# Module 2

## The Vega and UCSC Genome Browsers

### Using Web Browsers to View Genome Annotation

Jane Loveland PhD
Wellcome Trust Sanger Institute
Hinxton, UK

# Genome Browsers / Gene Sets

NCBI     http://www.ncbi.nlm.nih.gov/RefSeq/

e!Ensembl     http://www.ensembl.org

Vega     http://vega.sanger.ac.uk

UCSC     http://genome.ucsc.edu/

CCDS Database     http://www.ncbi.nlm.nih.gov/projects/CCDS/
CcdsBrowse.cgi (CDS only)

# NCBI Map Viewer

# NCBI Map Viewer

- Excellent integration with other NCBI resources

- Best "map" views of non-sequence maps (i.e. clone maps, genetic maps)

- Includes Celera assembly, alternate haplotypes, assemblies of everything available

- BLAST for sequence searching

# NCBI - RefSeq

- Non-redundant gene set
- Accessed via browsers or Entrez Gene
- Accessions for genomic DNA, transcripts and proteins
- Primarily protein-coding
- Semi-curated

|          | **Automated** | **Curated** |
|----------|-----------|----------|
| Genomic  | NC_12345  |          |
| mRNA     | XM_12345  | NM_12345 |
| ncRNA    | XR_12345  | NR_12345 |
| Protein  | XP_12345  | NP_12345 |

# http://www.ncbi.nlm.nih.gov/gene

# Vertebrate Genome Annotation Database



http://vega.sanger.ac.uk

# VEGA gene set

- Manually annotated using Otterlace/Zmap annotation software

- Based on direct pairwise alignment of mRNA, EST and protein evidence (including cross-species)

- Multiple biotypes, reflect confidence levels

- Includes additional data sources as DAS tracks (eg. CAGE tags, RNAseq)

# UCSC Genome Browser



Searching

Navigation

Display Controls

Annotations called "tracks"

# UCSC Genome Browser

- Straightforward display, easy navigation
- Third-party annotations
- Evolutionary conservation
- "Wiggle" tracks for continuous data
- Fast sequence searching with BLAT
- View your own data
- ENCODE annotations

# UCSC gene set

- Non-redundant gene set

- Automatic annotation based on BLAT alignments

- Transcripts require Genbank accession plus one other supporting feature (eg. Uniprot)

- Includes RefSeq models (require no additional support)

- Both protein-coding and non-coding

- Data hub for ENCODE data, displays GENCODE geneset

    http://genome.ucsc.edu/

# Common Functionality

- Navigational tools
  - Searching for markers by name or sequence (BLAST, BLAT)
  - Zooming in and out
  - Choose annotations to display
- Download of annotations
  - Whole genome
  - Specific regions
- Links to other resources
- View own data in browser

# Manual Analysis and Annotation pipeline: Otterlace



DAS=Distributed Annotation system

All annotation is supported by a combination of cDNA, EST and/or protein evidence
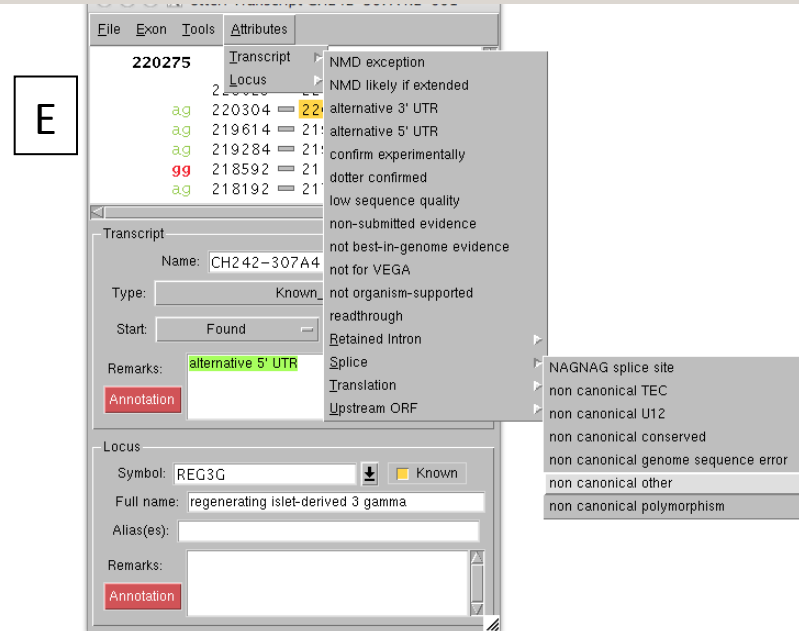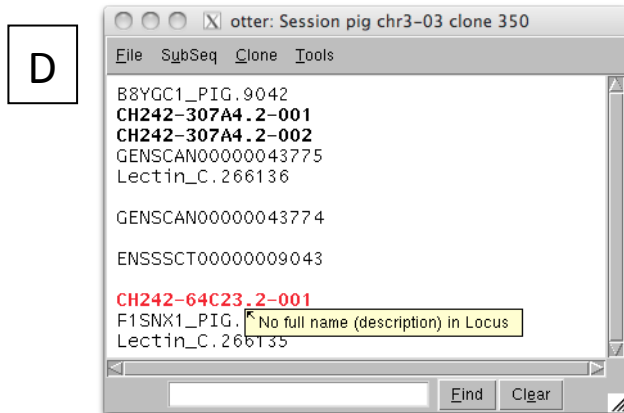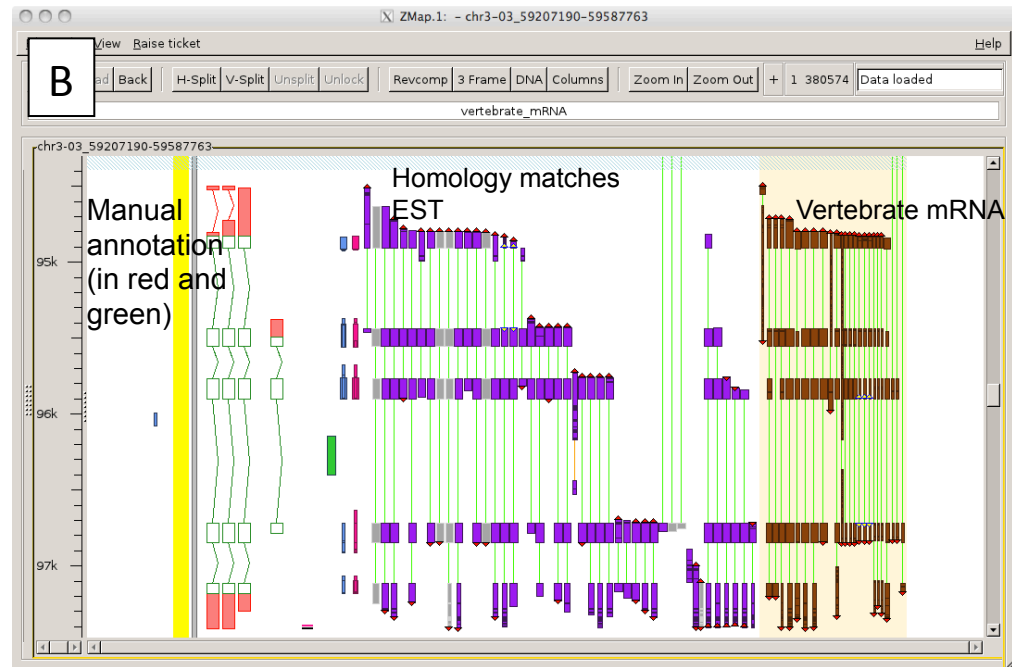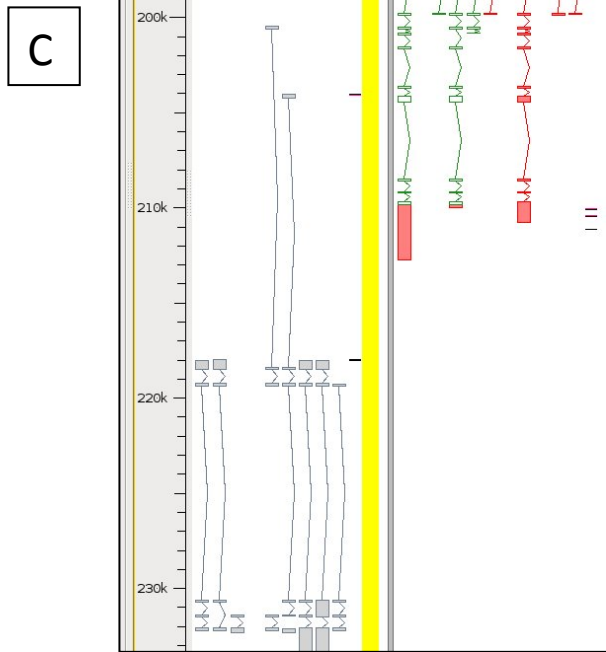
# Manual Annotation and Biotypes:
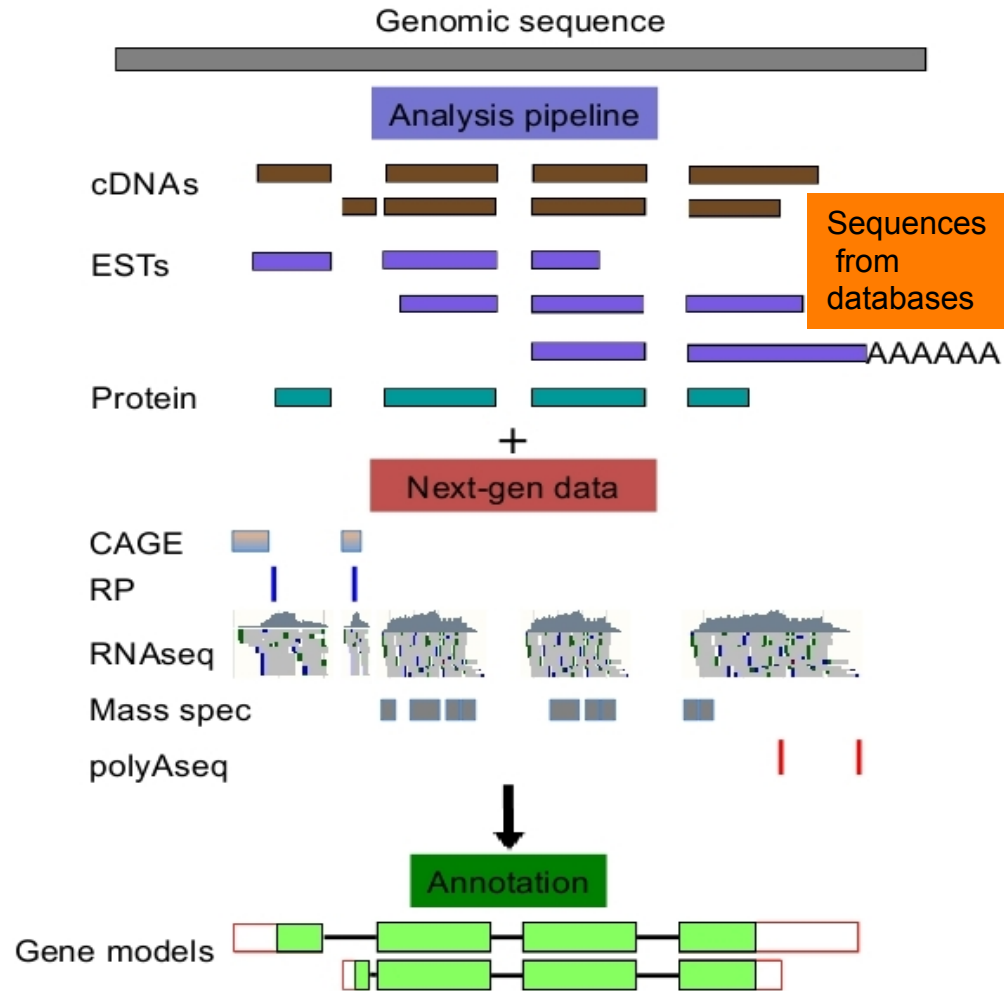
## Annotation:
### based on transcriptional evidence



## Biotypes

**Protein Coding**
- Known_CDS
- Novel_CDS
- Putative_CDS
- Nonsense_mediated_decay

**Transcript**   retained intron
putative

    **Non-coding**  lincRNA
        Antisense
        Sense_intronic
        Sense_overlapping
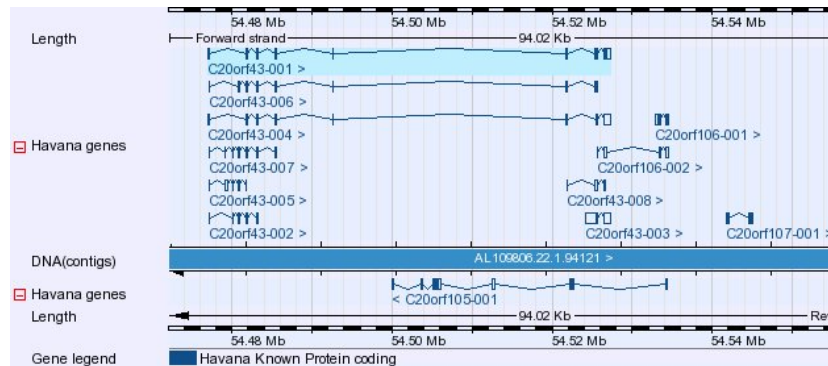        3'_overlapping_ncRNA

**Pseudogene**
- Processed
- Unprocessed
- Transcribed
- Translated
- Unitary
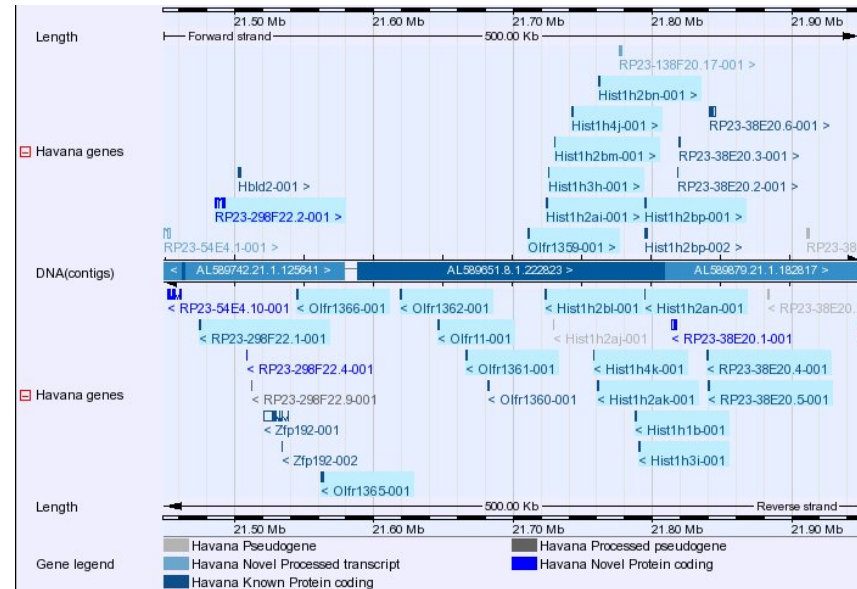- Polymorphic

**Immunoglobulin**
- IG_pseudogene
- IG_Gene
- TR_Gene

# Manual annotation is advantageous for:

- Overlapping genes
- Alternative splicing

- Pseudogenes
- Duplications/gene clusters





- Non-coding genes
- Complex loci e.g *GNAS*
- RNA seq data

- *Anything out of the ordinary*

# Classifying Functional Transcripts within a protein coding gene:

RefSeq



**TSC2 (chr16)**

Gencode

Partial CDS
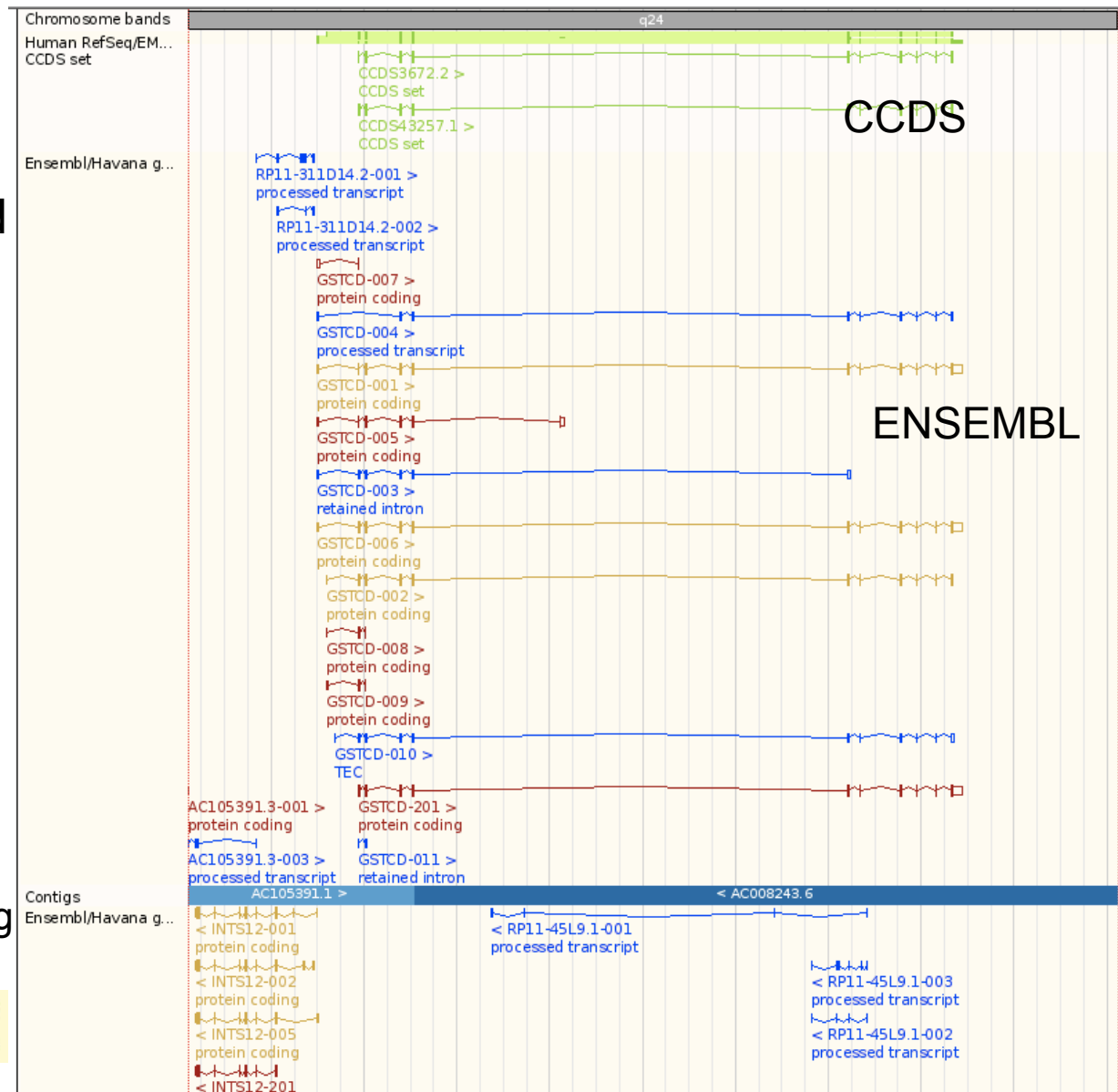
Nonsense-mediated decay

artifact

Retained intron

# Worked Example