pseudogenes:	553		
<b>Immunoglobulin/T-cell receptor gene segments</b>			
- protein coding segments:	370		
- pseudogenes:	192		



# Automatic Annotation vs Manual

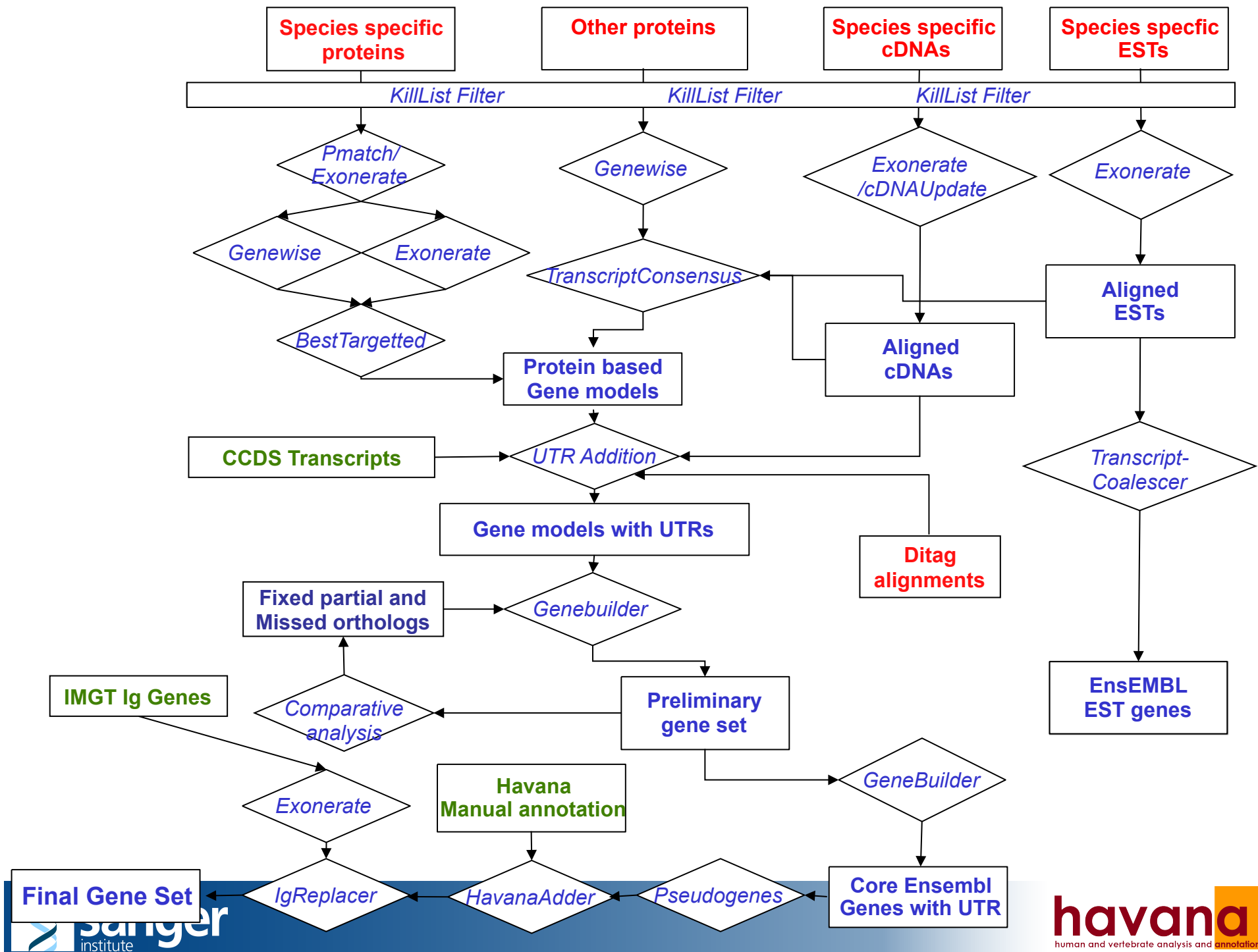


## Automatic Annotation

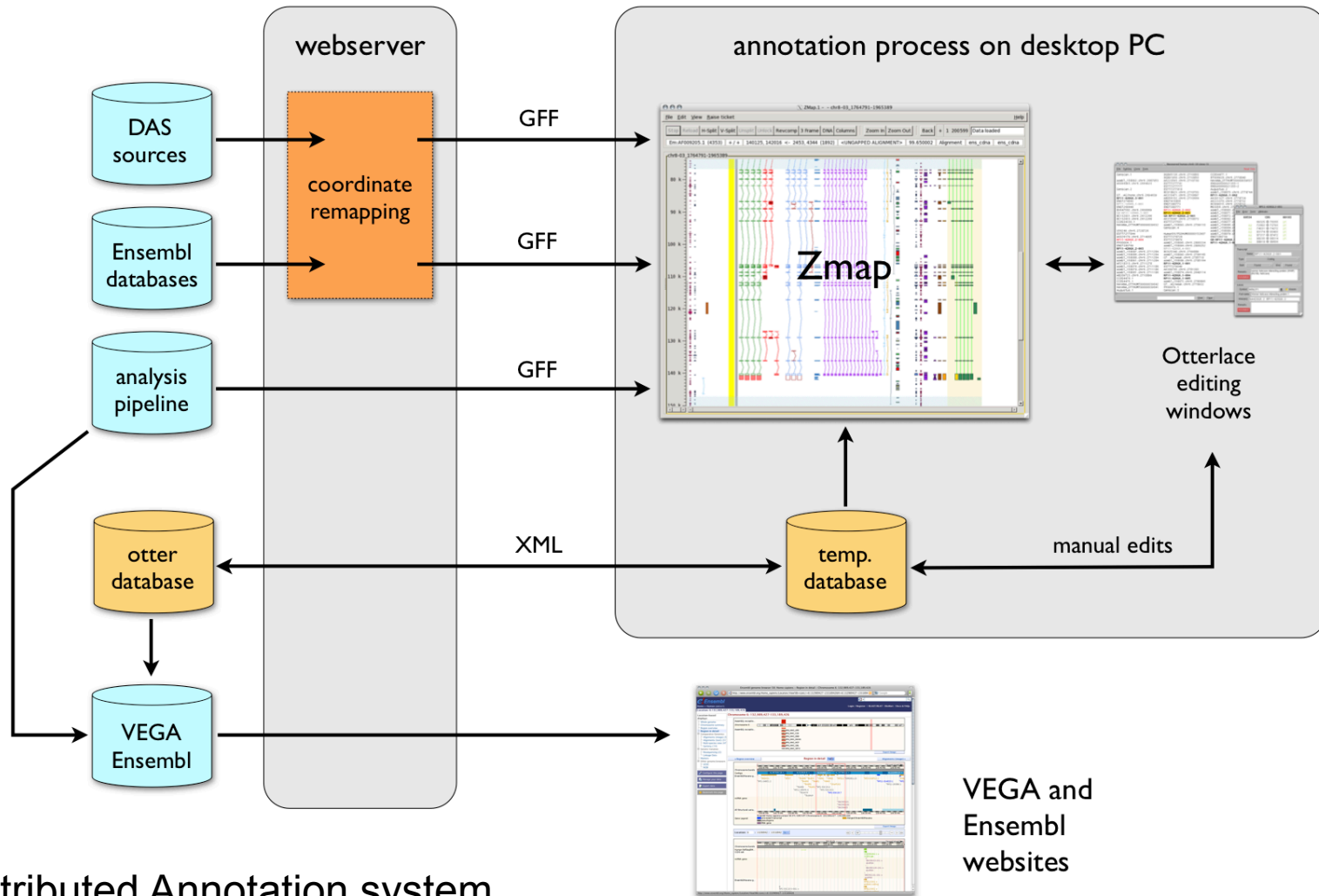
- Quick whole genome analysis ~ weeks
- Consistent annotation
- Use unfinished/illumina sequence/shotgun assembly
- No polyA sites/signals, pseudogene
- Predicts ~80% loci

## Manual Annotation

- Extremely slow~3 months Chr 6
- Need finished (high quality) seq
- Flexible, can deal with inconsistencies in data
- Most rules have exception
- Consult publications as well as databases



# Otterlace pipeline



DAS=Distributed Annotation system

# DAS source visible in Zmap

lace chr21-03, clones 486..490


File	SubSeq	Clone	Tools
AP001469.6-002			Launch ZMap Ctrl+Z
ESTT60480			Launch In A ZMap
augustus.2			Genomic Features Ctrl+G
ENST397708			Dotter Zmap hit Ctrl+.
AP001469.6-001			Exonerate Zmap hit/Column Ctrl+X
ENST291688			Rename locus Ctrl+Shift+L
ERI: AP001469.3-00			Re-authorize Ctrl+Shift+A
GD: AP001469.9-001			Load column data
MPI: AP001469.2-003	AP000471.3-010	ESTT60502	
genscan.5	AP000471.3-009	augustus.1	
CCDS13734.1	ENST310126	ESTT60498	
ESTT60605	CCDS13735.1	ESTT60499	
ESTT60603	ESTT60583	AP000337.2-005	
ESTT60600	ESTT60585	ESTT60495	
ESTT60599	ESTT60588	ESTT60496	
ESTT60607	AP000471.3-008	augustus.5	
AP001469.6-006	OTTHUMT00000207282	AP000337.2-004	
PF03399.1	AP000471.3-003	genscan.1	
	AP000471.3-004	AP000337.2-002	
ESTT60474	OTTHUMT00000207283	ESTT60491	
	ERI: AP000471.59-001	CCDS33592.1	
AP001469.6-008	AP000471.3-002	MPI: AP000337.1-001	
AP001469.6-007	AP000471.3-001	ERI: AP000337.1-001	
	MPI: AP000471.60-001	GD: AP000337.1-001	
PF02130.1	ESTT60587	OTTHUMT00000207336	
CCDS33591.1	AP000471.3-007	ENST337772	
ESTT60487	MPI: AP000471.60-003	AP000337.2-003	
OTTHUMT00000207272	OTTHUMT00000207286	AP000337.2-001	


Load column data


- augustus
- cpg
- das\_aspic
- das\_comparacons\_10way
- das\_congo\_exons
- das\_evigan
- das\_exonify
- das\_gerp\_23\_way\_constrelem
- das\_phastcons\_17way
- das\_phastcons\_28way
- das\_siepel\_novellocci
- das\_transmap\_mrna
- das\_transmap\_refseq
- das\_transmap\_splicedest
- das\_transmap\_ucscgenes
- das\_ucsc\_retroali3
- das\_washu\_human\_pasa\_ests
- das\_washu\_mrnas
- das\_washu\_nscan1
- das\_yale\_pseudogene
- ditag\_chip\_pet
- ditag\_gis\_pet
- ditag\_gis\_pet\_encode
- ens\_ccds\_from\_ensembl
- ens\_ensembl
- ens\_ensembl\_from\_ensembl\_havana
- ens\_ensembl\_havana
- ens\_estgenes
- ens\_ncrna
- ens\_separate\_ccds
- eponine
- est2genome\_human
- est2genome\_mouse
- est2genome\_other
- genscan
- halfwise
- refseq\_human
- repeatmasker
- trf
- uniprot\_sw
- uniprot\_tr
- vertrna

Load Cancel

# Manual annotation browser: Vega








 BLAST/BLAT | [Help & Documentation](#) [Login](#) · [Register](#)

 **wellcome trust sanger institute** The Vertebrate Genome Annotation (VEGA) database is a central repository for high quality manual annotation of vertebrate finished genome sequence. Human, mouse and zebrafish are in the process of being completely annotated, whereas for other species the annotation is only of specific genomic regions of particular biological interest. The majority of the annotation is from the [HAVANA](#) group at the [Wellcome Trust Sanger Institute](#)

 **havana**  
human and vertebrate analysis and annotation

The website is built upon code from the [Ensembl](#) project.

### Browse a genome

 <b>Human</b> [19-05-2011] <small>Ensembl</small>	 <b>Wallaby</b> [30-03-2009] <small>Ensembl</small>
 <b>Mouse</b> [01-02-2011] <small>Ensembl</small>	 <b>Pig</b> [16-05-2007] <small>Ensembl</small>
 <b>Zebrafish</b> [19-05-2011] <small>Ensembl</small>	 <b>Dog</b> [14-02-2005] <small>Ensembl</small>
 <b>Gorilla</b> [30-03-2009] <small>Ensembl</small>	

### What's New in Release 43 (19 May 2011)

- [Update to human annotation](#) (Human)
- [Update to zebrafish annotation](#) (Zebrafish)
- [Vega search updates](#) (all species)
- [Website enhancements](#) (all species)
- [Schema change](#) (all species)

[More news...](#)

### What's New in Release 42 (23 March 2011)

- [Update to human annotation](#) (Human)

[More news...](#)

### What's New in Release 41 (1 February 2011)

- [Update to human annotation](#) (Human)

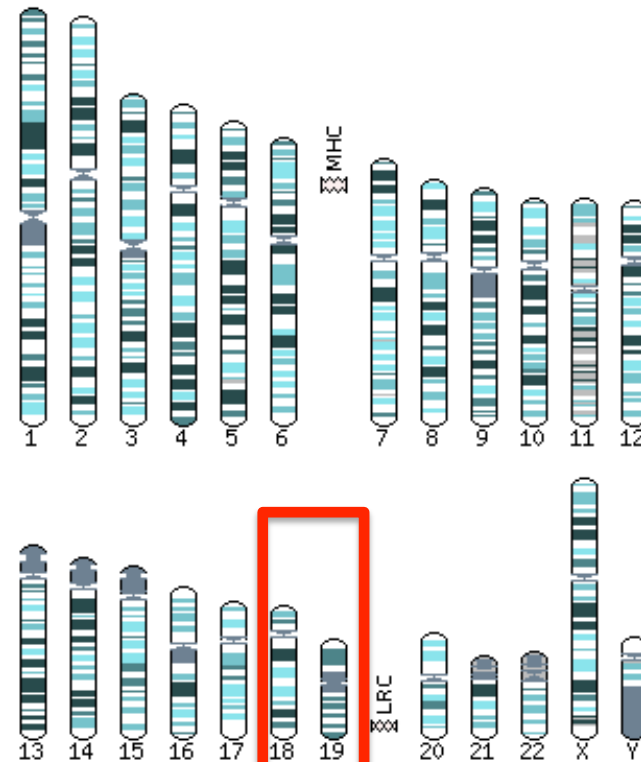
# Human annotation update

## Statistics

Last update	3 April 2012
Annotated <a href="#">Vega genes</a> :	46,369
Protein coding	19,557
Processed transcripts	12,165
Pseudogenes	13,362
IG & TR Genes	631
Annotated transcripts	166,571
Annotated exons	1,017,199
Total Bases	3,274,812,713
Golden Path Length	3,096,295,592

## Annotation progress

Click on a chromosome to browse:

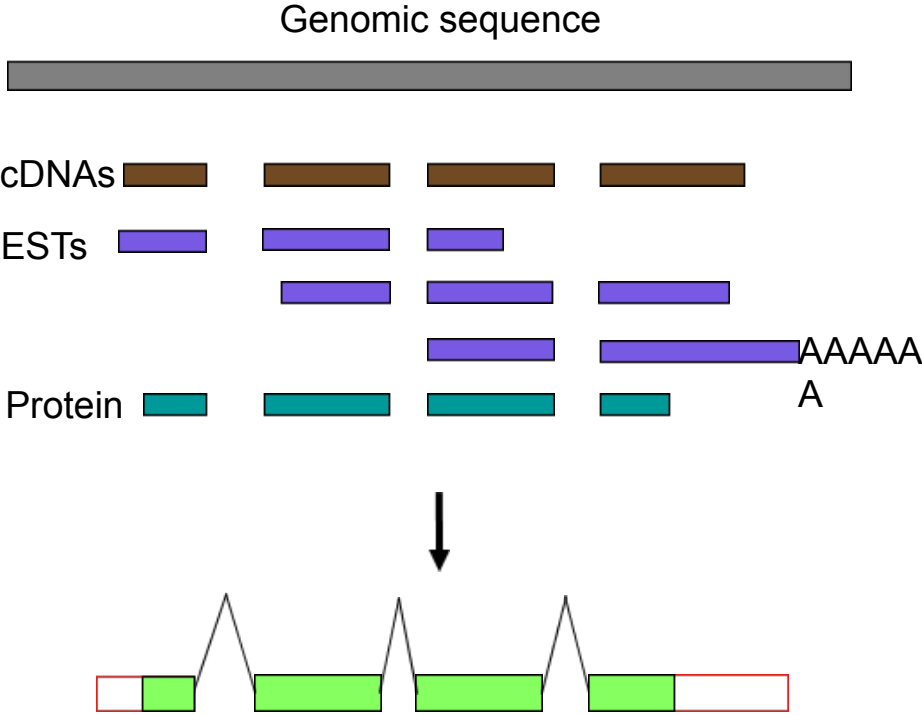


Shading indicates the centromere bands

To be annotated by Havana

# Manual Annotation and Biotypes

## Annotation based on transcriptional evidence.



### Protein Coding

- Known\_CDS
- Novel\_CDS
- Putative\_CDS
- Nonsense\_mediated\_decay

### Transcript

- Non\_coding
- Antisense
- Retained\_intron
- Putative
- Artefact

### Pseudogene

- Processed\_pseudogene
- Unprocessed\_pseudogene
- Transcribed\_processed
- Transcribed\_unprocessed
- Unitary\_pseudogene
- Polymorphic\_pseudogene

# Havana Biotypes now in Ensembl

Show/hide columns		Search: <input type="text"/>				
Name	Transcript ID	Length (bp)	Protein ID	Length (aa)	Biotype	CCDS
SHQ1-001	<a href="#">ENST00000325599</a>	2879	<a href="#">ENSP00000315182</a>	577	Protein coding	<a href="#">CCDS33788</a>
SHQ1-004	<a href="#">ENST00000463369</a>	2065	<a href="#">ENSP00000417452</a>	549	Protein coding	-
SHQ1-005	<a href="#">ENST00000482785</a>	502	<a href="#">ENSP00000418398</a>	111	Protein coding	-
SHQ1-003	<a href="#">ENST00000444040</a>	2844	<a href="#">ENSP00000402447</a>	50	Nonsense mediated decay	-
SHQ1-006	<a href="#">ENST00000471526</a>	553	<a href="#">ENSP00000417739</a>	63	Nonsense mediated decay	-
SHQ1-002	<a href="#">ENST00000468371</a>	4024	No protein product	-	Processed transcript	-
SHQ1-007	<a href="#">ENST00000468347</a>	134	No protein product	-	Processed transcript	-
SHQ1-008	<a href="#">ENST00000475558</a>	666	No protein product	-	Processed transcript	-

NB NMD variants coding unlike RefSeq

Status (known/novel/putative CDS/transcript) not taken from Havana currently gives some indication of confidence.

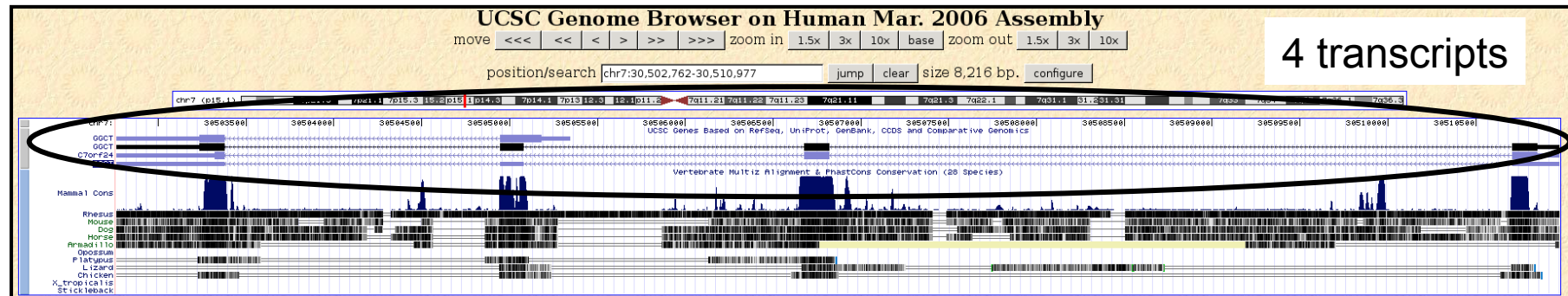
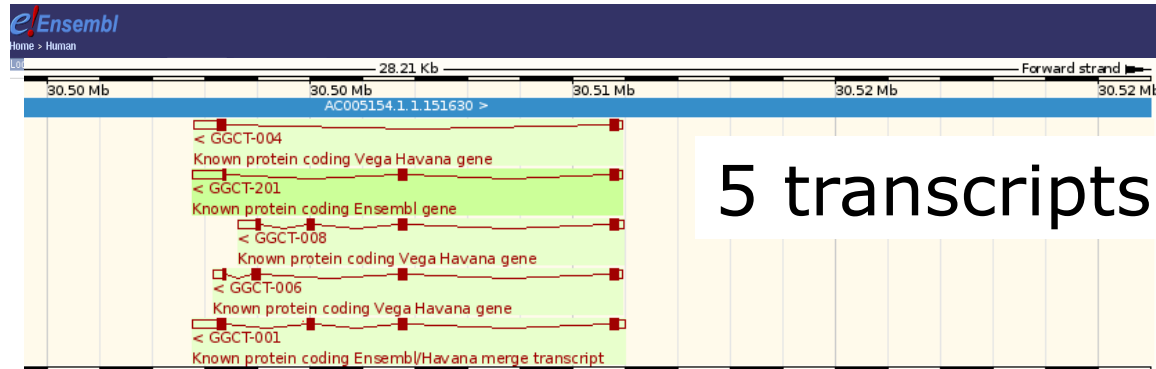
New confidence levels to be investigated



# Different gene sets ensembl/ UCSC/NCBI(Refseq)

	human	Mouse
RefSeq	20,669	23,090
Ensembl	21,297	21,111
UCSC	20,930	21,637

# GGCT in Ensembl/Refseq/UCSC



# What is CCDS

- Consensus coding sequence project
- UCSC, Ensembl, RefSeq and Havana
- Produce reference CDS set ATG-STOP on human and mouse genome must agree
- 1st rel 2005 13,142 genes 14 795 IDs
- Rel sept 2009 18,177 genes 23 739 IDs
- High quality but few alt splice variants and no UTRs, slow to increase

# CCDS website: GATA3 gene

## ATG->STOP

**Report for CCDS ID CCDS15674.1**

CCDS	Status	Species	Chrom.	Gene	NCBI Builds	Links
15674.1	Public	<i>Mus musculus</i>	2	Gata3	36.1 - 37.1	<a href="#">H</a> <a href="#">G</a> <a href="#">G</a>

**Sequence IDs included in CCDS 15674.1**

Original	Current	Source	Nucleotide ID	Protein ID	Status in CCDS	Seq. Status	Links
✓	✓	EBI,WTSI	ENSMUST00000102976	ENSMUSP00000100041	Accepted	alive	<a href="#">N</a> <a href="#">P</a> <a href="#">N</a> <a href="#">P</a>
✓	✓	EBI,WTSI	OTTMUST00000026063	OTTMUSP00000011932	Accepted	alive	<a href="#">N</a> <a href="#">P</a> <a href="#">N</a> <a href="#">P</a>
✓		NCBI	NM_008091.2	NP_032117.1	Updated	not alive	<a href="#">N</a> <a href="#">P</a> <a href="#">N</a> <a href="#">P</a>
	✓	NCBI	NM_008091.3	NP_032117.1	Accepted	alive	<a href="#">N</a> <a href="#">P</a> <a href="#">N</a> <a href="#">P</a> <a href="#">B</a>

**Chromosomal Locations for CCDS 15674.1**

On '-' strand of Chromosome 2 (NC\_000068.6)

Genome Browser links: [N](#) [U](#) [E](#) [V](#)

Chromosome	Start	Stop	Links
2	9779997	9780281	<a href="#">N</a> <a href="#">U</a> <a href="#">E</a> <a href="#">V</a>
2	9784722	9784847	<a href="#">N</a> <a href="#">U</a> <a href="#">E</a> <a href="#">V</a>
2	9790388	9790533	<a href="#">N</a> <a href="#">U</a> <a href="#">E</a> <a href="#">V</a>
2	9796016	9796552	<a href="#">N</a> <a href="#">U</a> <a href="#">E</a> <a href="#">V</a>
2	9798979	9799216	<a href="#">N</a> <a href="#">U</a> <a href="#">E</a> <a href="#">V</a>

### CCDS Sequence Data

Blue highlighting indicates alternate exons.

Red highlighting indicates amino acids encoded across a splice junction.

Mouse over the nucleotide or protein sequence below and click on the highlighted codon or residue to select the pair.

### Nucleotide Sequence (1332 nt):

```

ATGGAGGTGACTCGGGACCAGCCGGCTGGGGTGAGCCACCATCACCCCGGGTCTCAACGGTCAGCACC
CAGACACGCCACCACCCGGGCTCGGCCATTCTGACATGGAAGCTCAGTATCCGCTGACGGGAAGAGGTGGA
CGTACTTTTAAACATCGATGGTCAAGGCAACCACGTCCTTACTACGGAACCTCCGTCAGGGCTACG
GTGCAGAGGTATCCTCCGACCACACGGAGCCAGGTATGCCGCCCGCTCTGCTGCACGGATCTCTGC
CCTGGCTGGATGGCGGCAAGCCCTGAGCAGCCACCACCCGCTCGCCCTGGAACCTCAGCCCTTCTC
CAAGACGTCCATCCACCACGGCTCTCGGGGCTCTGTCCGTTTACCCTCCGGCTTCATCTCTCTCTG
CGGGCCGGCCACTCCAGTCTCATCTTTCACCTTCCCGCCACCCCGCCGAAAGACGTCTCCCCAGACC
CGTCGCTGTCCACCCCGGATCCCGGGTCCGGCAGGCAAGATGAGAAAGAGTGCCTCAAGTATCAGGT
GCAGCTGCCAGATAGCATGAAGCTGGAGACGTCTCACTCTCGAGGCAGCATGACCACCTGGTGGGGCC
TCATCTCAGCCACCACCCATACCCATCCGCTTATGTGCCGAGTACAGCTCTGGACTCTTCC
CACCAGCAGCTGTGGGAGGATCCCTACCGGTTCCGGATGTAAGTCGAGGCCAAGGCAGCATCCAG
CACAGAAGCAGGAGTGTGTAAGTCCGGGCAACCTTACCCACTGTGGCCGAGATGTTACCCGG
CACTACCTTTGCAATGCCTCGGACTCTACCATAAAATGAATGGGCAGAACCCGCCCTTATCAAGCCCA
AGCGAAGGCTGTCCGGCAGCAAGGAGCAGGACATCTCGCGGAACCTGCAGACCCACCCACCCCT
CTGGAGGAGGAACGCTAATGGGGACCCGCTGTGCAATGCCTGTGGGCTGTACTACAAGCTTCATAAT
AACAGACCCCTGACTATGAAGAAAGAAGGCATCCAGACCCGAAACCGGAAGATGTCTAGCAAAATCGAAA
AGTGCAAAAAGGTGCATGACCGCTGGAGGACTTCCCAAGAGCAGCTCCTCAACCCGGGCTCTCTC
CAGACACATGTATCCCTGAGCCACATCTCCTTCCAGCCTCCAGCCACATGTGACACACCCGAGC
CCCATGCATCCGCCCTCCGGCTCTCCTTCGGACCTACCACCTTCCAGCATGGTCACCCCATGGGTT
AG
    
```

### Translation (443 aa):

```

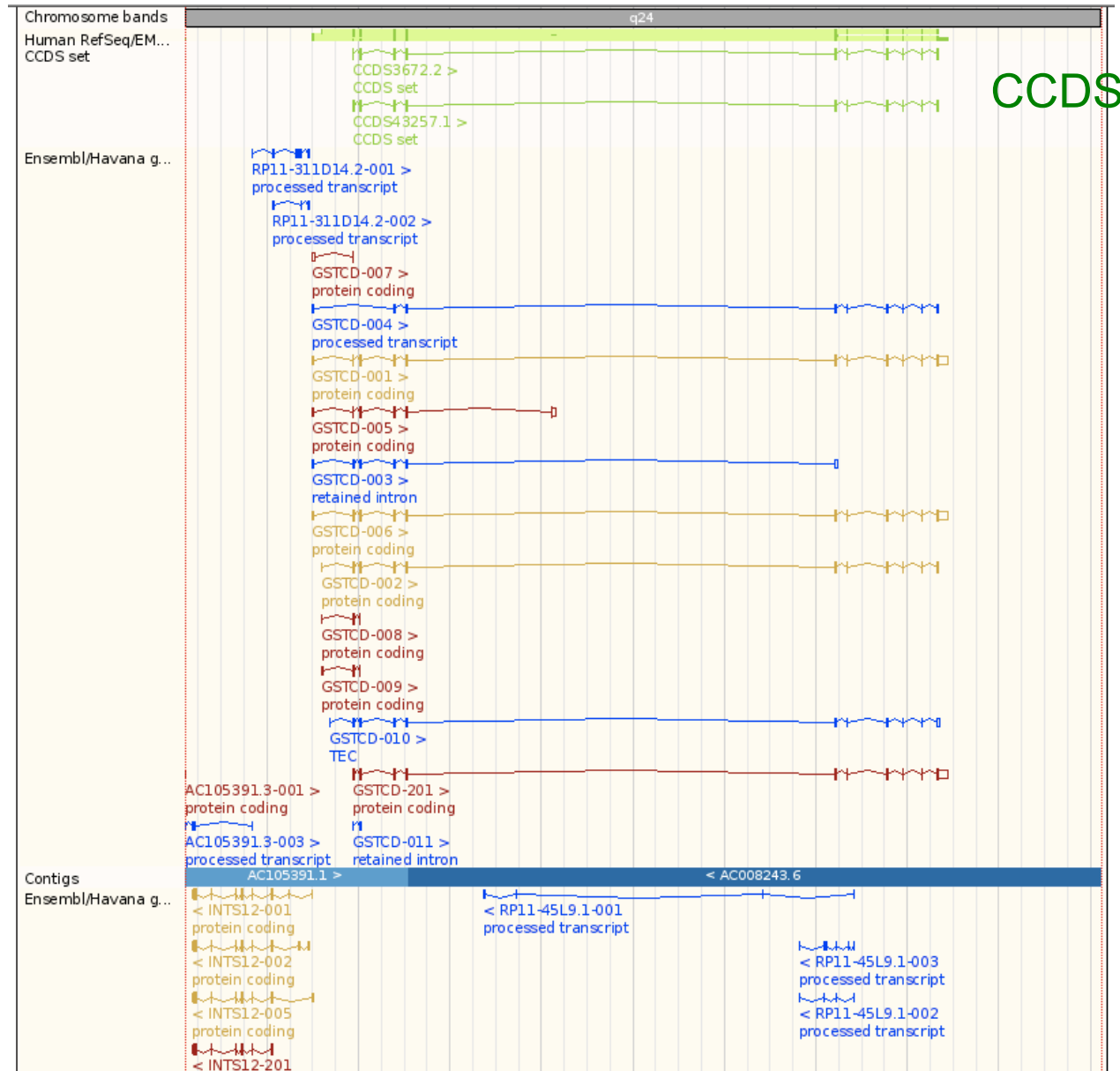
MEVTADQPRWVSHHHPAVLNGQHPDTHHPGLGHSYMEAQYPLTEEVDVLFNIDGQCQNHVPSYVNSVRAT
VQRYPTHHSQVCRPPLLHGSLPWLDDGKALSSHHTASPNLSPFKTSIHGSPGLSVYPASSSSL
AAGHSSPHLFTFPPTPKDVSPPSLSTPGSAGSARQDEKELKYQVQLPDSMKLETSHSRGSMTTLGG
SSSAHPITTPYVPEYSSGLFPSSLLGGSPTFGCKSRPKARSSREGRECVNCGATSTPLWRRDGTG
HYLCNACGLYHKMNGQNRPLIKPKRRLSAARRAGTSCANCQTFTTFLWRRNANGDPVCNACGLYKLNLI
NRPLTMKKEGIQTRNRKMSKSKKCKKVHDALEDFPKSSSFNPAALSRHMSLSLHISPFSSHMLTPT
PMHPPSGLSFGPHHPSMVTAMG
    
```

# Ensembl view

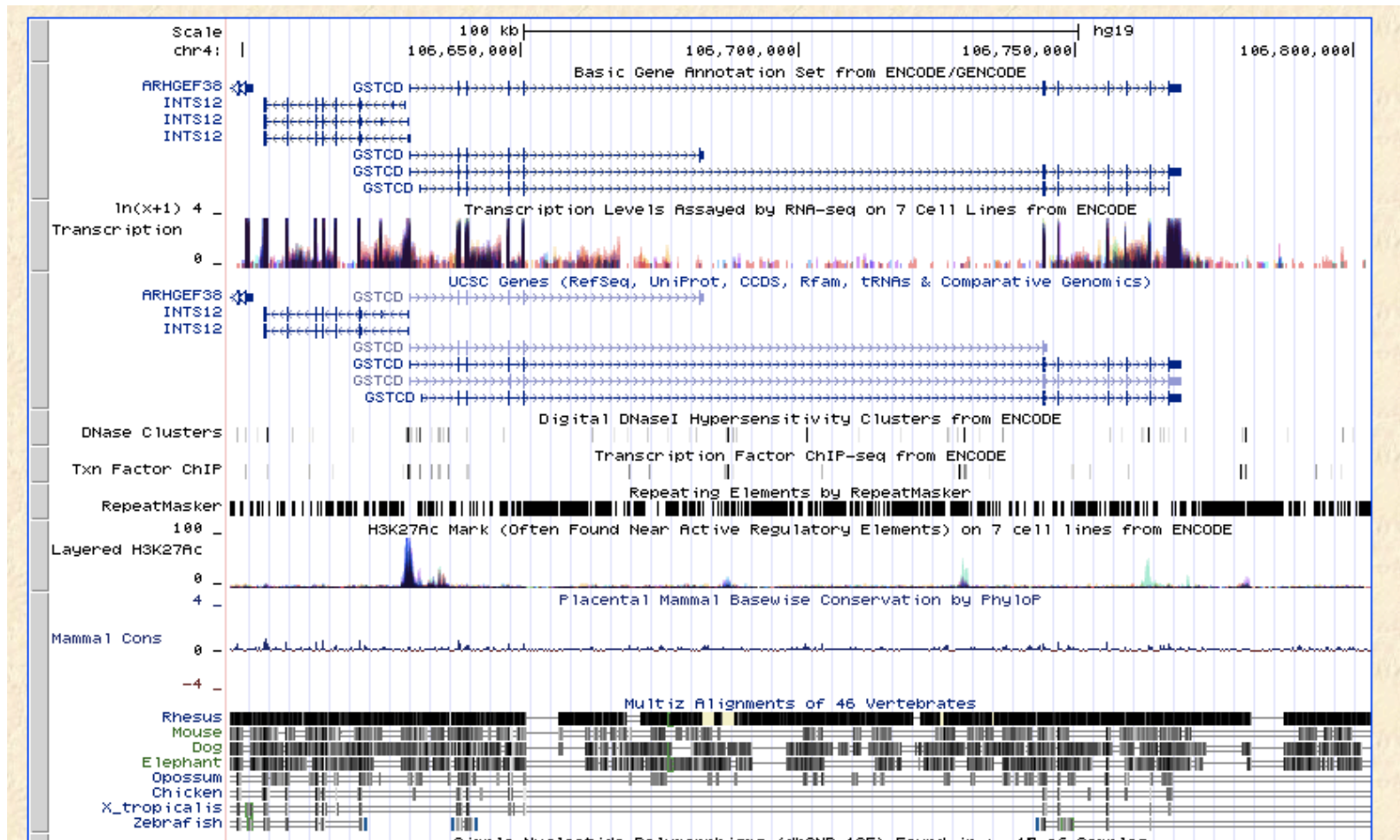
Gold: agreed  
ensembl/havana

Red : coding (001  
havana, 201  
ensembl)


Blue : non-coding



# UCSC View of GENCODE genes



# Changing default settings

 **Gene Annotations from ENCODE/GENCODE Version 11** ([All Genes and Gene Prediction Tracks](#))

Maximum display mode:    [Reset to defaults](#)

Select view ([help](#)):

Select all subtracks

List subtracks:  only selected/visible  all

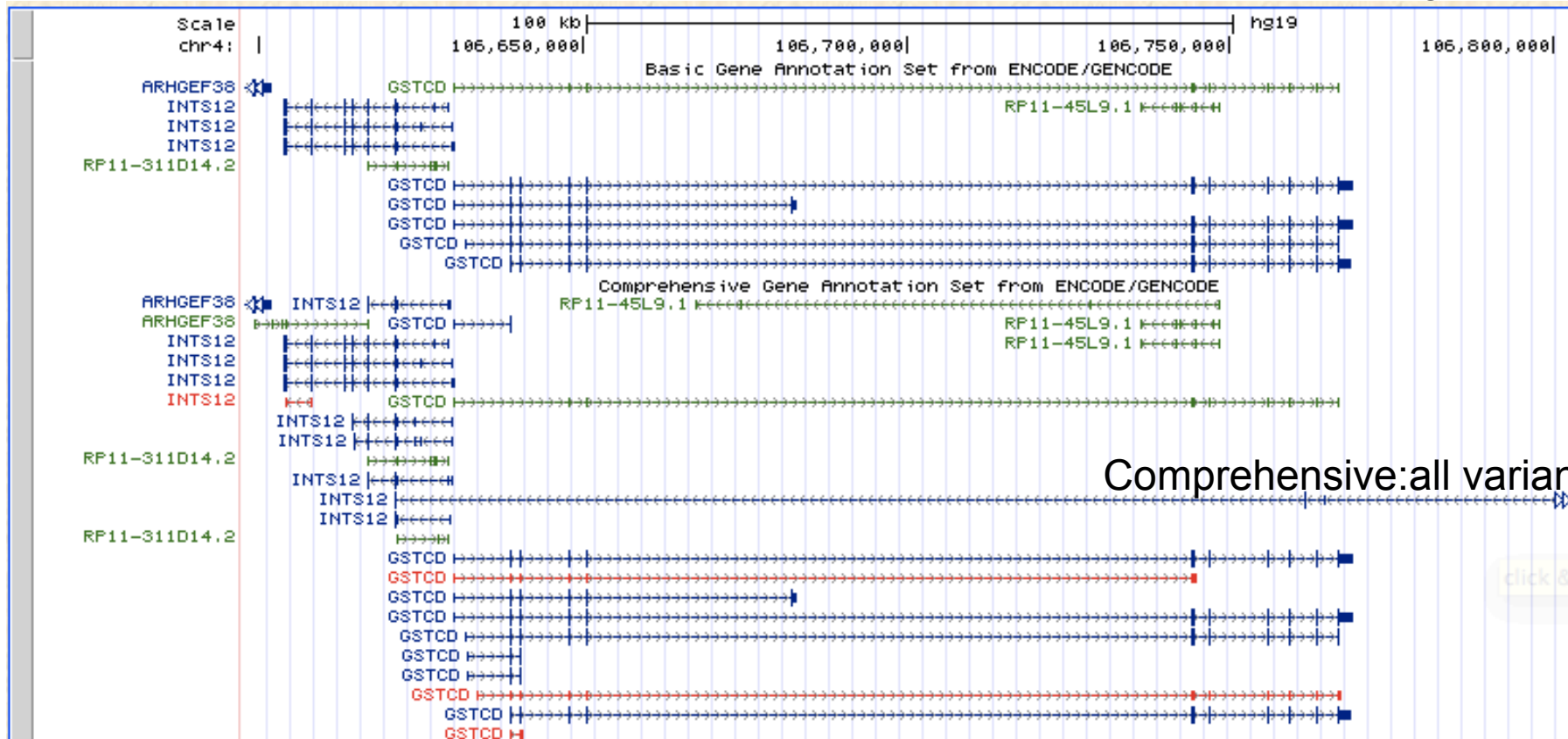
	Name <sup>1</sup>	View <sup>2</sup>	Track Name <sup>3</sup>							
<input checked="" type="checkbox"/>	<input type="text" value="full"/> <input type="button" value="Basic"/>	Genes	Basic Gene Annotation Set from ENCODE/GENCODE Version 11 ▾	<a href="#">schema</a>						
<div style="border: 2px solid black; padding: 5px;"><p>Label: <input checked="" type="radio"/> gene <input type="radio"/> accession <input type="radio"/> both <input type="radio"/> none</p><p>Color track by codons: <input type="text" value="genomic codons"/> <a href="#">Help on codon coloring</a></p><p>Show codon numbering: <input type="checkbox"/></p><p>Filter items by: (select multiple categories and items - <a href="#">help</a>)</p><table border="1"><thead><tr><th>Transcript Class</th><th>Transcript Annotation Method</th><th>Transcript Biotype</th></tr></thead><tbody><tr><td><input type="text" value="All"/></td><td><input type="text" value="All"/></td><td><input type="text" value="All"/></td></tr></tbody></table></div>					Transcript Class	Transcript Annotation Method	Transcript Biotype	<input type="text" value="All"/>	<input type="text" value="All"/>	<input type="text" value="All"/>
Transcript Class	Transcript Annotation Method	Transcript Biotype								
<input type="text" value="All"/>	<input type="text" value="All"/>	<input type="text" value="All"/>								
<input checked="" type="checkbox"/>	<input type="text" value="full"/> <input type="button" value="Comprehensive"/>	Genes	Comprehensive Gene Annotation Set from ENCODE/GENCODE Version 11 ▾	<a href="#">schema</a>						
<input type="checkbox"/>	<input type="text" value="full"/> <input type="button" value="Pseudogenes"/>	Genes	Pseudogene Annotation Set from ENCODE/GENCODE Version 11 ▾	<a href="#">schema</a>						
<input type="checkbox"/>	<input type="text" value="hide"/> <input type="button" value="2-way Pseudogenes"/>	2-way	2-way Pseudogene Annotation Set from ENCODE/GENCODE Version 11 ▾	<a href="#">schema</a>						
<input type="checkbox"/>	<input type="text" value="hide"/> <input type="button" value="PolyA"/>	PolyA	PolyA Transcript Annotation Set from ENCODE/GENCODE Version 11 ▾	<a href="#">schema</a>						

Filter by biotype



# Basic vs comprehensive GENCODE

BASIC:full length



Comprehensive:all variants

Manual and automatic coding non-coding pseudogene problem



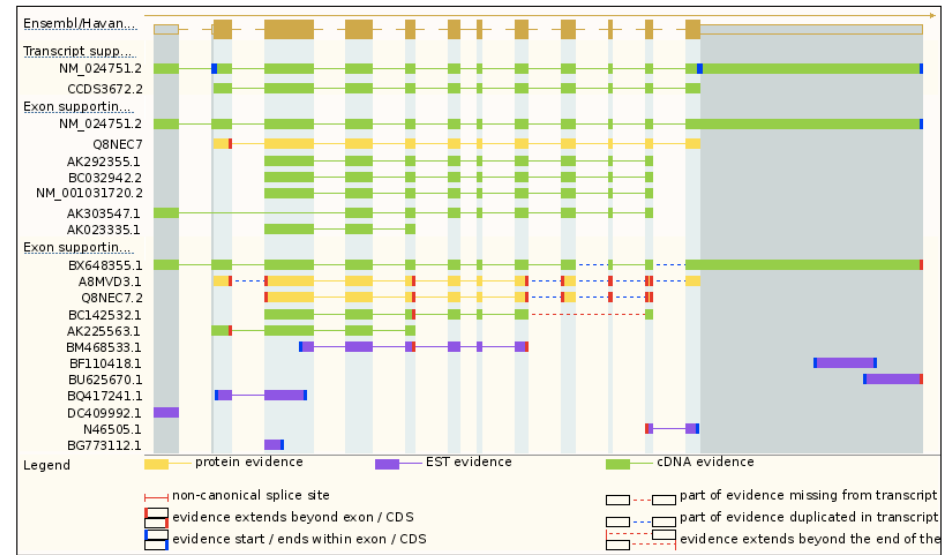
# Evidence for annotation

	Transcript	Gene
<b>Gencode id</b>	<a href="#">ENST00000394730.3</a>	<a href="#">ENSG00000138780.9</a>
<b>HAVANA manual id</b>	<a href="#">OTTHUMT00000253947.3</a>	<a href="#">OTTHUMG00000131211.4</a>
<b>Position</b>	<a href="#">chr4:106629941-106768882</a>	<a href="#">chr4:106629941-106768882</a>
<b>Strand</b>	+	
<b>Biotype</b>	protein_coding	protein_coding
<b>Status</b>	KNOWN	KNOWN
<b>Annotation Level</b>	manual (2)	
<b>Annotation Method</b>	manual & automatic	manual & automatic
<b>HUGO gene</b>	GSTCD	
<b>CCDS</b>	CCDS3672.2	

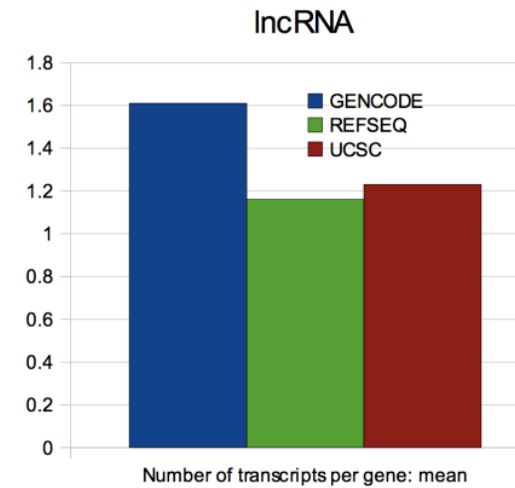
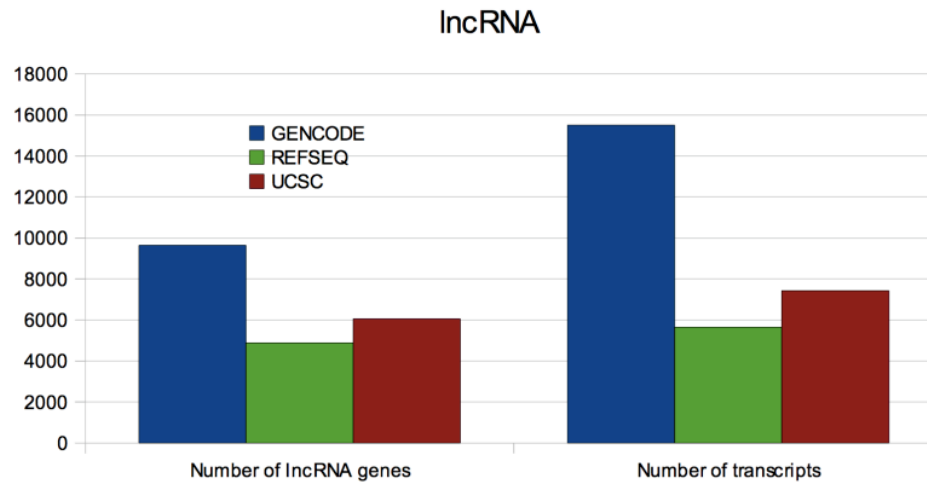
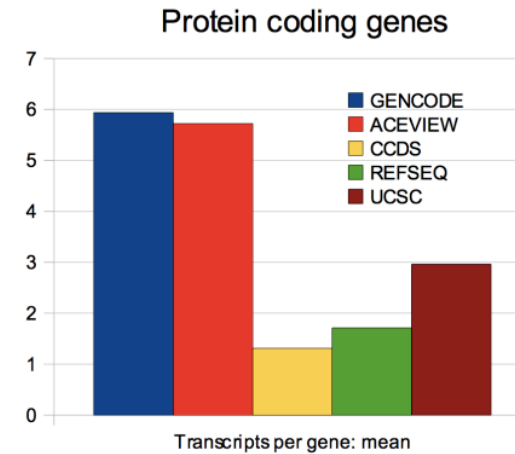
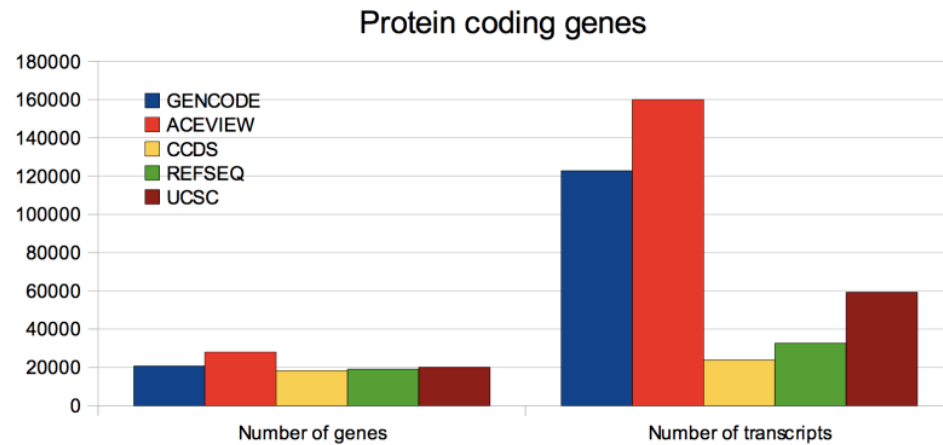
Show **All** entries      Show/hide columns      Filter

Name	Transcript ID	Length (bp)	Protein ID	Length (aa)	Biotype	CCDS
GSTCD-001	<a href="#">ENST00000394730</a>	4043	<a href="#">ENSP00000378218</a>	546	Protein coding	<a href="#">CCDS3672</a>
GSTCD-002	<a href="#">ENST00000360505</a>	2154	<a href="#">ENSP00000353695</a>	633	Protein coding	<a href="#">CCDS43257</a>
GSTCD-005	<a href="#">ENST00000507281</a>	2091	<a href="#">ENSP00000422858</a>	340	Protein coding	-
GSTCD-006	<a href="#">ENST00000515279</a>	4273	<a href="#">ENSP00000422354</a>	633	Protein coding	<a href="#">CCDS43257</a>
GSTCD-007	<a href="#">ENST00000512828</a>	572	<a href="#">ENSP00000423639</a>	40	Protein coding	-
GSTCD-008	<a href="#">ENST00000510865</a>	578	<a href="#">ENSP00000423792</a>	166	Protein coding	-
GSTCD-009	<a href="#">ENST00000509336</a>	718	<a href="#">ENSP00000423779</a>	167	Protein coding	-
GSTCD-201	<a href="#">ENST00000394728</a>	4049	<a href="#">ENSP00000378216</a>	633	Protein coding	<a href="#">CCDS43257</a>
GSTCD-010	<a href="#">ENST00000505640</a>	2425	No protein product	-	TEC	-
GSTCD-004	<a href="#">ENST00000515255</a>	1373	No protein product	-	Processed transcript	-
GSTCD-003	<a href="#">ENST00000484843</a>	2228	No protein product	-	Retained intron	-
GSTCD-011	<a href="#">ENST00000503409</a>	653	No protein product	-	Retained intron	-

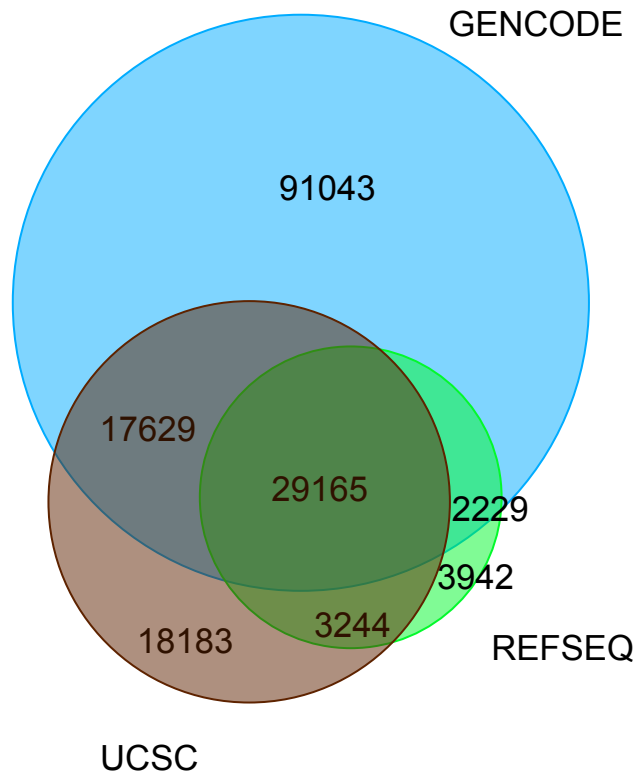
Supporting Evidence			
Source	Sequence	Source	Sequence
CCDS	CCDS3672.2	EMBL	AK023335.1
EMBL	AK225563.1	EMBL	AK292355.1
EMBL	AK303547.1	EMBL	BC032942.2
EMBL	BC142532.1	EMBL	BF110418.1
EMBL	BG773112.1	EMBL	BM468533.1
EMBL	BQ417241.1	EMBL	BU625670.1
EMBL	BX648355.1	EMBL	DC409992.1
EMBL	N46505.1	RefSeq_dna	NM_001031720.2
RefSeq_dna	NM_024751.2	Uniprot/SPTREMBL	A8MVD3.1
Uniprot/SWISSPROT	Q8NEC7	Uniprot/SWISSPROT	Q8NEC7.2



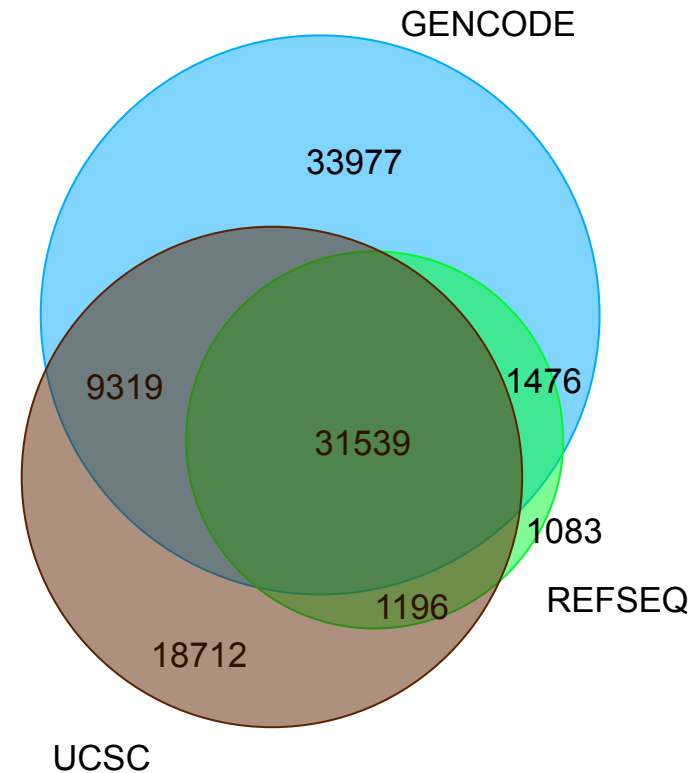
# Overview of GENCODE



# Overlap of exact match transcripts



**Transcript**



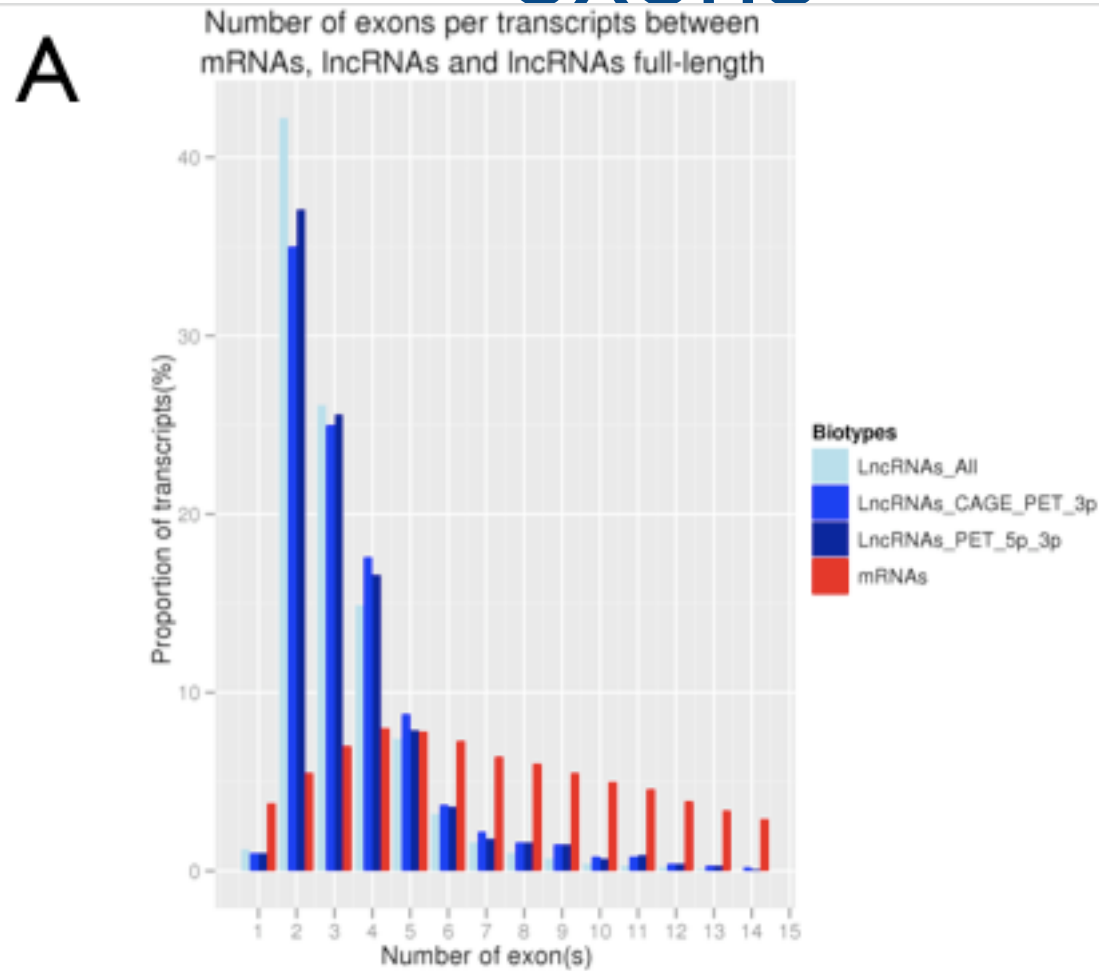
**CDS**

# Long noncoding RNA biotypes (HGNC and John Mattick)

- **lincRNA:**
  - Intergenic >200bp spliced (chromatin signatures observed but not mandatory)
- **Antisense:**
  - for transcripts overlapping any part of the genome within 5kb of the start of the CDS and 30 kb of the end of the CDS of a coding locus on the opposite strand.
- **3'\_overlapping\_ncRNA:**
  - for transcripts where ditag and/or published experimental data strongly supports the existence of short non-coding transcripts transcribed from the 3' UTR.
- **Sense\_overlapping:**
  - for transcripts that contains a coding gene in their introns on the same strand.
- **Sense\_intronic:**
  - for transcripts that are in introns of coding genes and do not overlap any exons.



# lncRNA shorter than coding exons



# Genome reference consortium (GRC)

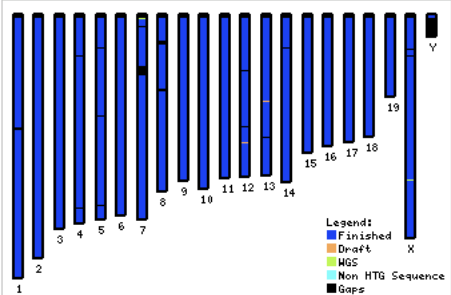
## Genome Reference Consortium

[GRC Home](#) | [Human](#) | [Mouse](#) | [Help](#) | [Report an issue](#) | [Contact Us](#) | [Curators Only](#)

[Overview](#) | [Issues under Review](#) | [Assembly Data](#) | [Report a problem](#)

### Mouse Genome Overview

Information concerning the continuing improvement of the mouse genome.



**Mouse Genome Build 37:** A graphical representation of the mouse genome in Build 37. The genome is colored with respect to the genomic component used to build the genome assembly at that location. While >95% of the genome is constructed using HTGS phase 3 (finished) sequence, small bits of unfinished and WGS contribute as well.

The most recent assembly for the mouse is Build 37, which was produced by the Mouse Genome Sequencing Consortium (MGSC). This assembly is based on DNA from a single inbred strain (C57BL/6J) and is largely composed of finished clone sequences. There are, however, 105 unspanned gaps and just over 700 spanned gaps remaining in the assembly. Most of the unspanned gaps result from the inclusion of Whole Genome Shotgun (WGS) contigs in the assembly. The inclusion of this sequence was critical in order to preserve genes not represented in the BAC tiling path. Much of the WGS sequence has been placed on the chromosome, but some remains either [unplaced](#) or [unlocalized](#).

The GRC is now actively working to close these gaps. Additionally libraries from C57BL/6J have become available and are being screened in an effort to identify clones that span these gaps. Many of the remaining gaps are associated segmental duplication as described in [She et al., 2008](#). If you know of any remaining issues, please contact us using the links above.

#### GRC News and Updates

**GRCh37 is now available in Map Viewer**  
Fri, 14 Aug 2009

NCBI has annotated and released the latest version of the public human genome assembly (GRCh37).

**GRCh37 now available at Ensembl**  
Fri, 14 Aug 2009

Ensembl has annotated and released the latest version of the public human genome assembly (GRCh37). This is available as part of Ensembl release 55.

[see all](#)

#### Recently Resolved Mouse Issues

**Mouse (MG-40)** Oct 8, 2009  
AC102264.26 is redundant. Will change the status of AC102264.26 to redundant on sequence curation page.

**Mouse (MG-3500)** Oct16, 2009  
Duplicate of MG-3558.

[see all](#)

#### References

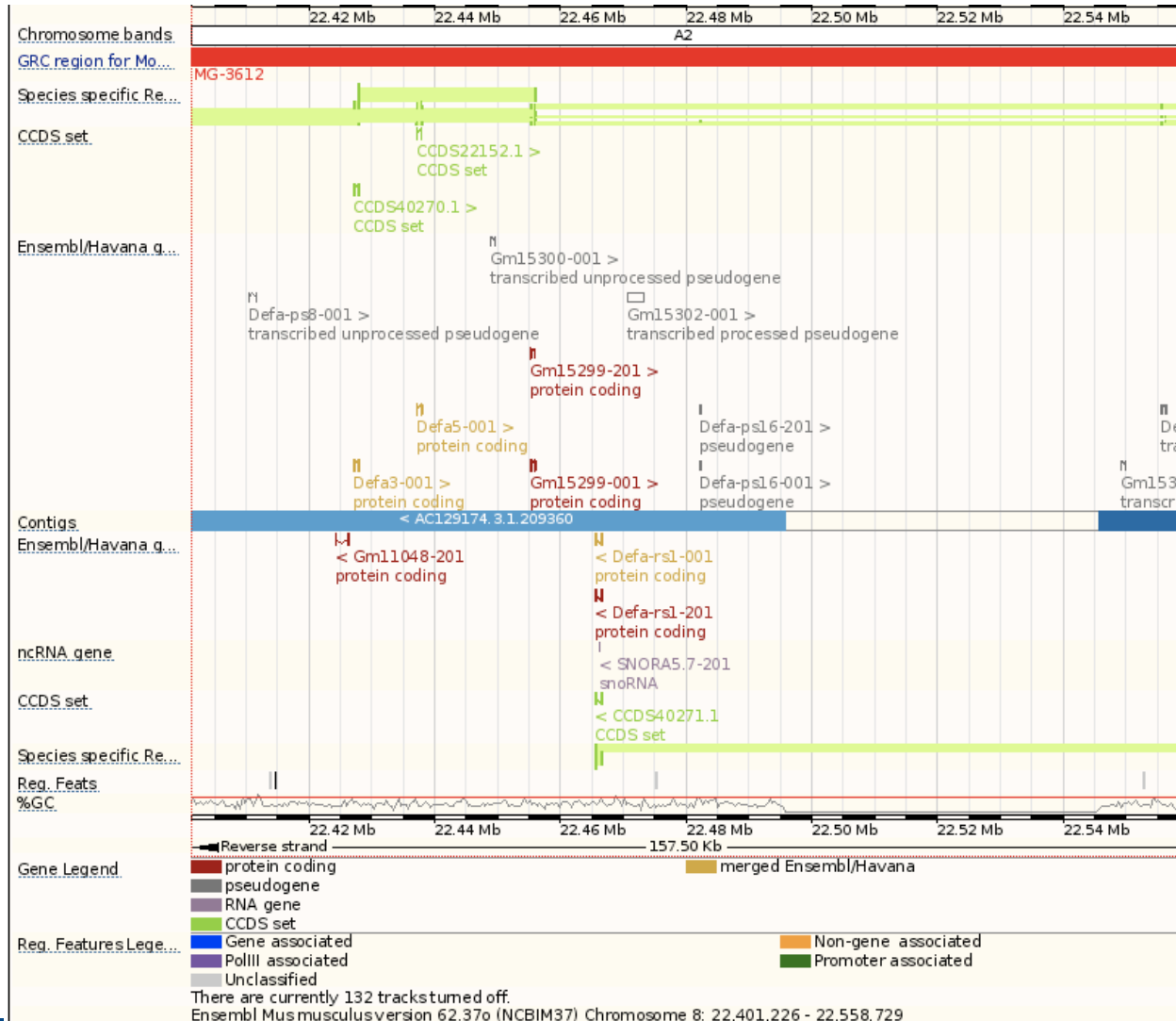
**Whole Genome Papers**

- [The Mouse Genome WGS Assembly](#)
- [The Mouse Genome: Clone based assembly](#)

FTP | NHGRI | The Wellcome Trust | HHS | NIH | Accessibility | Page last updated: Jun 4, 2009

Gaps and assembly issues highlighted which can Affect annotation

# GRC problem regions in Ensembl



GRC DAS track  
Has link to GRC  
database

Gap in assembly



# Information about assembly errors

<a href="#">GRC Home</a>	<a href="#">Data</a>	<a href="#">Help</a>	<a href="#">Report an Issue</a>	<a href="#">Contact Us</a>	<a href="#">Credits</a>	<a href="#">Curators Only</a>
<a href="#">Mouse Overview</a>	<a href="#">Mouse Issues Under Review</a>	<a href="#">Mouse Assembly Data</a>	<a href="#">Report a problem</a>			

## Issue Report for MG-3612

Category: Gap

Report type: TPF Analysis

Status: Resolved

Description: There is a type 2 reference assembly gap present between components AC129174.3 and AC133094.4

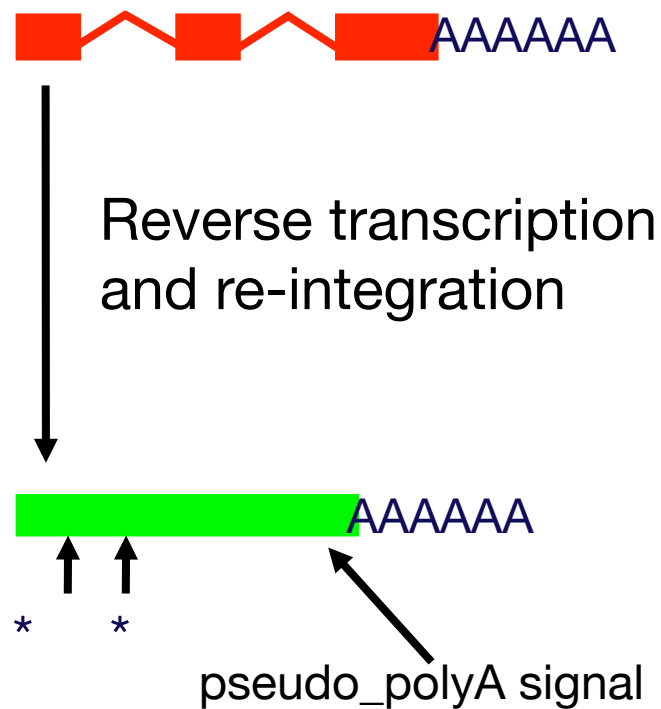
Resolution: RP23-16O24 AC240396.3 and CH36-223O13 AC239604.3 have been selected, sequenced, and submitted, and they close this gap.

### Assembly Information

MGSCv37 chr8: 22,336,657-22,729,648 (View Region: [Ensembl](#) | [NCBI](#) | [UCSC](#) )

# Exercises

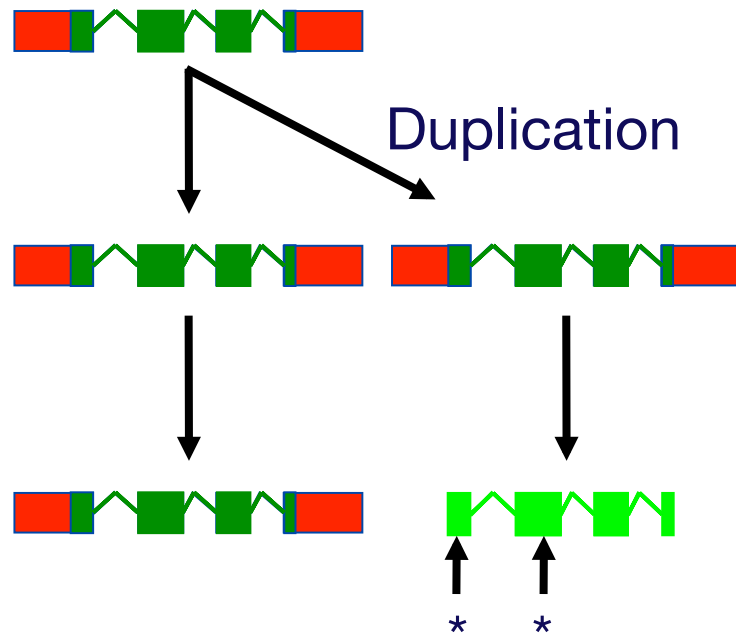
# Retrotransposed pseudogene



- mRNA transcript reverse transcribed back into DNA and inserted into chromosomal DNA
- Inserted randomly into genome; introns spliced out; considered “dead on arrival”
- often found by automatic prediction algorithms

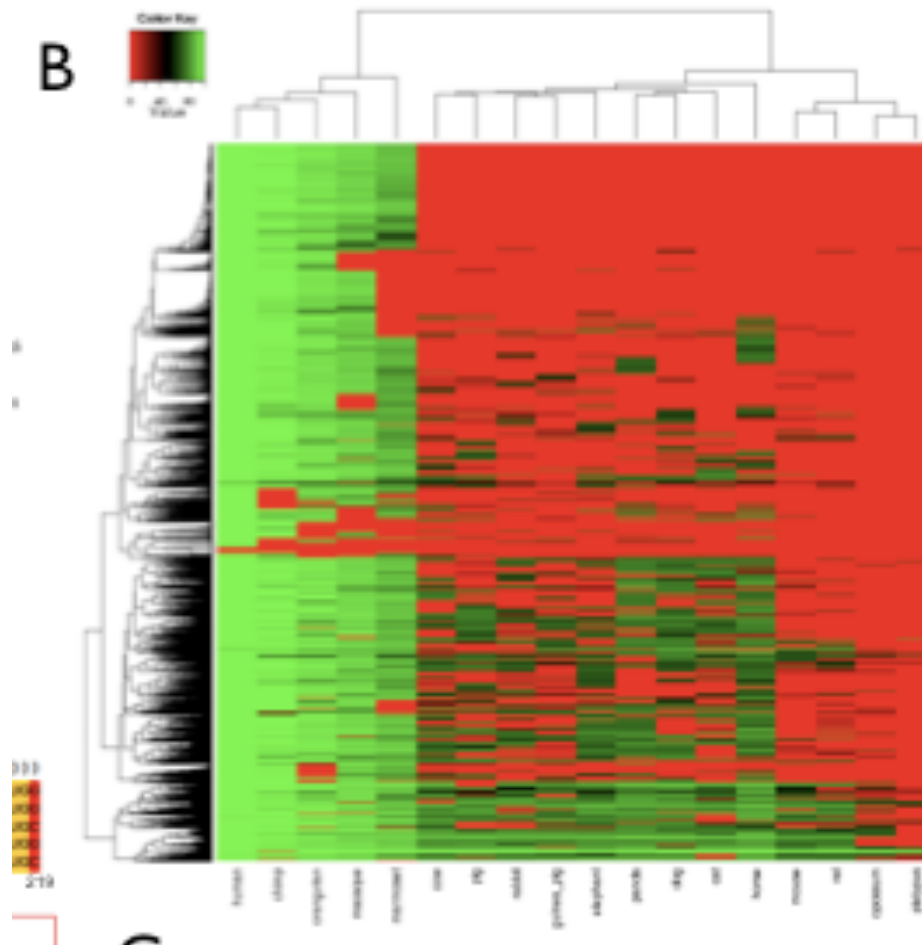


# Unprocessed pseudogene



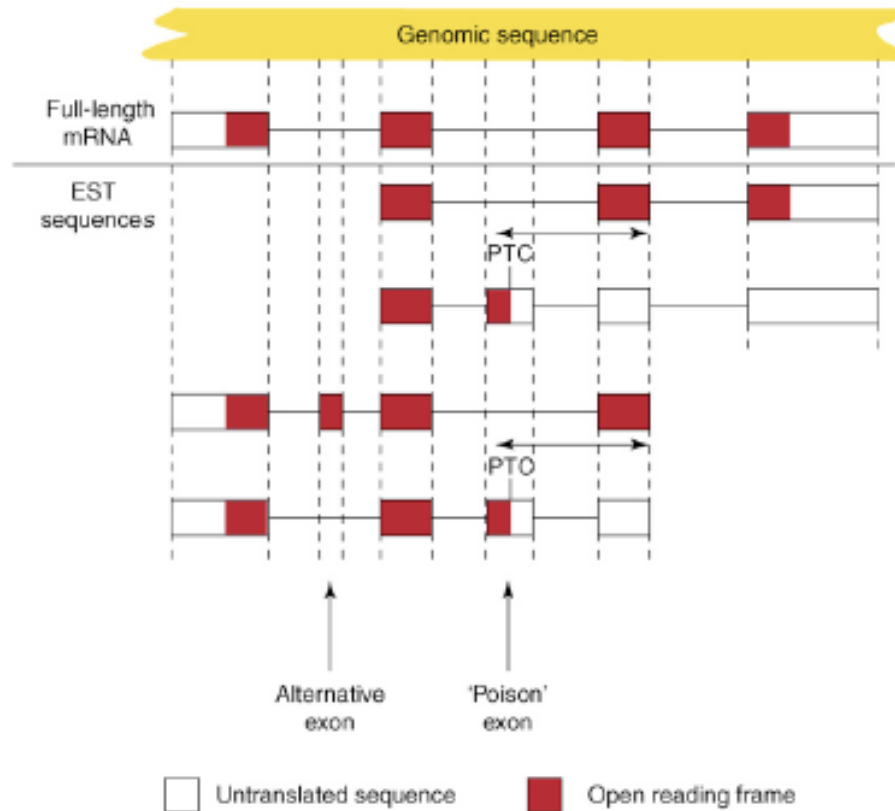
- Gene duplication
- Intact exon-intron structure
- Acquire mutations that result in loss of function
- Identified by lack of CDS

# Are lncRNAs conserved?



Guigo's group generated Heat map using blast and exonerate analysis of 9,200 loci from gencode 7 to examine conservation in mammals (platypus/opossum)

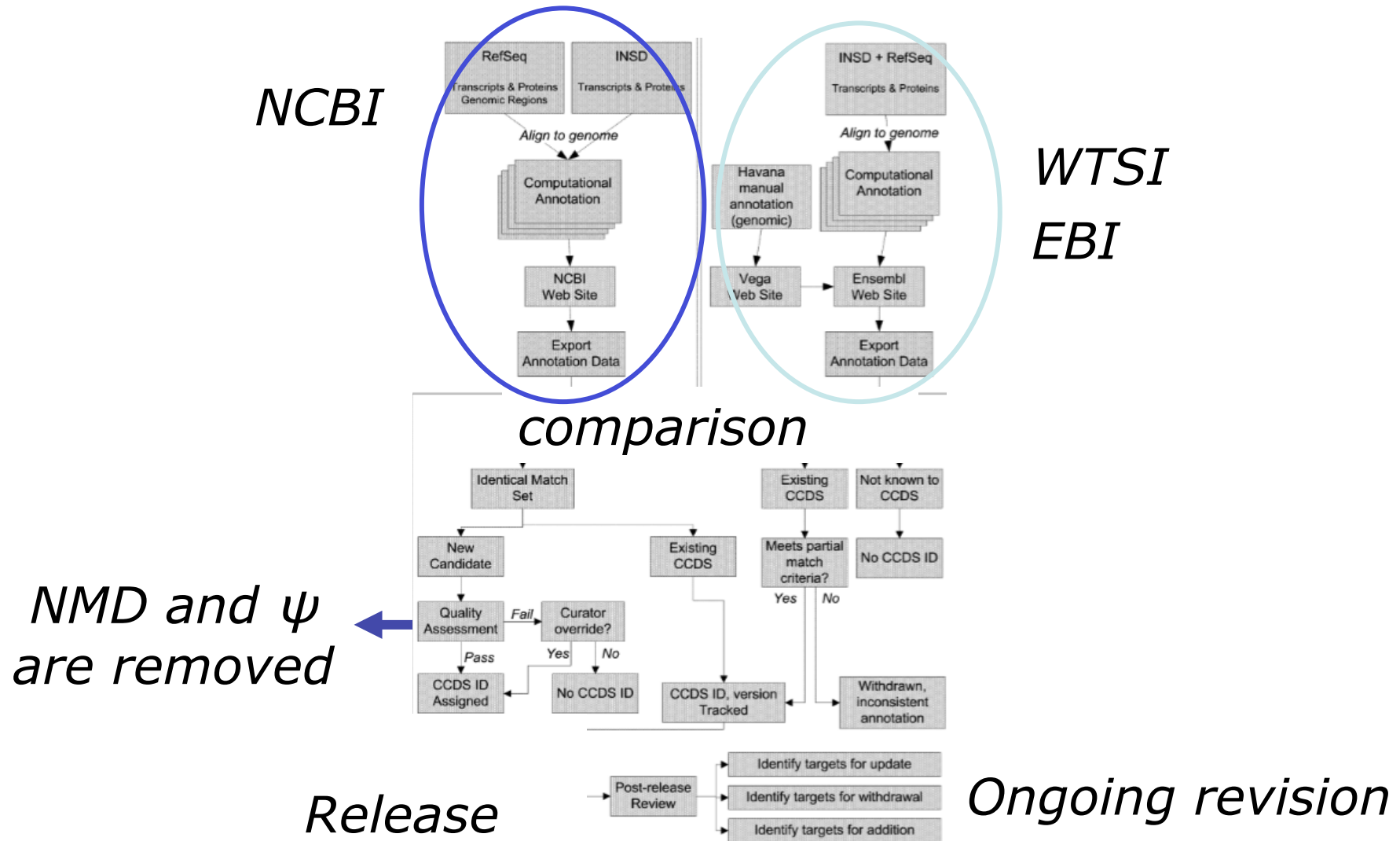
# Identification of Nonsense mediated decay (NMD) (Stephen Brenner)



PTC=  
Premature  
Termination  
Codon

TIBs Vol 33:8

# CCDS pipeline: producing consensus



# Human and Vertebrate analysis and Annotation (HAVANA)group

