# Tools for Community Genome Annotation:
# Zmap and Otter

## Jane Loveland

Havana group, Wellcome Trust Sanger Institute,
Hinxton, Cambridge, UK

**Havana: Human and vertebrate analysis and annotation**

- Manual annotation of human, mouse, zebrafish, pig and rat whole chromosomes or genomes

- Human GENCODE annotation and working on mouse GENCODE annotation

- Annotation of specific regions: human MHC & LRC haplotypes, multiple species MHCs & LRCs,

**Vega: Vertebrate Genome Annotation**

- Ensembl derived browser focusing on manual annotation

# Overview

- **Manual annotation process**

  – Tools, data, biotypes

- **Community Manual Annotation**

  – Mouse, Swine autosomes (IRAG), Rat, Chicken

- **New data and projects**

- **Manual annotation process**

  – Tools, data, biotypes

- Community Manual Annotation

  – Mouse, Swine autosomes (IRAG), Rat, Chicken

- New data and projects

# Automatic Annotation vs Manual

**Automatic Annotation**

- Quick whole genome analysis ~ weeks
- Consistent annotation
- Use unfinished/illumina sequence/shotgun assembly
- No polyA sites/signals, pseudogenes, lncRNAs
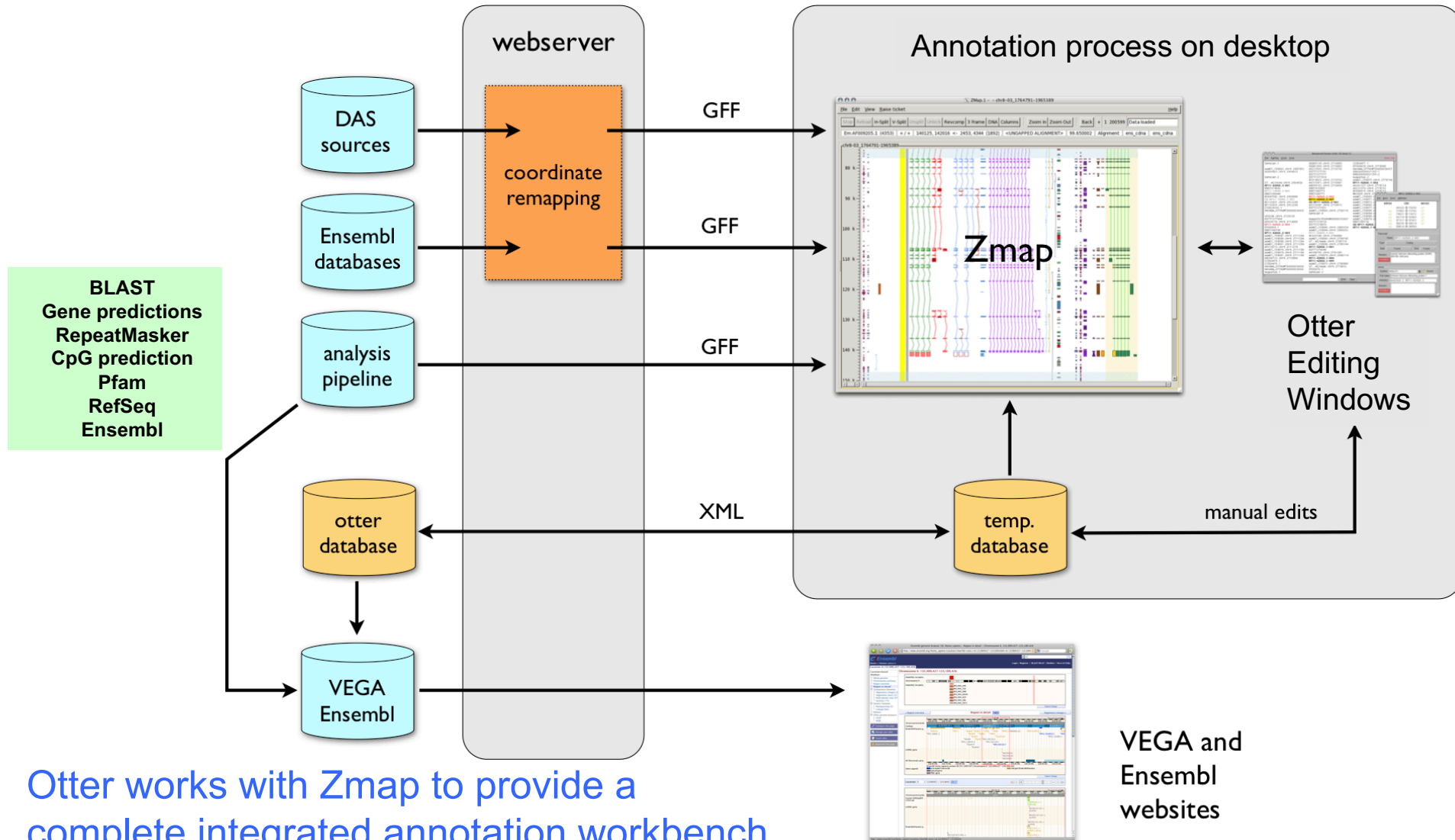- Limited functional annotation
- Predicts ~75% loci

**Manual Annotation**

- Slow~3 months per chromosome
- Prefer finished (high quality) sequence
- Flexible, can deal with inconsistencies in data
- Most rules have exception
- Consult publications as well as databases
- Extensive Biotypes:
  - Excellent functional annotation
  - e.g. pseudogenes, lncRNA

Automated annotation alone is not sufficient for researchers needs
GENCODE geneset

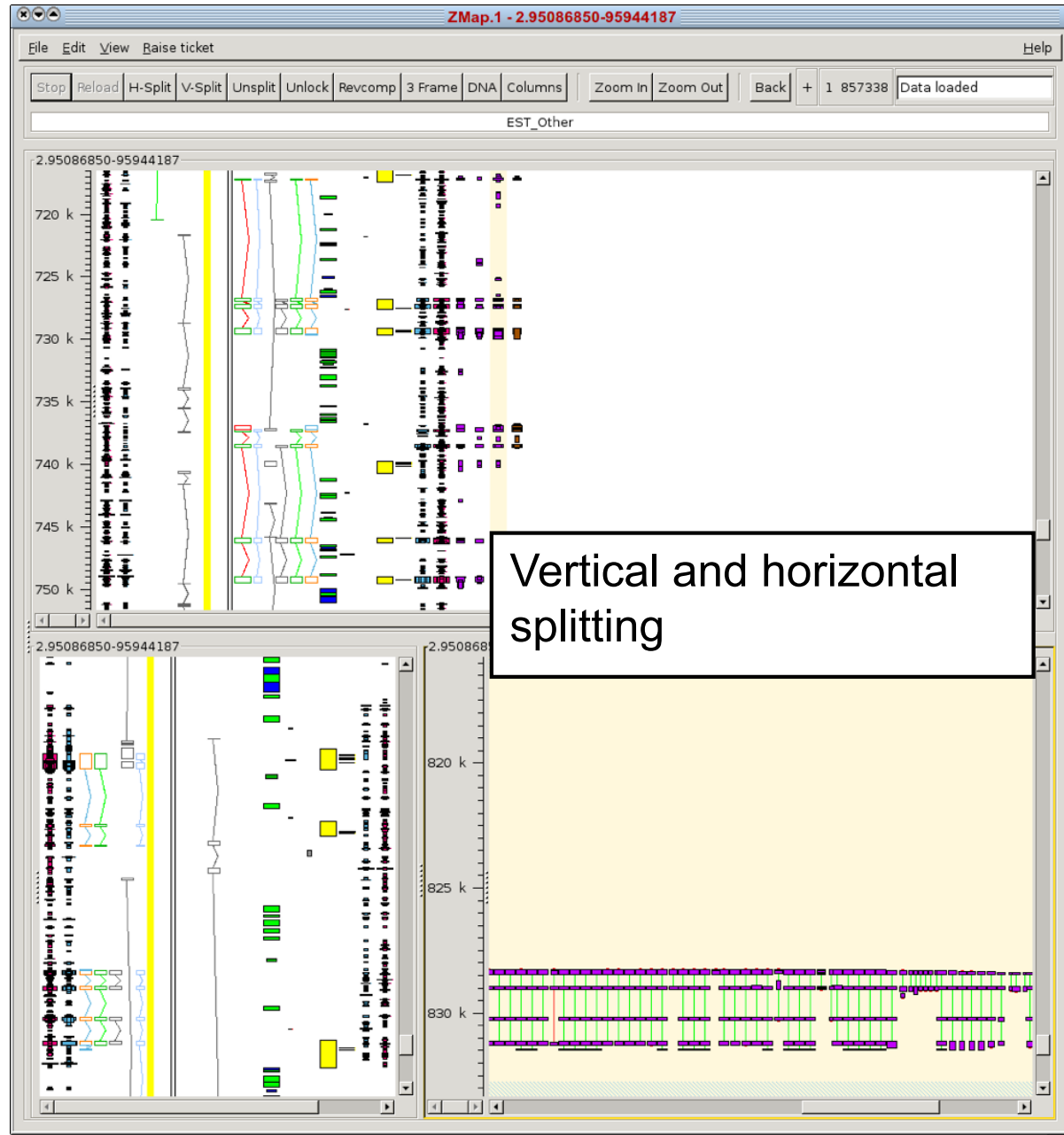# Analysis and Annotation pipeline: Otter/ZMap

Annotators can work anywhere as software communicates with server over HTTP
The server controls data access and allows multi-user annotation



Otter works with Zmap to provide a
complete integrated annotation workbench

# Annotation tools: Zmap



Vertical and horizontal splitting

# Annotation tools: Otter

# Annotation tools: Blixem

# Annotation tools: Dotter

# Manual Annotation: 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

Set of guidelines to help make annotation decisions

# Alternative Splicing:

Reference model

Skipped exon

Retained intron

Alternative splice donor

Alternative splice acceptor

Alternative first exon

Alternative final exon

Mutually exclusive

# 5' end annotation:
## *GPR56* (Human G protein-coupled receptor 56 gene)

# RNAseq data to extend 3'UTRs
## *GRIN2B*



Ensembl
Bodymap

Breast          Brain

polyA

Extended 3' end

# Improvements of lncRNA annotation: understanding functionality

# HAVANA Pseudogene Loci:



Processed

Unprocessed

AAAAAA

Reverse transcription and re-integration

Duplication

AAAAAA

*   *

*   *

# Sequencing error, pseudogene or polymorphic pseudogene?

## Polymorphic pseudogene



## Loss of function gene

Havana team move to EBI April 2017
VEGA is being archived and final release will be February 2017

# Ensembl view: GENCODE geneset

Gold (merged): agreed
ensembl/havana

Red: coding (001 Havana,
201 Ensembl)

Blue: non-coding

http://www.gencodegenes.org

# Update genes:

**Havana Rat update genes**

Show [All ▼] entries          Show/hide columns          Filter [          ]

| Gene Name | Biotype | Vega ID | Chromosome | Location | Modified date | New / Updated |
|---|---|---|---|---|---|---|
| Ak6 | protein_coding | OTTRNOG00000000242 | 2 | 49824275-49828050 | 2013-10-10 | new |
| Taf9 | protein_coding | OTTRNOG00000000241 | 2 | 49823853-49836100 | 2013-10-10 | new |
| Rn50_10_0878.8 | processed_pseudogene | OTTRNOG00000000239 | 10 | 87714193-87714513 | 2013-09-10 | new |
| Rn50_10_0877.2 | processed_pseudogene | OTTRNOG00000000238 | 10 | 87698915-87699931 | 2013-09-10 | new |
| Rn50_10_0877.1 | processed_pseudogene | OTTRNOG00000000237 | 10 | 87683090-87684368 | 2013-09-10 | new |
| Rn50_10_0876.2 | processed_pseudogene | OTTRNOG00000000236 | 10 | 87581818-87582174 | 2013-09-10 | new |
| Krt34 | processed_transcript | OTTRNOG00000000235 | 10 | 87736318-87740372 | 2013-09-10 | new |
| Rn50_10_0878.5 | processed_pseudogene | OTTRNOG00000000234 | 10 | 87732934-87733011 | 2013-09-10 | new |
| Rn50_10_0878.6 | processed_pseudogene | OTTRNOG00000000233 | 10 | 87722438-87723009 | 2013-09-10 | new |
| Rn50_10_0878.7 | processed_pseudogene | OTTRNOG00000000232 | 10 | 87708503-87708948 | 2013-09-10 | new |
| Krtap16-1 | processed_pseudogene | OTTRNOG00000000231 | 10 | 87705416-87706898 | 2013-09-10 | new |
| Krt32 | protein_coding | OTTRNOG00000000230 | 10 | 87785761-87792744 | 2013-09-10 | new |
| Krt31 | protein_coding | OTTRNOG00000000229 | 10 | 87745725-87749163 | 2013-10-09 | new |
| Krt36 | protein_coding | OTTRNOG00000000228 | 10 | 87809278-87812643 | 2013-09-06 | new |
| Krt35 | protein_coding | OTTRNOG00000000227 | 10 | 87798389-87801786 | 2013-09-06 | new |
| Ka11 | protein_coding | OTTRNOG00000000226 | 10 | 87934150-87937029 | 2013-09-06 | new |
| Krt9 | protein_coding | OTTRNOG00000000225 | 10 | 87895711-87899956 | 2013-09-06 | new |
| Rn50_10_0879.5 | processed_transcript | OTTRNOG00000000224 | 10 | 87871753-87877455 | 2013-09-06 | new |
| Rn50_10_0879.6 | processed_transcript | OTTRNOG00000000223 | 10 | 87868731-87870238 | 2013-09-06 | new |

New Track Hub:



ftp://ngs.sanger.ac.uk/production/gencode/update_trackhub/hub.txt

- Manual annotation process
  - Tools, data, biotypes

- **Community Manual Annotation**
  - Mouse, Swine autosomes (IRAG), Rat, Chicken

- New data and projects

# Community Annotation:

- Part of IKMC with EUCOMM annotation in mouse:
  - KOMP and NorCOMM annotation ("Blessed Annotator")

- Jamborees for species with strong community interest ("Gatekeeper"):
  - *Xenopus tropicalis* 2005 (cDNA)
  - Cow 2007 (Genomic WGS)

  - Pig
  2008 (Genomic WGS)
  2010 - 2013
    - IR genes in Pig (~1300 genes) manually annotated by community

Dawson et al. BMC Genomics 2013, 14:332
http://www.biomedcentral.com/1471-2164/14/332

BMC Genomics

RESEARCH ARTICLE                    Open Access

Structural and functional annotation of the porcine immunome

Many transcript variants
Found gene expansions and duplications
Co-expression clustering analysis: some exhibited accelerated evolution

- Rat manual annotation 2013, 2015 (BBSRC)

- Chicken MHC 2016

wellcome trust
sanger
institute

havana
human and vertebrate analysis and annotation

# Community Annotation Approaches:

The value of a genome is only as good as its annotation
**Rat** whole genome annotation of Rnor 6.0
**Chicken** MHC

**Otter/Zmap Annotation Software**
Authentication:
Sanger single sign-on account (email)
Registered email for Otter permitted users:
Access to our data and and analysis pipeline

**Mac and Linux:** Platforms of choice
Monthly updates/bugfixes

**Windows:**
Virtual machine image installed and run using VirtualBox
Runs an Ubuntu desktop with a bespoke Otter release

Over 2600 genes manually annotated that have been chosen by the rat community and RGD. Targeted annotation.

Final rat Vega and Ensembl merge in progress

- Manual annotation process

  – Tools, data, biotypes

- Community Manual Annotation

  – Mouse, Swine autosomes (IRAG), Rat, Chicken

- New data and projects

Long-read data:
Better for discovering novel alternatively spliced transcripts and full-length transcripts

Targeted sequencing for gene discovery and quantification using RNA CaptureSeq

Tim R Mercer, Michael B Clark, Joanna Crawford, Marion E Brunck, Daniel J Gerhardt, Ryan J Taft, Lars K Nielsen, Marcel E Dinger & John S Mattick

PacBio-CaptureSeq
(human: brain, testis, heart, liver, HeLa, K562)
(mouse: brain, testis, heart, liver, E7, E15)

SLR-RNAseq
(human/mouse brain)

ARTICLES

Comprehensive transcriptome analysis using synthetic long-read sequencing reveals molecular co-association of distant splicing events

Hagen Tilgner[1,3], Fereshteh Jahanbani[1,3], Tim Blauwkamp[2], Ali Moshrefi[2], Erich Jaeger[2], Feng Chen[2], Itamar Harel[1], Carlos D Bustamante[1], Morten Rasmussen[1] & Michael P Snyder[1]

wellcome trust sanger institute

havana
human and vertebrate analysis and annotation

Computationally combining evidence:

Havana annot.

PolyA-seq

CAGE

Merged models
SLR  PacBio

PacBio reads

SLR-seq reads

EDEM3

wellcome trust
sanger
institute

havana
human and vertebrate analysis and annotation

## Mouse strains:



### gEVAL

Tools | Help & Document

**Browse the Mouse Genome**

Click on a link below to go to the assembly's home page.

**Interim builds on current paths**

browse current path, 24th Septermber 2014

**AGP builds on public reference releases**

**AGP Viewer** - reference assembly, finished clones
browse AGP build **GRCm38p3**

**AGP Viewer** - reference assembly, finished clones
browse AGP build **NCBIm37**

**Other mouse assemblies**

**WGS MGSCv3 viewer** - Whole Genome Shotgun Supe
browse

**WGS Celera build viewer** - Whole Genome Shotgun
browse

**WGS C57BL/6J** - Whole Genome Shotgun Supercontigs
browse

### Browse the Mouse Strain Assemblies

**Commonly viewed genomes**

| | |
|---|---|
| **Mouse** 129S1_SvImJ_R | **Mouse** A_J_R |
| **Mouse** AKR_J_R | **Mouse** BALB_cJ_R |
| **Mouse** C3H_HeJ_R | **Mouse** C57BL_6NJ_R |
| **Mouse** CAST_EiJ_R | **Mouse** CBA_J_R |
| **Mouse** DBA_J_R | **Mouse** FVB_NJ_R |
| **Mouse** LP_J_R | **Mouse** NOD_ShiLtJ_R |
| **Mouse** NZO_HlLtJ_R | **Mouse** PWK_PhJ_R |
| **Mouse** SPRET_EiJ_R | **Mouse** WSB_EiJ_R |
| **Mouse** CAROLI_EiJ | **Mouse** Pahari_EiJ |

Search gEVAL...

n assemblies are viewable via the dedicated mice gEVAL

ath mouse chr2 05.02.2015

# Mouse strain annotation reveals new genes:



Itgb3

Cage tags = 268
5808aa

Ens_RNASeq_ENCODE_canonical_introns_plus

200k

300k

400k

500k

600k

Efcab3
Cage tags = 2181

Mettl2

Ruth Bennett

wellcome trust
sanger
institute

havana
human and vertebrate analysis and annotation

# Mouse strain annotation reveals strain specific coding transcripts: *Ifi214*



393aa NMD    420aa coding

Reference C57BL/6          129S1-SvImJ

>Reference longest NMD transcript 393 AA.
MVNEYKRIVLLTGLMGINDHDFRMVKSLLSKELKLNKMQDEYDRVKI
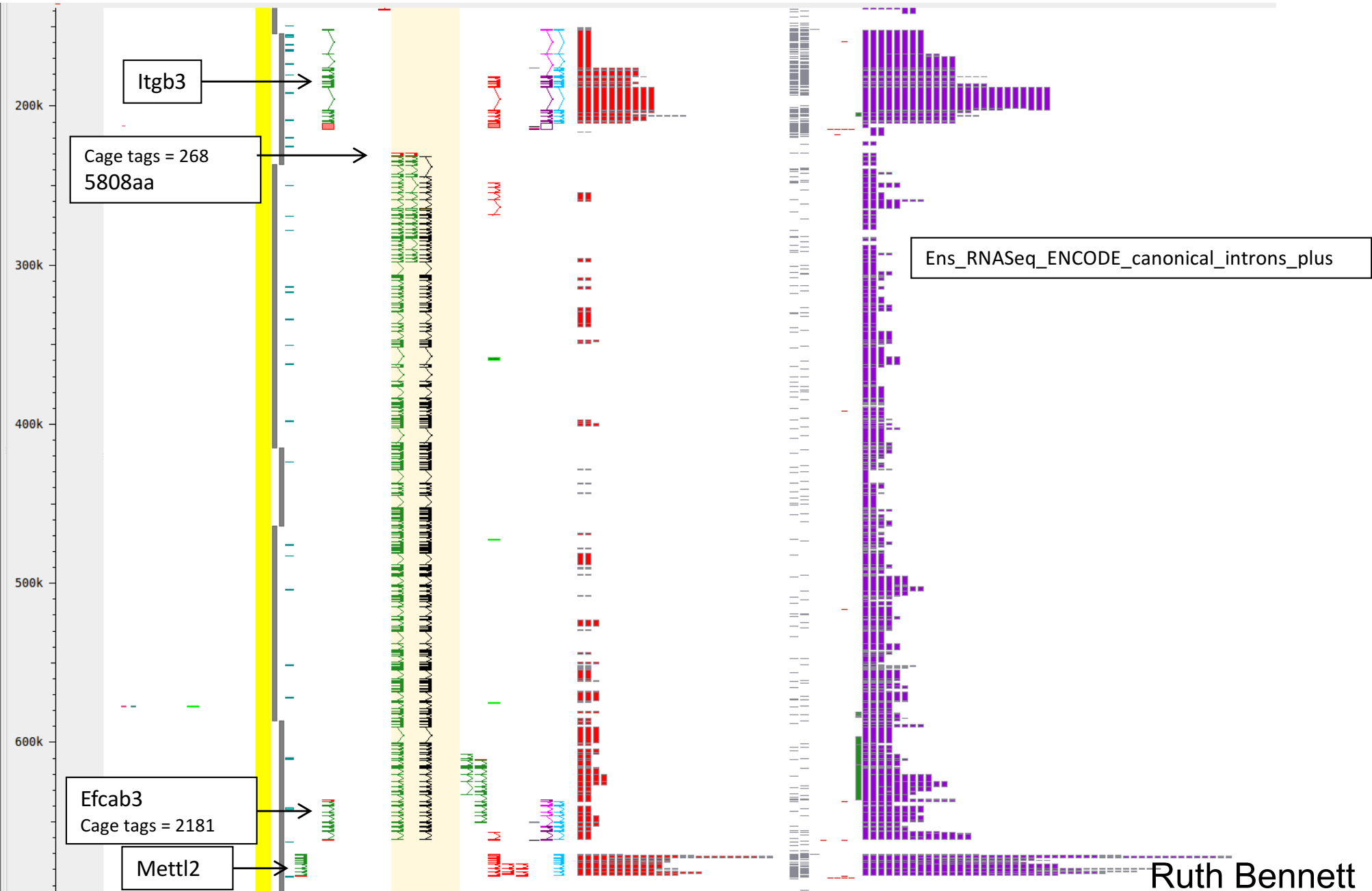ADLMEDKFPKDAGVVQLIKLYKQIPGLGDIANKLKNEKAKAKRKGKG
KRKTAAKRQRQEEPSTSQPMSTTNEDAEPESGRSTPDTQVAQLSLPT
ASRRNQAIQISPTIASSSGQTSSRSSETLQSIIQSPETPTRSSSRIL
DPPVSPGTAYSSAQALGVLLATPAKRQRLKNVPKEPSEENGYQQGSK
KVMVLKVTEPFAYDMKGEKMFHATVATETEFFRVKVFDIVLKEKFIP
NKVLTISNYVGCNGFINIYSASSVSEVNDGEPMNIPLSLRKSANRTP
KINYLCSKRRGIFVNGVFTVCKKEERGYYICYEIGDDTGMMEVEVYG
RLTNIACNPGDKLRLML*

Stop codon isn't a SNP. Caused by a much larger disruption.

>129 Patch long coding transcript (equivalent to ref NMD) 420 AA.
MVNEYKRIVLLTGLMGINDHDFRMVKSLLSKELKLNKMQDEYDRVKI
ADLMEDKFPKDAGVVQLIKLYKQIPGLGDIANKLKNEKAKAKRKGKG
KRKTAAKRQRQEEPSTSQPMSTTNEDAEPESGRSTPDTQVAQLSLPT
ASRRNQAIQISPTIASSSGQTSSRSSETLQSIIQSPETPTRSSSRIL
DPPVSPGTAYSSAQALGVLLATPAKRQRLKNVPKEPSEENGYQLGSK
KVMVLKVTEPFAYDMKGEKMFHATVATETEFFRVKVFDIVLKEKFIP
NKVLTISNYVGCNGFINIYSASSVSEVNDGEPMNIPLSLRKSANRTP
KINYLCSKRRGIFVNGVFTVCKKEERGYYICYEIGDDTGMMEVEVYG
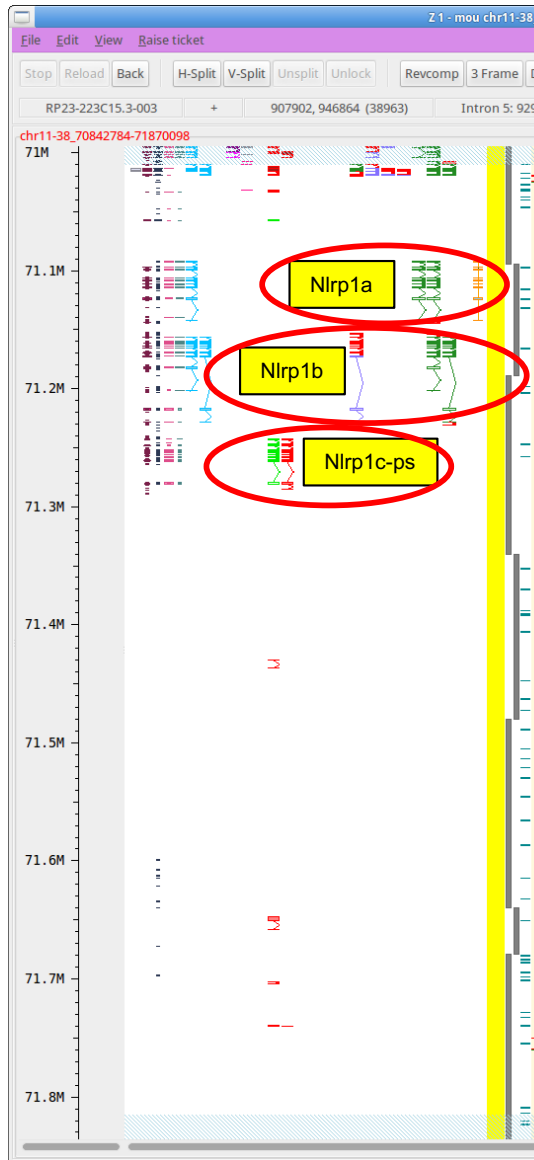RLTNIACNPGDKLRLICFELTPDEETAWLRSTTHSNMQVIKARN*

Ruth Bennett

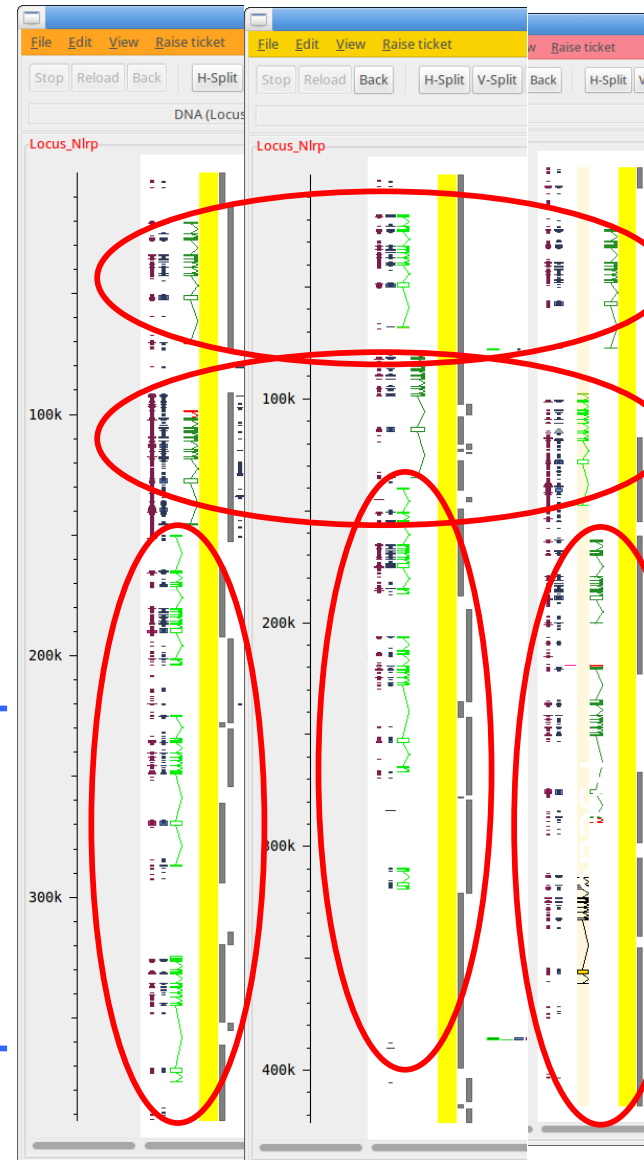Mouse strain annotation reveals strain specific expansions and biotype differences: Nlrp locus  Reference: C57BL/6    Strains: PWK/PhJ   WSB/EiJ   CAST/EiJ
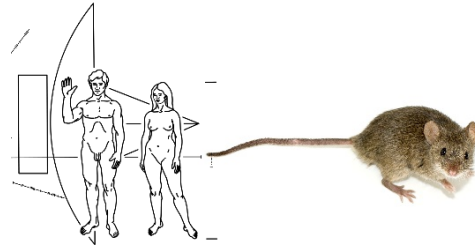
Ruth Bennett

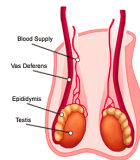# Clinical data: RNA capture-seq across a range of tissues



| HUMAN | MOUSE |
|---|---|
| 1. Brain | 1. Brain |
| 2. Testis | 2. Testis |
| 3. Heart | 3. Heart |
| 4. Liver | 4. Liver |
| 5. HeLa | 5. Embryo 7d |
| 6. K562 | 6. Embryo 15d |

# Clinical Data
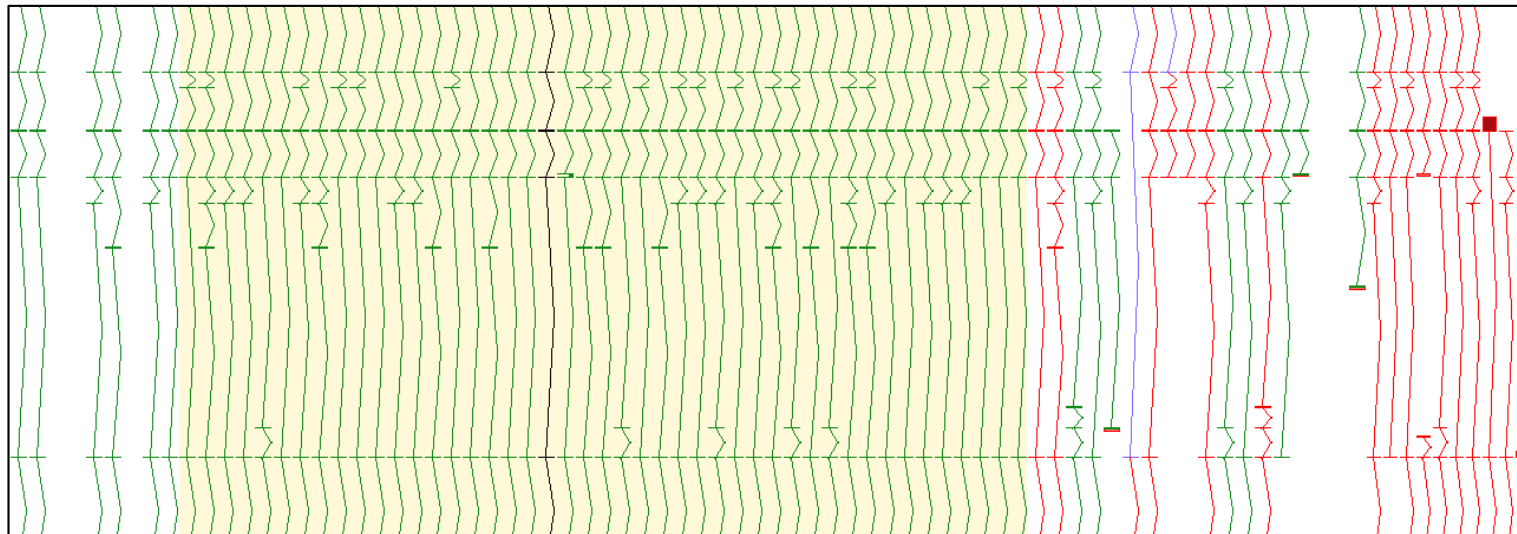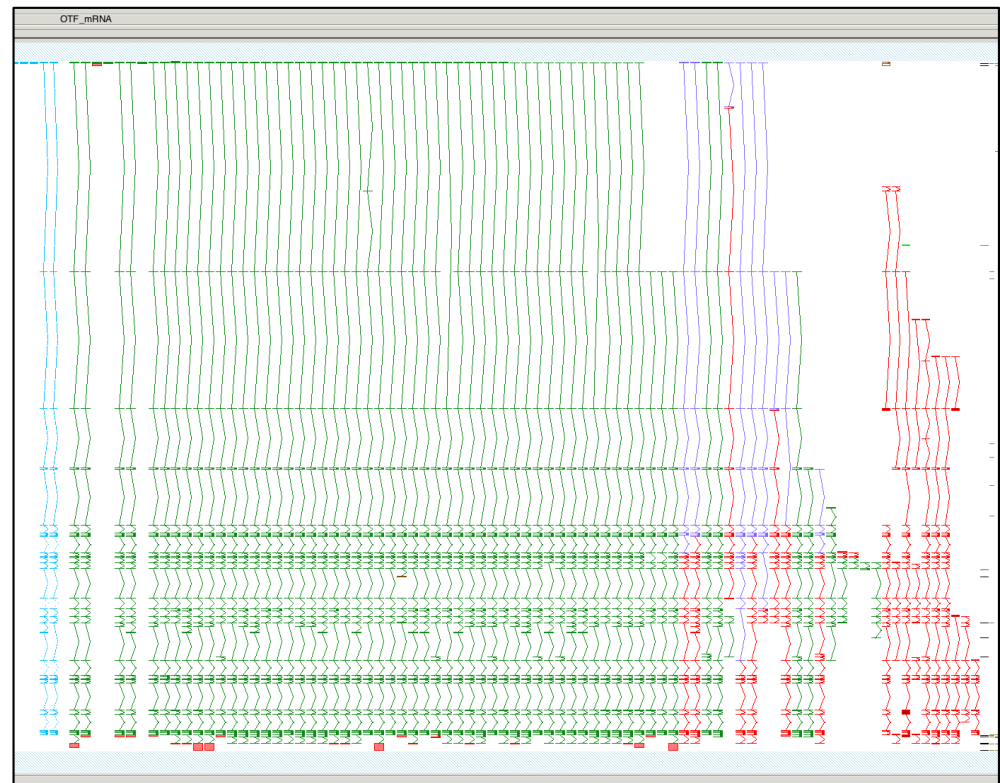
## Human *KCNMA1*:

22 original transcripts

Now 92 transcripts

25 from SLR-Seq

SLR-Seq also extended other transcripts

Micro exons, NAGNAG, alternate exon use

Marie-Marthe Suner

# Final Comments

Zmap/otter tool
       Can handle large datasets
       Training, QC and feedback

Let the pipelines take the strain
       Targeted manual annotation
       Spend time on the tricky things

Feedback genome quality
       Report errors
       Improve the assembly

wellcome trust
**sanger**
institute

**havana**
human and vertebrate analysis and annotation

# Acknowledgements

**Havana:**
Adam Frankish
If Barnes
Ruth Bennett
Andrew Berry
Alex Bignell
Claire Davidson
Gloria Despacio-Reyes
Sarah Donaldson
Matthew Hardy
Toby Hunt
Mike Kay
Jane Loveland
Deepa Manthravadi
Gaurab Mukherjee
Jonathan Mudge
Marie-Marthe Suner
Mark Thomas

**Gencode:**
Jose Gonzalez
Stephen Fitzgerald

**Annosoft:**
Ed Griffiths
Steve Miller

**Vega:**
Stephen Trevanion
Dan Sheppard

ftp://ngs.sanger.ac.uk/production/gencode/update_trackhub/hub.txt

wellcome trust
**sanger**
institute

**havana**
human and vertebrate analysis and annotation