

Can we apply lessons learned from manual annotation in human and mouse to wheat?

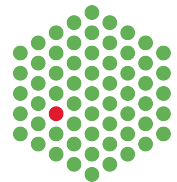
Jane Loveland

Annotation Project leader
Ensembl-HAVANA

PAG XXVI, 16th January 2018



EMBL-EBI



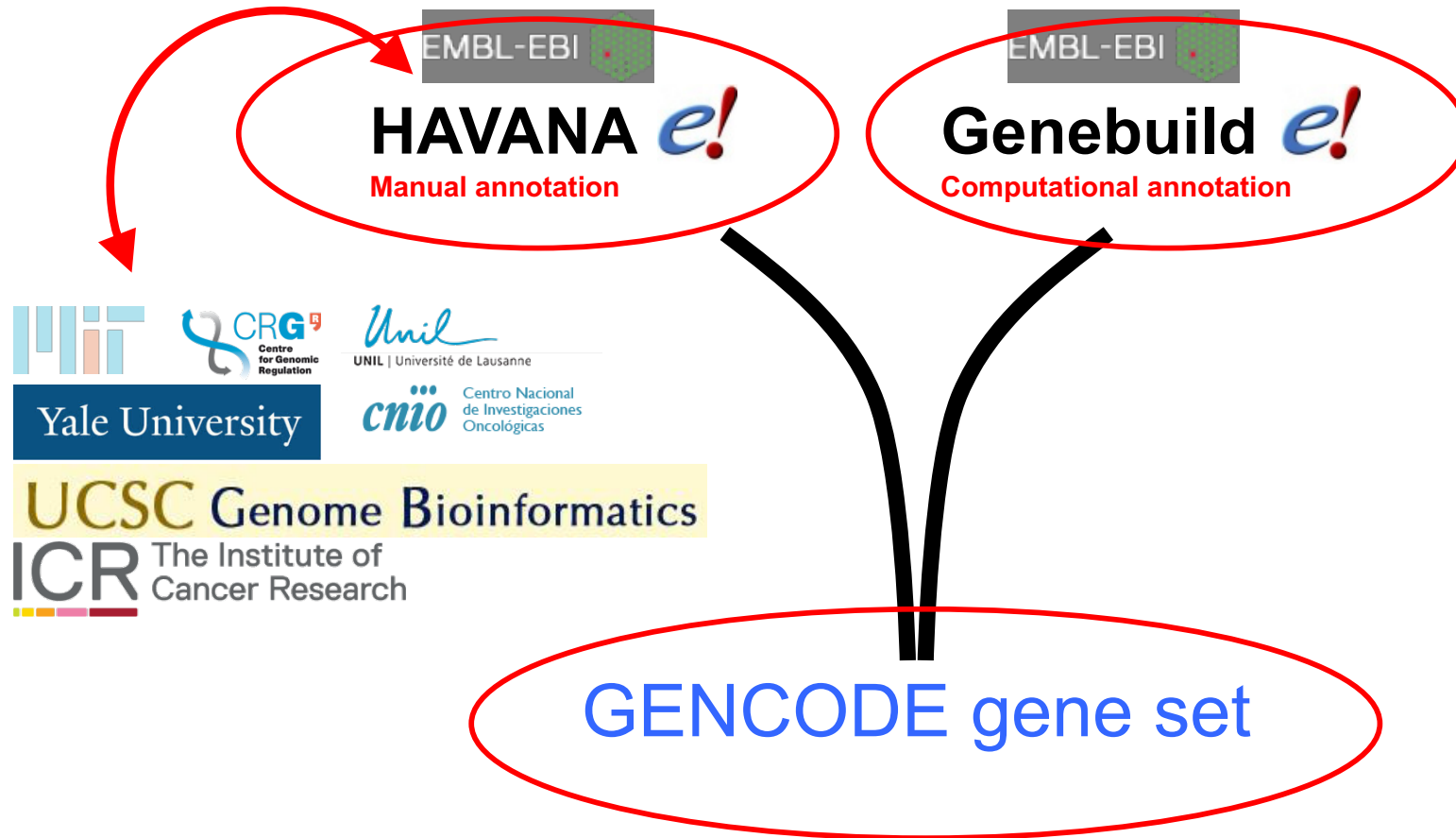
*e!*Ensembl



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,

The GENCODE consortium



The HAVANA team

GENCODE



Whole Genome
or chromosome



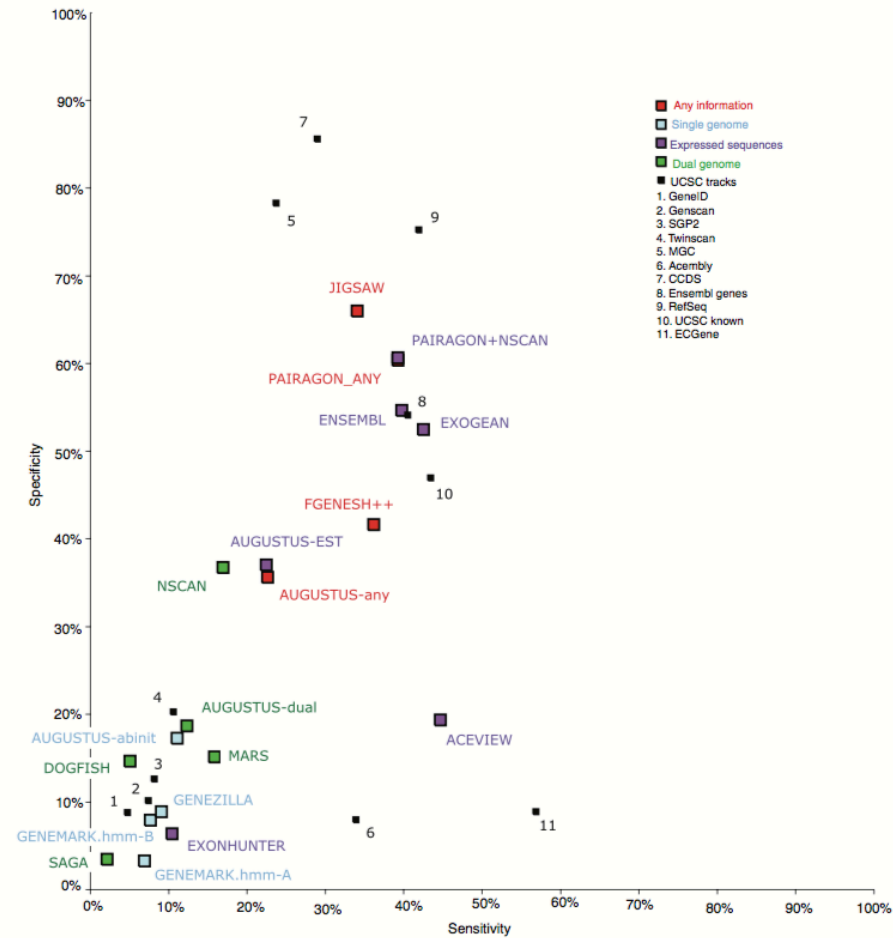
Targeted regions
or genes



Community projects



Requirement for manual annotation



EGASP - Guigó et al. Genome Biol. 7 Suppl 1:S2.1-31 (2006)

Requirement for manual annotation



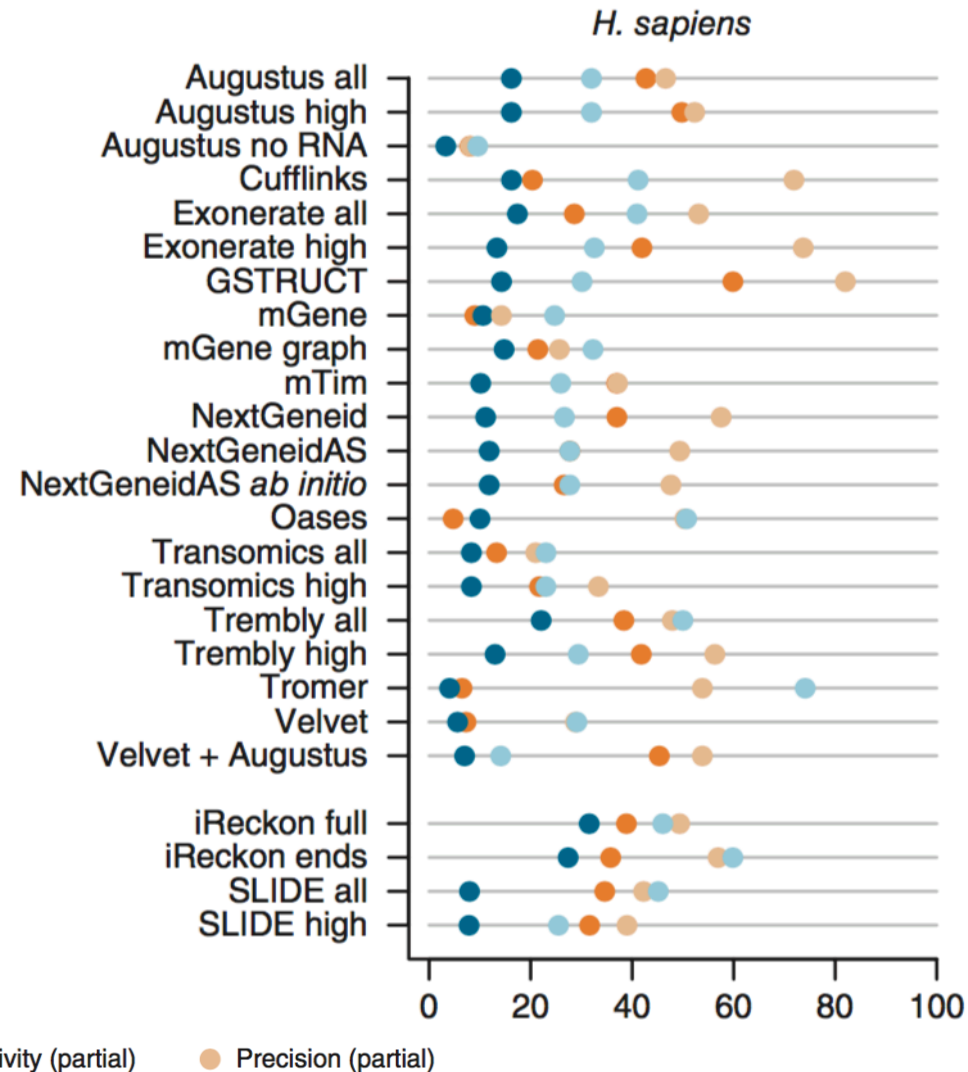
Best methods:

Sensitivity < 50%

Specificity < 70%

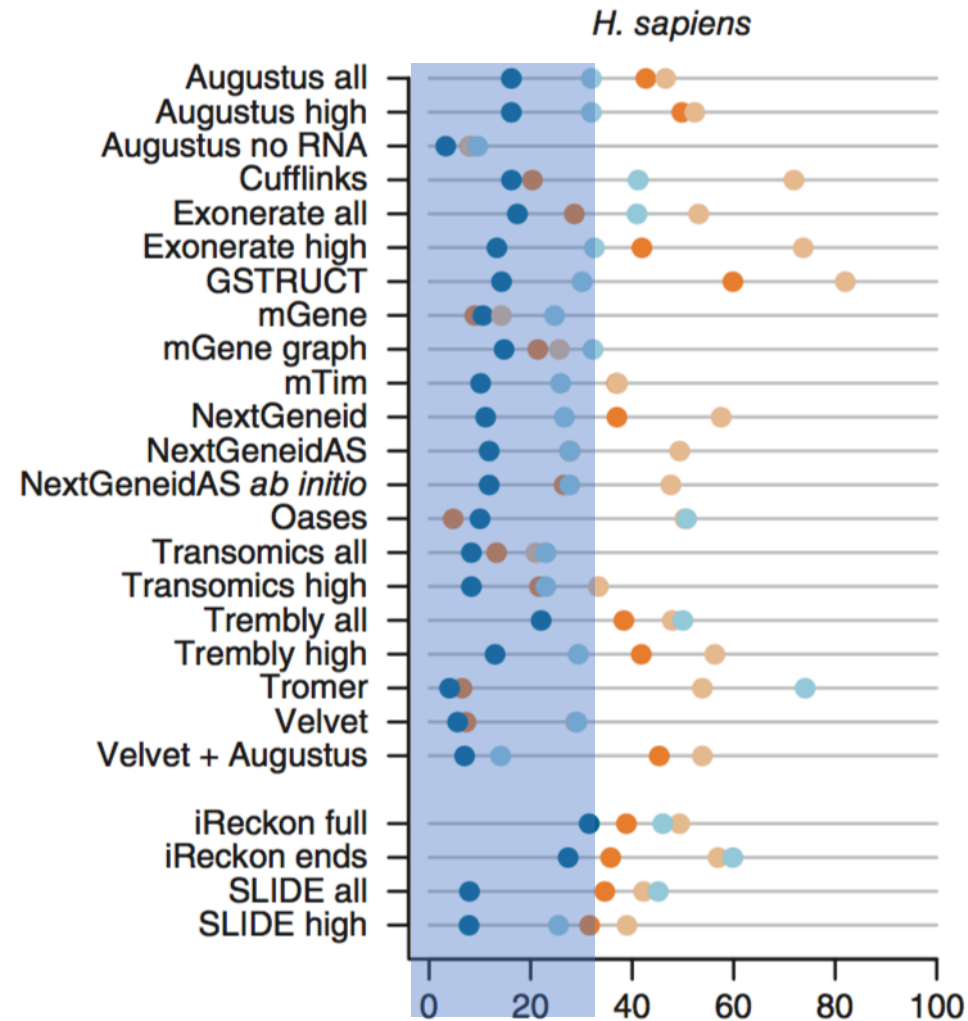
EGASP - Guigó et al. Genome Biol. 7 Suppl 1:S2.1-31 (2006)

Requirement for manual annotation



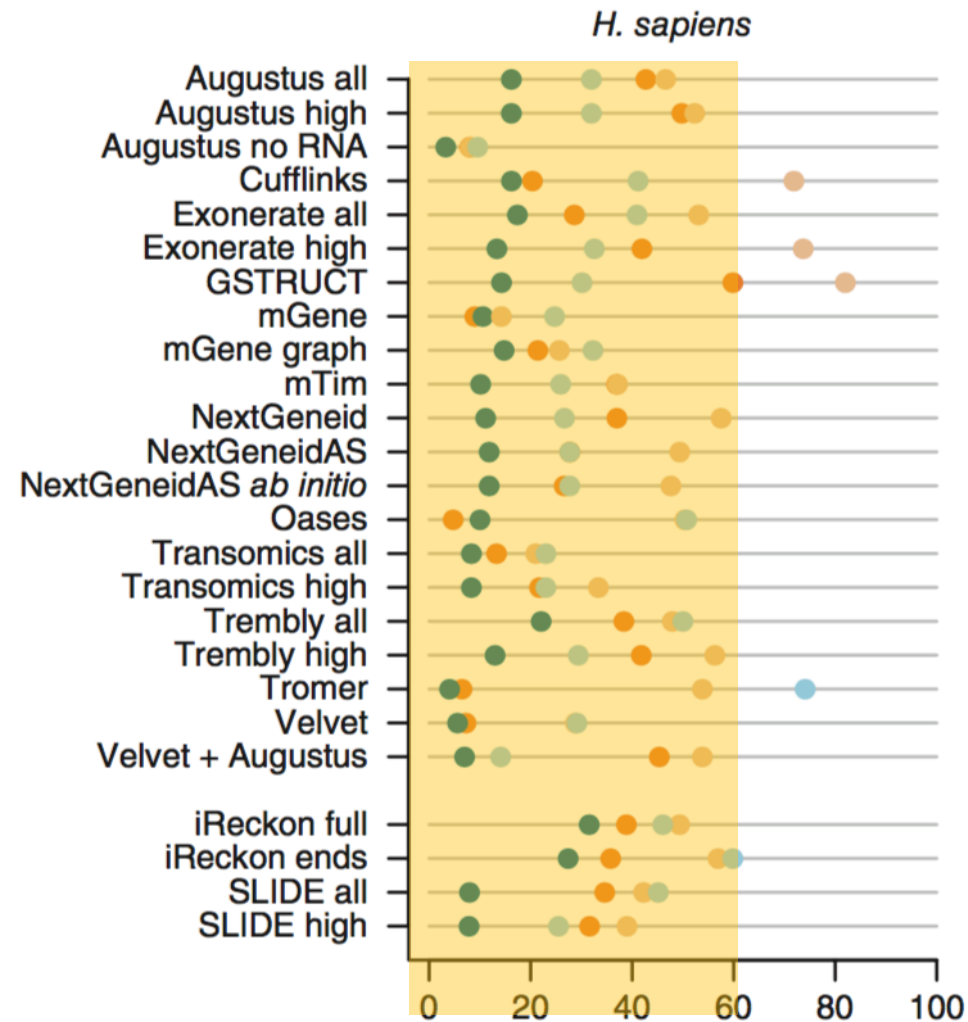
RGASP - Steijger et al. Nat Methods. 10(12):1177-84 (2013)

Requirement for manual annotation



RGASP - Steijger et al. Nat Methods. 10(12):1177-84 (2013)

Requirement for manual annotation



RGASP - Steijger et al. Nat Methods. 10(12):1177-84 (2013)

Manual annotation supports automated annotation

Gene: ENSBTAG00000000573

Location [Chromosome 19: 42,260,882-42,285,621](#) reverse strand.
UMD3.1:GK000019.2

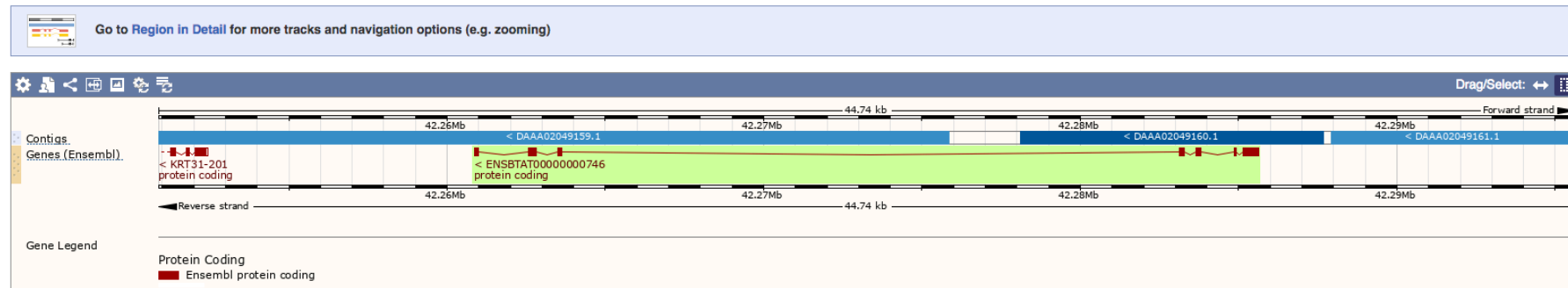
About this gene This gene has 1 transcript ([splice variant](#)), [74 orthologues](#), [12 paralogues](#) and is a member of [1 Ensembl protein family](#).

Transcripts [Hide transcript table](#)

Name	Transcript ID	bp	Protein	Translation ID	Biotype	UniProt	Flags
-	ENSBTAT00000000746.5	1359	452aa	ENSBTAP00000000746	Protein coding	E1BFG1	

Summary

Ensembl version ENSBTAG00000000573.5
Gene type Protein coding
Annotation method Annotation produced by the Ensembl [genebuild](#).



Manual annotation supports automated annotation

Gene: ENSBTAG00000000573

Location [Chromosome 19: 42,260,882-42,285,621](#) reverse strand.
UMD3.1:GK000019.2

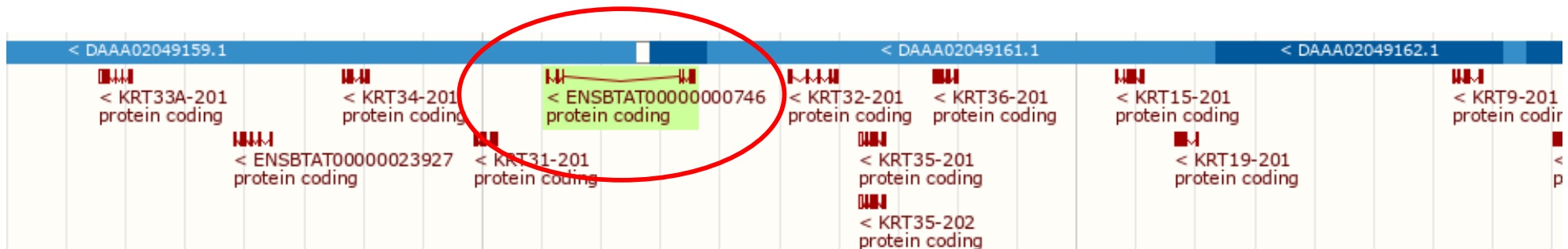
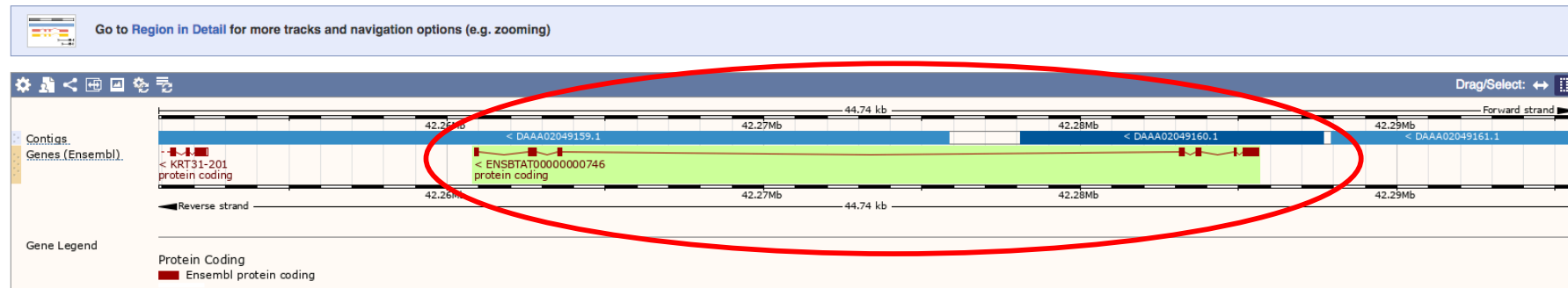
About this gene This gene has 1 transcript ([splice variant](#)), [74 orthologues](#), [12 paralogues](#) and is a member of [1 Ensembl protein family](#).

Transcripts [Hide transcript table](#)

Name	Transcript ID	bp	Protein	Translation ID	Biotype	UniProt	Flags
-	ENSBTAT00000000746.5	1359	452aa	ENSBTAP00000000746	Protein coding	E1BFG1	

Summary

Ensembl version ENSBTAG00000000573.5
Gene type Protein coding
Annotation method Annotation produced by the Ensembl [genebuild](#).



Automated annotation problems:

Transcript: ENSOART00000000372.1

Location: [Chromosome 1: 215,188,012-215,188,827](#) reverse strand.

About this transcript: This transcript has [5 exons](#).

Gene: This transcript is a product of gene [ENSOARG00000000352](#) [Hide transcript table](#)

Name	Transcript ID	bp	Protein	Translation ID	Biotype	Flags
-	ENSOART00000000372.1	540	No protein	-	Pseudogene	

Summary

< ENSOART00000000372 pseudogene
Reverse strand ————— 816 bp

Statistics: Exons: 5, Coding exons: 0, Transcript length: 540 bps,
Version: ENSOART00000000372.1
Type: Pseudogene
Annotation Method: Annotation produced by the Ensembl [genebuild](#).
Frameshift introns: [Frameshift introns](#) occur at intron number(s) [1](#), [3](#), [4](#).

Ensembl release 91 - Dec 2013

Frameshift introns are the length of 1, 2, 4, or 5 basepairs. They are introduced by the Ensembl [genebuild](#) in order to fit the cDNA sequence to the genome.

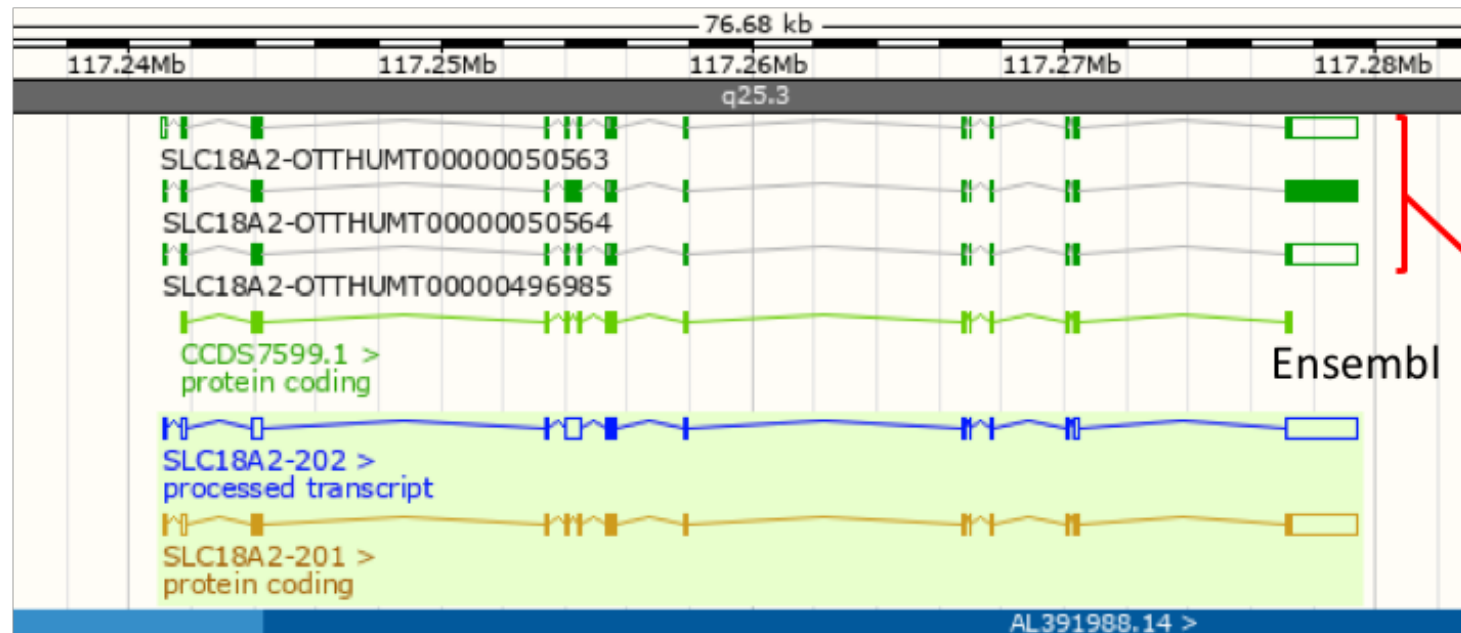
[Permanent link](#) - [View in archive site](#)

About Us Get help Our sister sites Follow us

Annotation issues:

Gaps, clusters and duplications, haplotype vs duplication, concatenation, overlapping genes, expansions and deletions, pseudogenes vs lof genes: biotype uncertainty, nomenclature, artefacts, transposons.....

The GENCODE update trackhub: Ensembl and UCSC genome browsers

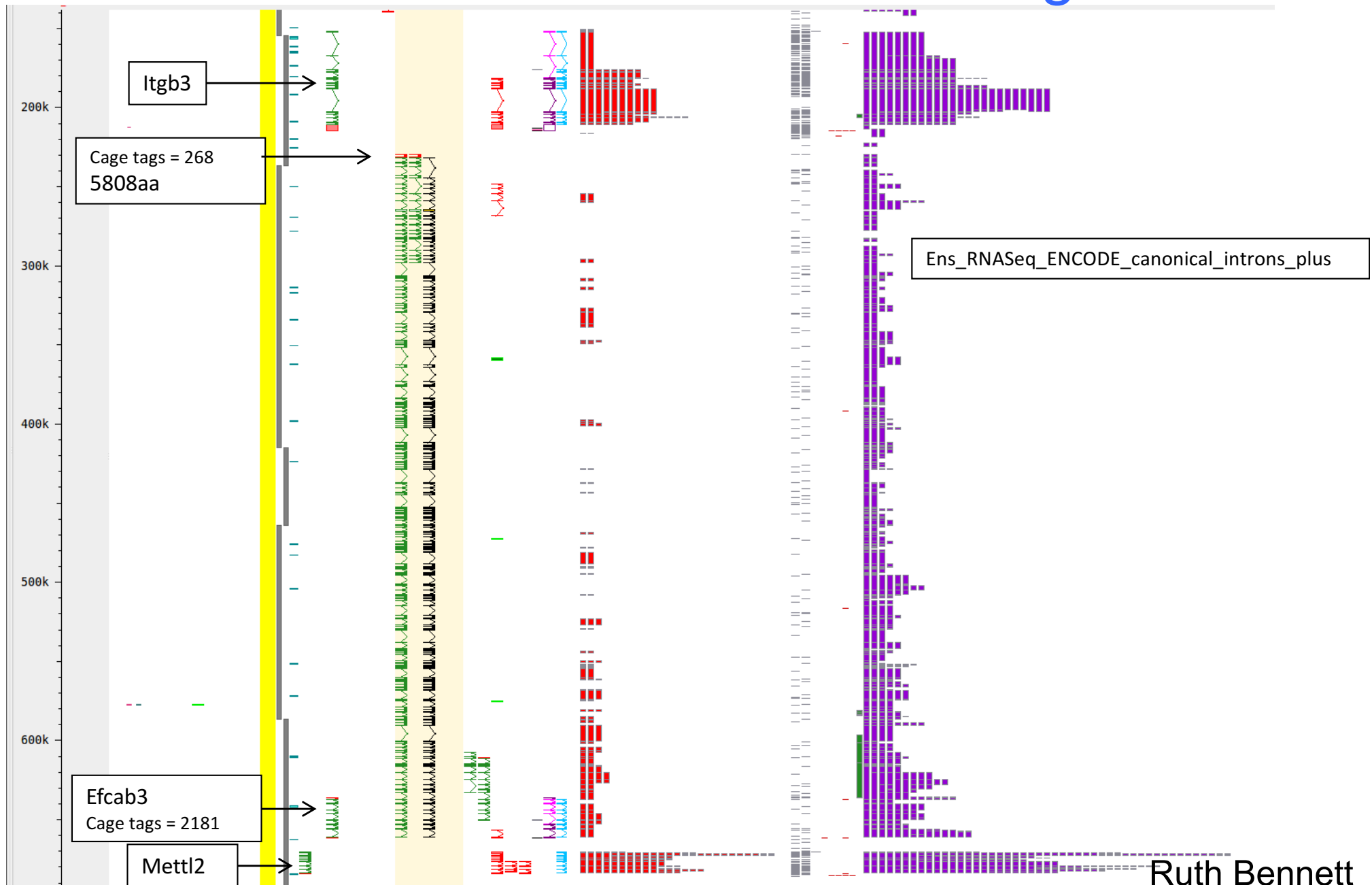


UCSC



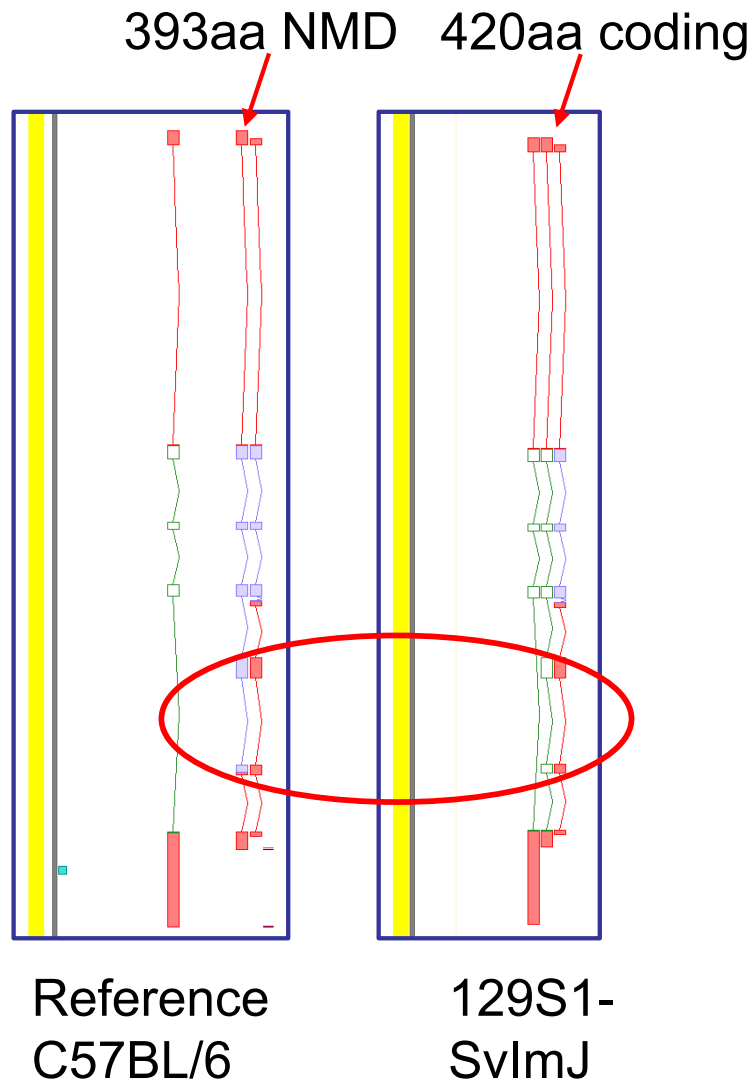
Update trackhub, in green, shows annotation updated in the last 24 hours. Annotation updated more than 24 hours ago will be shown in red.

Mouse strain annotation reveals new genes:



Ruth Bennett

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



>Reference longest NMD transcript 393 AA.

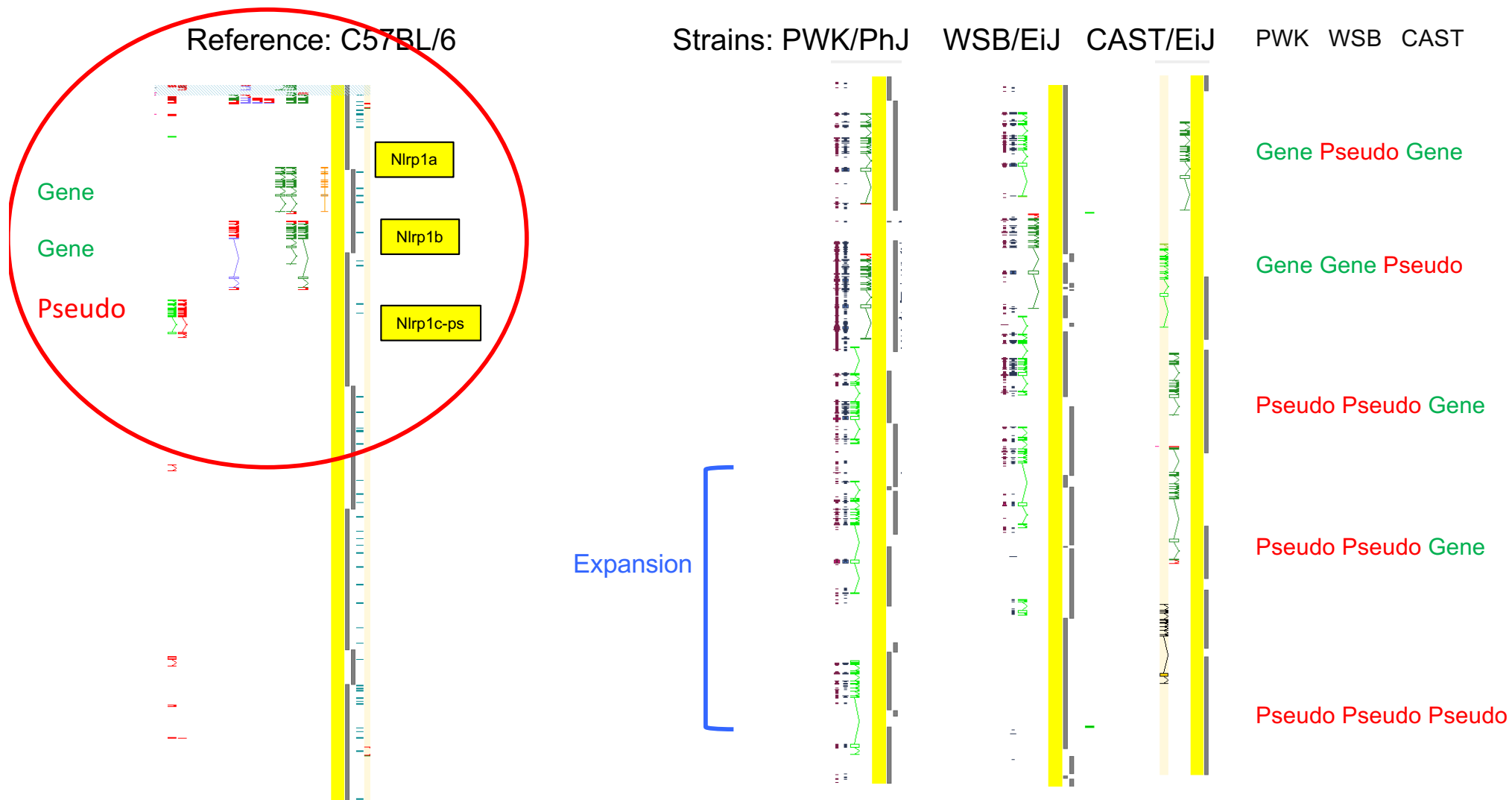
```
MVNEYKRIVLLTGLMGINDHDFRMVKSLLSKEKLNKMQDEYDRVKI
ADLMEDKFPKDAGVVQLIKLYKQIPGLGDIANKLKNEKAKAKRKGKG
KRKTAAKRQRQEEPSTSQPMSTTNEAEPESEGRSTPDTQVAQLSLPT
ASRRNQAIQISPTIASSSGQTSSRSSETLQSI IQSPETPTRSSSRIL
DPPVSPGTAYSSAQALGVLLATPAKRQRLKNVPKEPSEENGYQQGSK
KVMVLKVTEPFAYDMKGEKMFHATVATETEFFRVKVFDIVLKEKFIP
NKVLTISNYVGCNGFINIYSASSVSEVNDGEPMNIPLSLRKSANRTP
KINYLC SKRRGIFVNGVFTVCKKEERGY YICYEIGDDTGMMEVEVYG
RLTNIACNPGDKLRLML*
```

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

>129 Patch long coding transcript (equivalent to ref NMD) 420 AA.

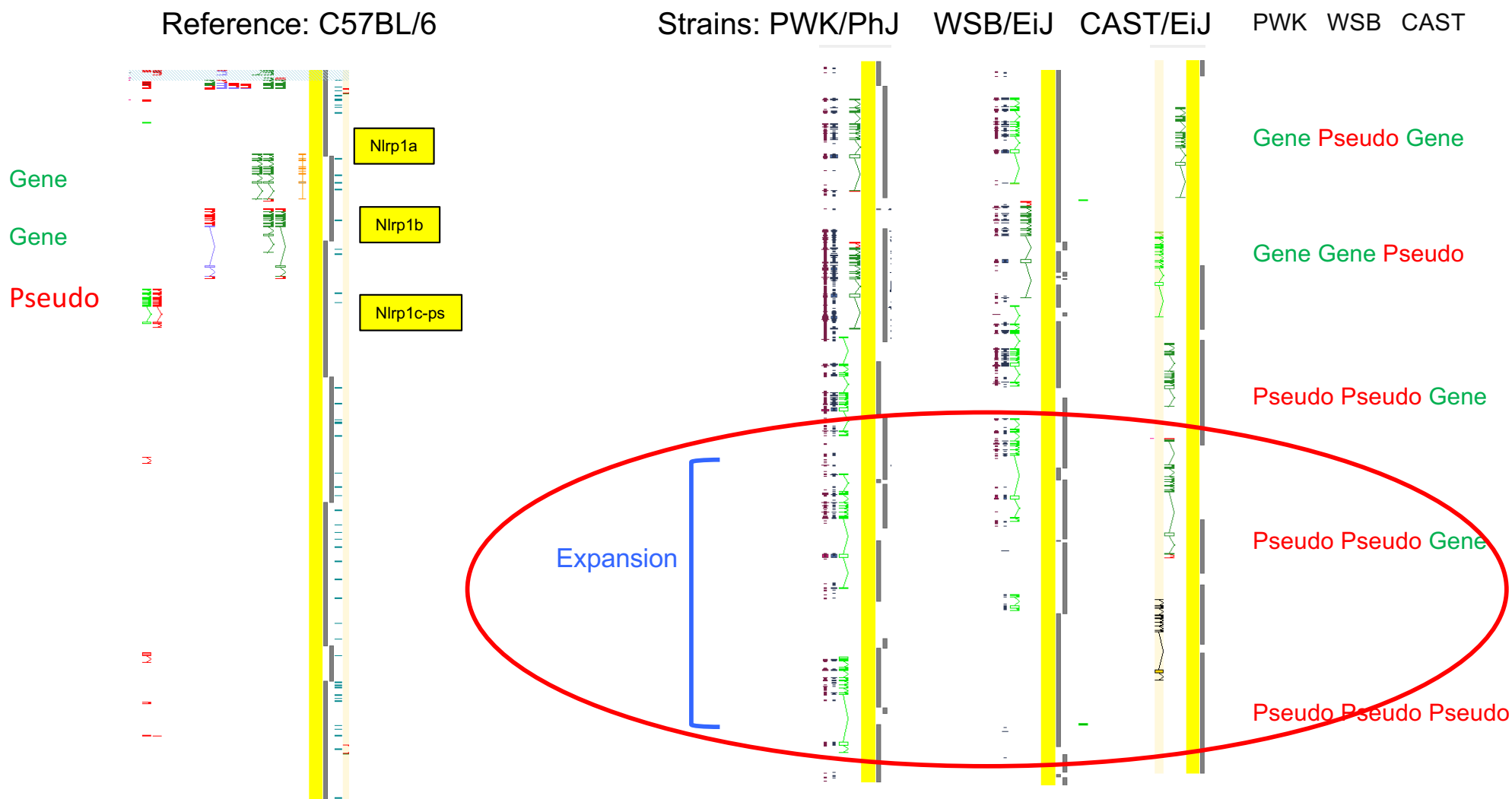
```
MVNEYKRIVLLTGLMGINDHDFRMVKSLLSKEKLNKMQDEYDRVKI
ADLMEDKFPKDAGVVQLIKLYKQIPGLGDIANKLKNEKAKAKRKGKG
KRKTAAKRQRQEEPSTSQPMSTTNEAEPESEGRSTPDTQVAQLSLPT
ASRRNQAIQISPTIASSSGQTSSRSSETLQSI IQSPETPTRSSSRIL
DPPVSPGTAYSSAQALGVLLATPAKRQRLKNVPKEPSEENGYQLGSK
KVMVLKVTEPFAYDMKGEKMFHATVATETEFFRVKVFDIVLKEKFIP
NKVLTISNYVGCNGFINIYSASSVSEVNDGEPMNIPLSLRKSANRTP
KINYLC SKRRGIFVNGVFTVCKKEERGY YICYEIGDDTGMMEVEVYG
RLTNIACNPGDKLRLICFELTPDEETAWLRSTHNSMQVIKARN*
```

Mouse strain annotation reveals strain specific expansions at the Nlrp locus



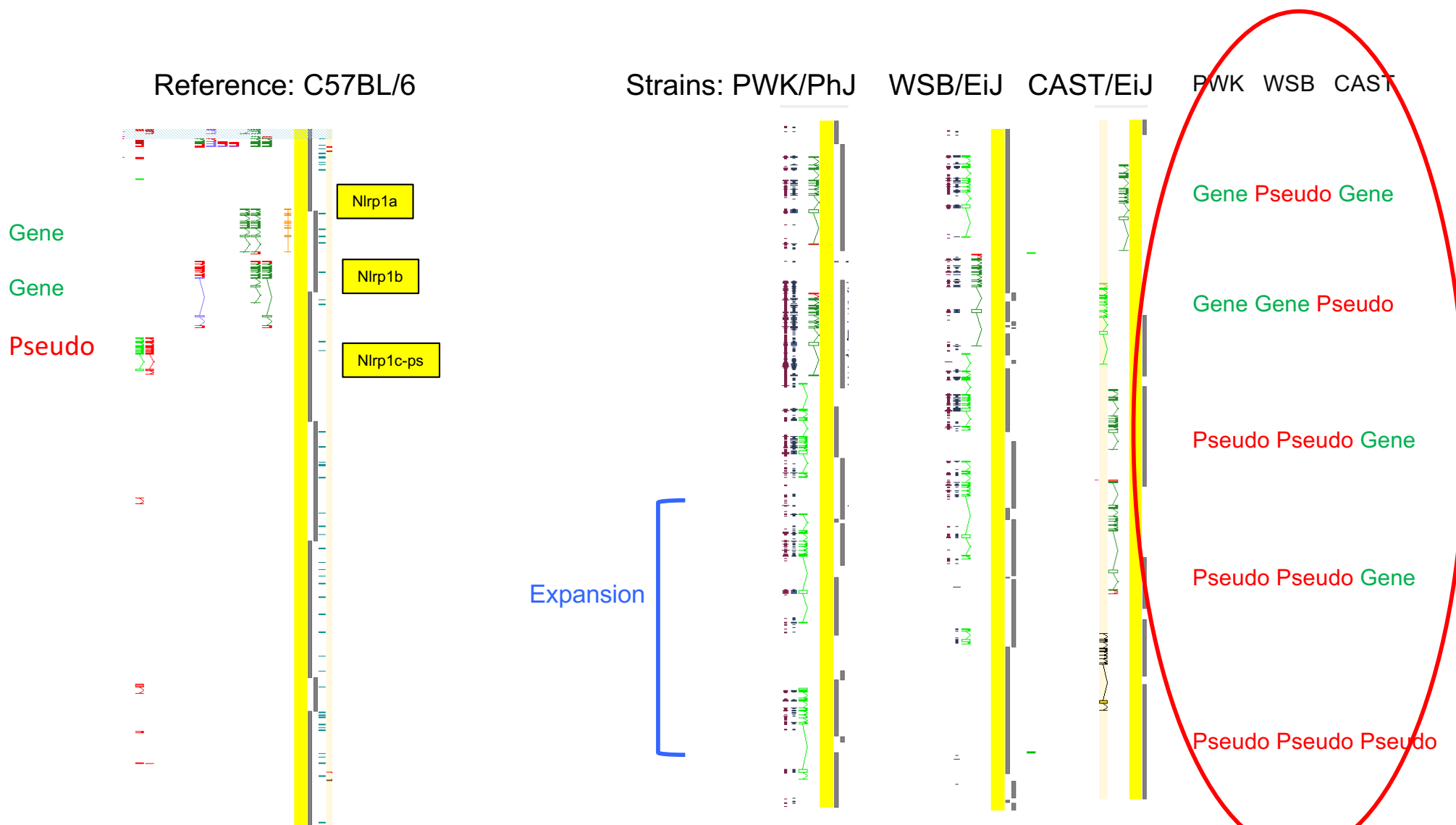
Ruth Bennett

Mouse strain annotation reveals strain specific expansions at the Nlrp locus



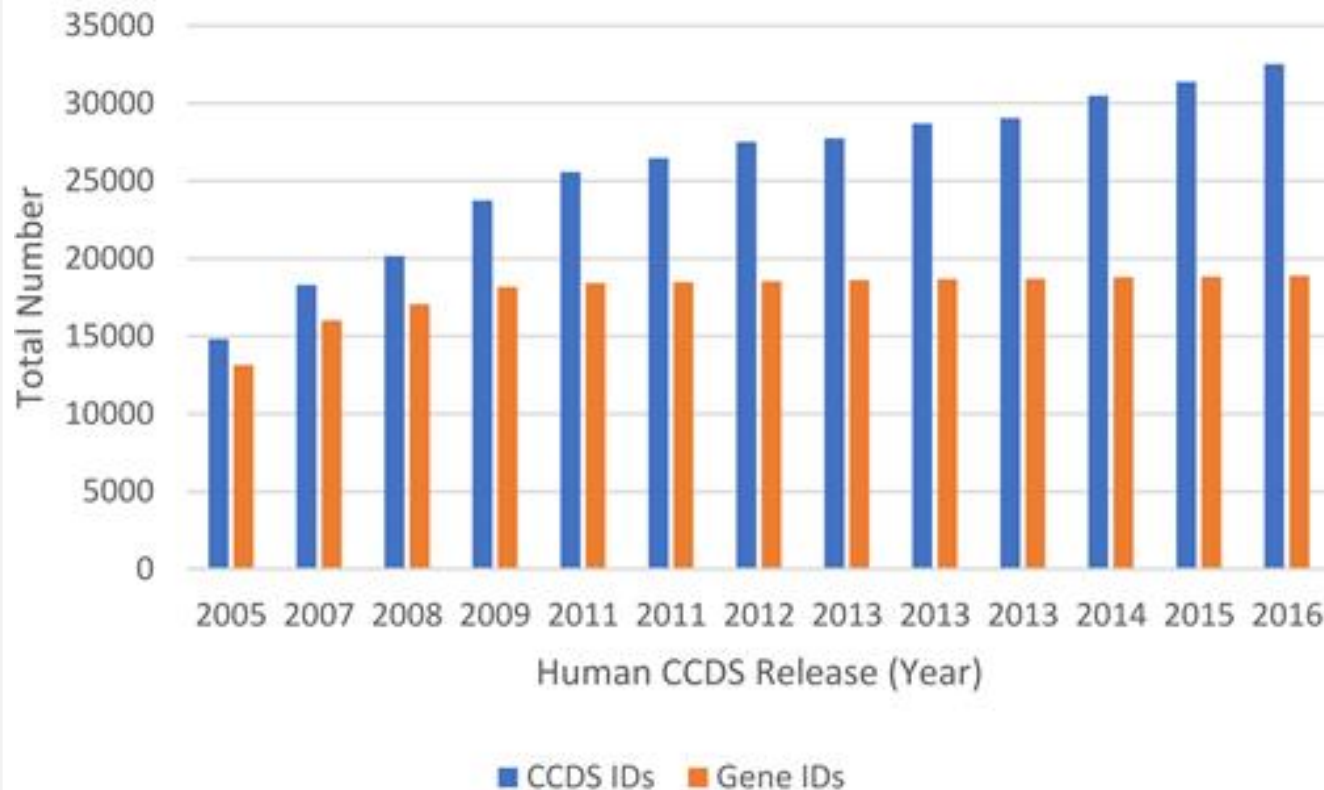
Ruth Bennett

Mouse strain annotation reveals strain specific expansions at the Nlrp locus



Ruth Bennett

Differences in manual annotation approach matter



CCDS Totals

Category	Count
CCDS IDs	32,524
Gene IDs	18,894
Sequence IDs	79,349

Status

Public	32,443
Reviewed, update pending	18
Reviewed, withdrawal pending	34
Under review, update	6
Under review, withdrawal	23

e! 19,766 protein coding

From: Consensus coding sequence (CCDS) database: a standardized set of human and mouse protein-coding regions supported by expert curation
Nucleic Acids Res. Published online November 06, 2017. doi:10.1093/nar/gkx1031

Differences in manual annotation approach matter

McCarthy *et al. Genome Medicine* 2014, **6**:26
<http://genomemedicine.com/content/6/3/26>



RESEARCH

Open Access

Choice of transcripts and software has a large effect on variant annotation

Davis J McCarthy^{1,2*}, Peter Humburg², Alexander Kanapin², Manuel A Rivas², Kyle Gaulton², The WGS500 Consortium, Jean-Baptiste Cazier³ and Peter Donnelly^{1,2}

Frankish *et al. BMC Genomics* 2015, **16**(Suppl 8):S2
<http://www.biomedcentral.com/1471-2164/16/S8/S2>



RESEARCH

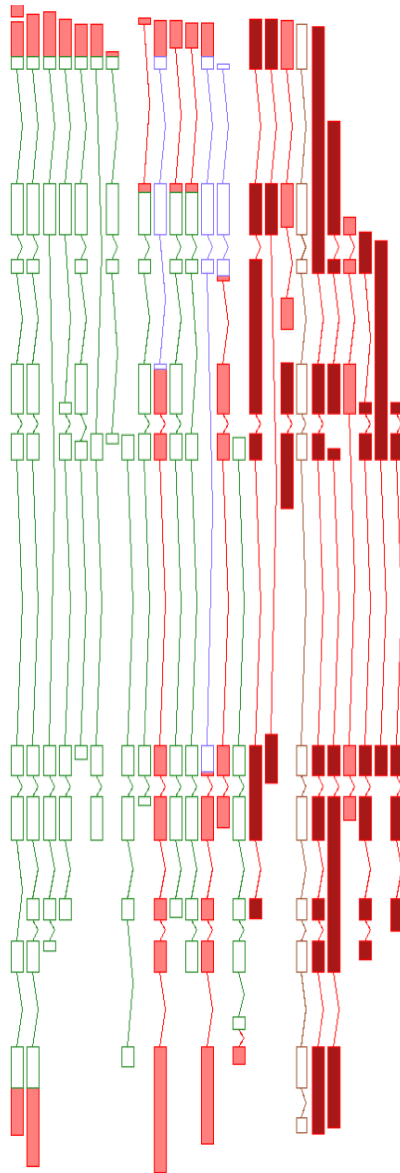
Open Access

Comparison of GENCODE and RefSeq gene annotation and the impact of reference geneset on variant effect prediction

Adam Frankish^{1*}, Barbara Uszczyńska², Graham RS Ritchie^{1,3}, Jose M Gonzalez¹, Dmitri Pervouchine^{2,4}, Robert Petryszak³, Jonathan M Mudge¹, Nuno Fonseca³, Alvis Brazma³, Roderic Guigo², Jennifer Harrow^{1*}



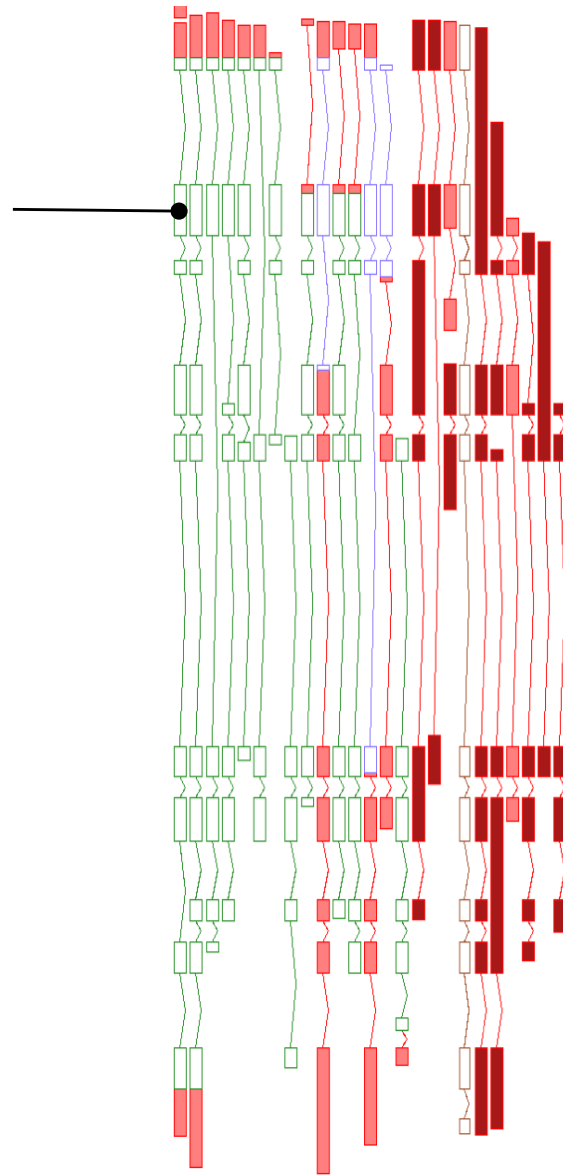
HAVANA annotation guidelines: protein-coding genes



TM7SF2

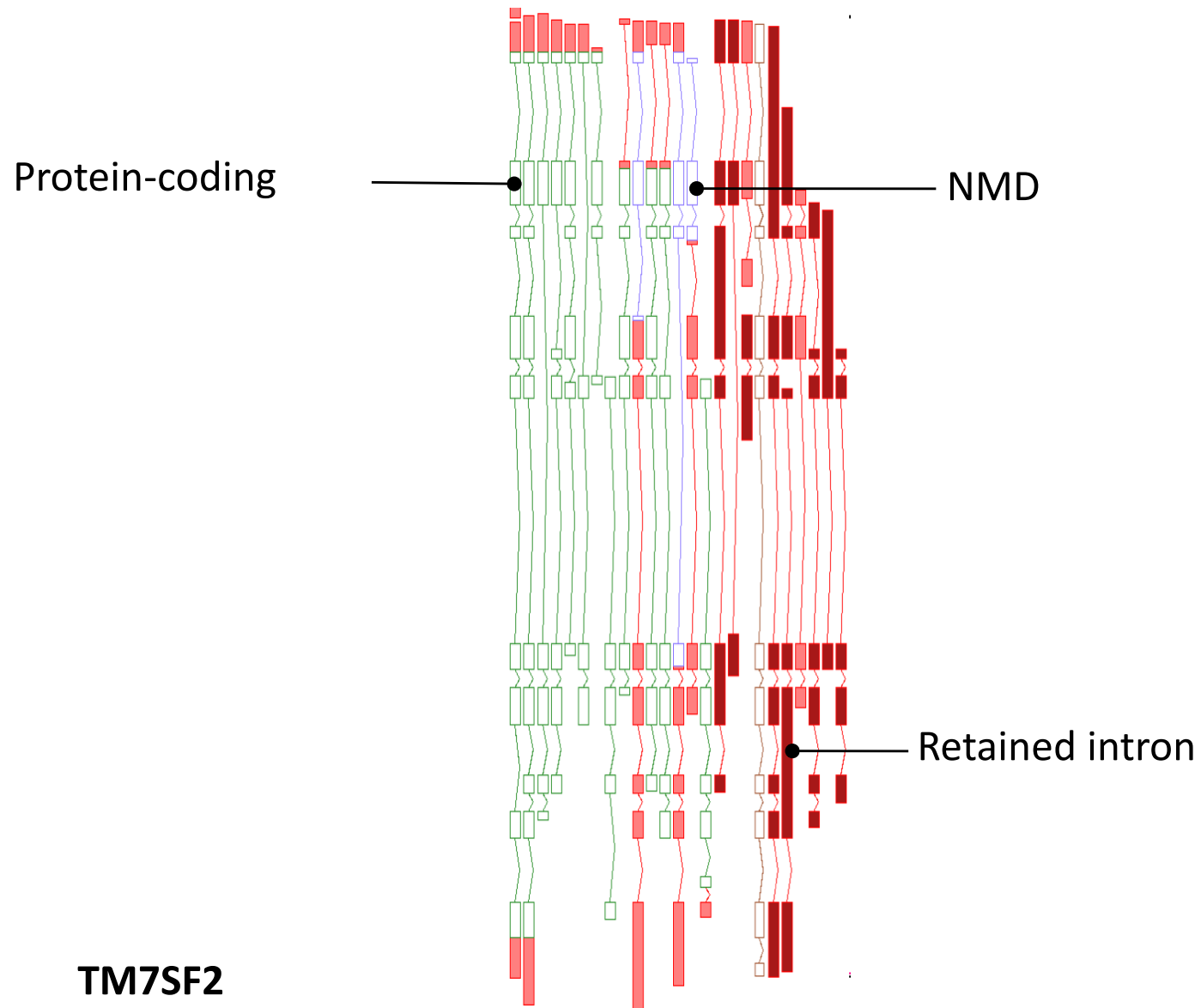
HAVANA annotation guidelines: protein-coding genes

Protein-coding

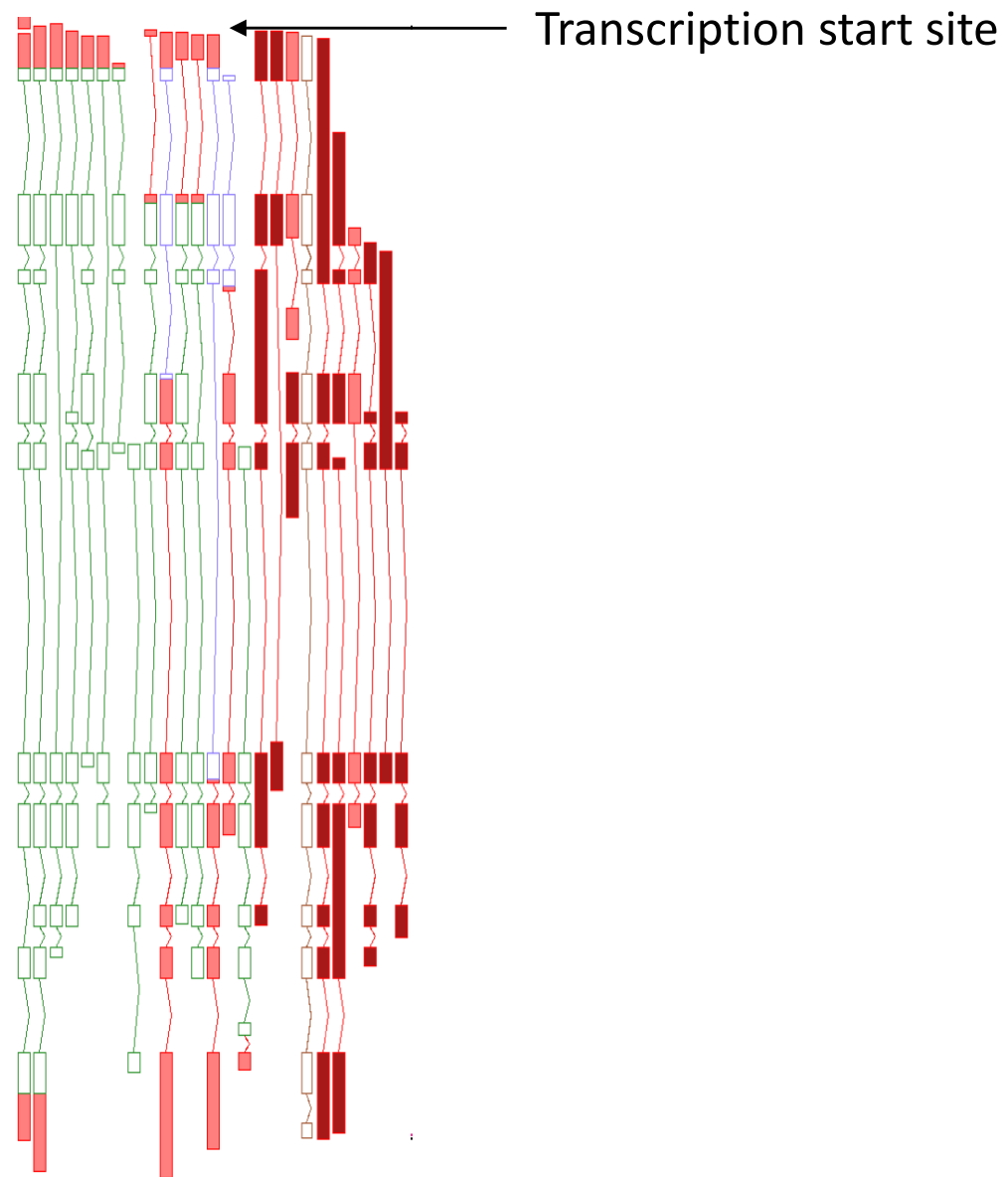


TM7SF2

HAVANA annotation guidelines: protein-coding genes

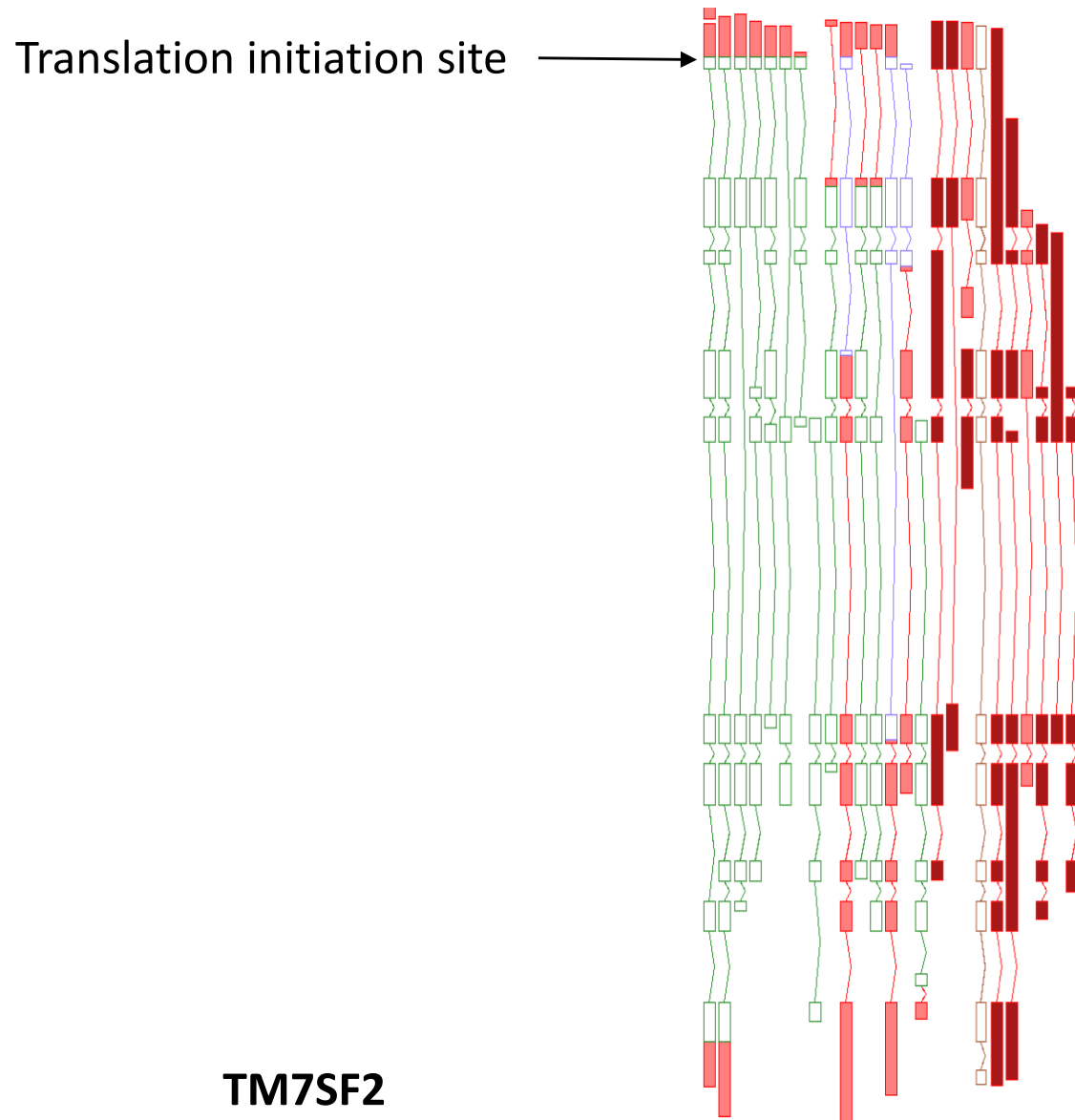


HAVANA annotation guidelines: protein-coding genes

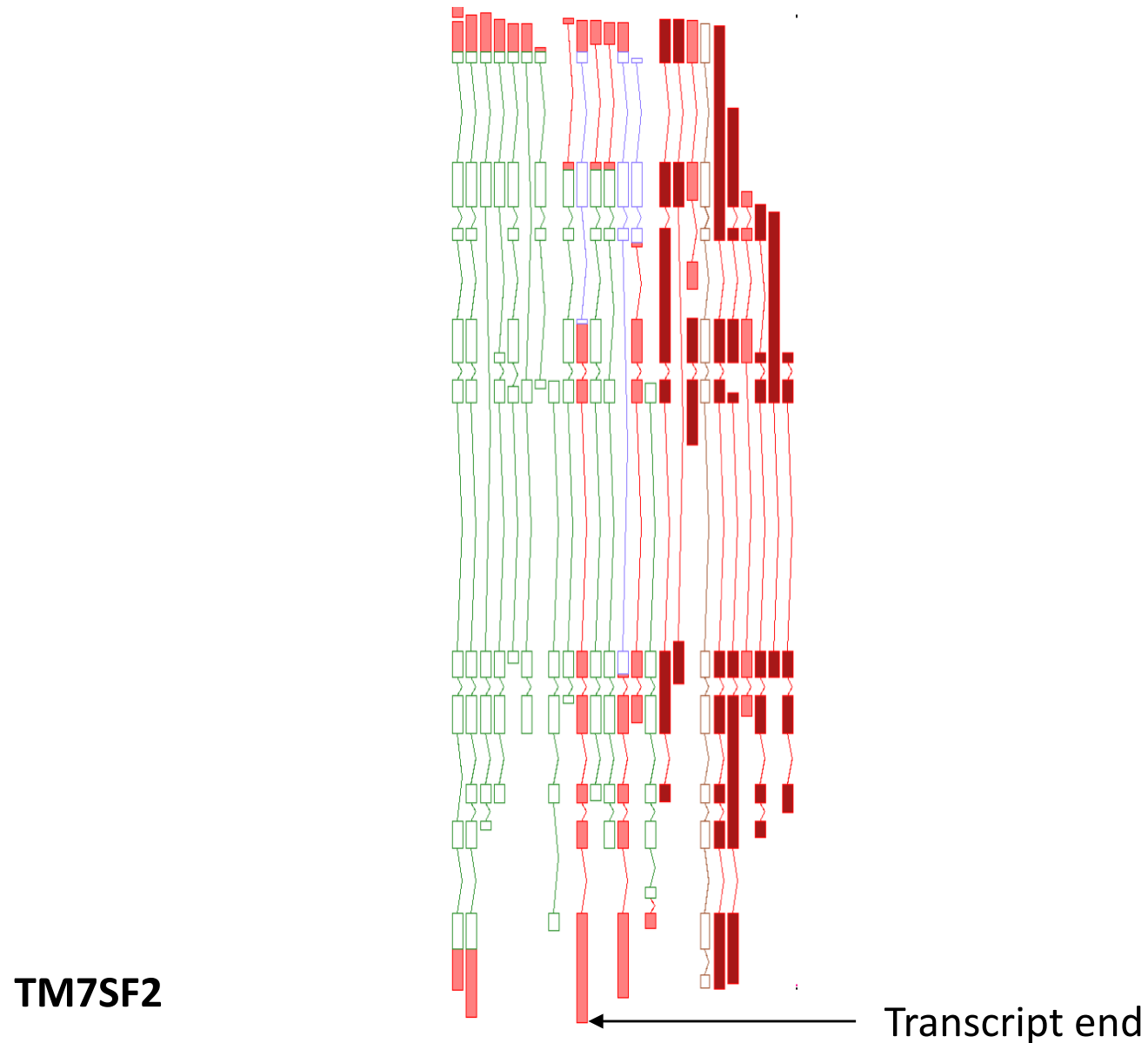


TM7SF2

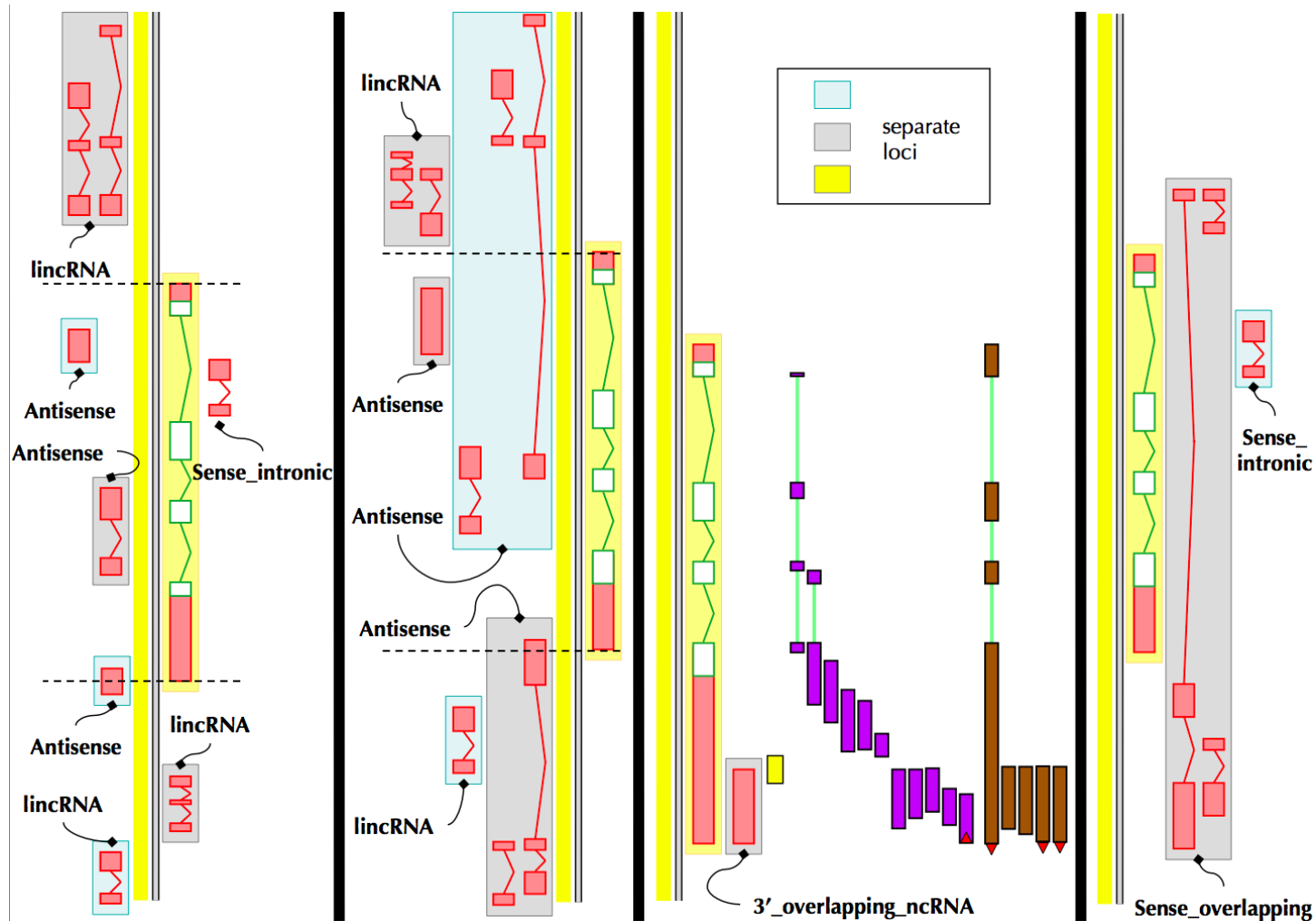
HAVANA annotation guidelines: protein-coding genes



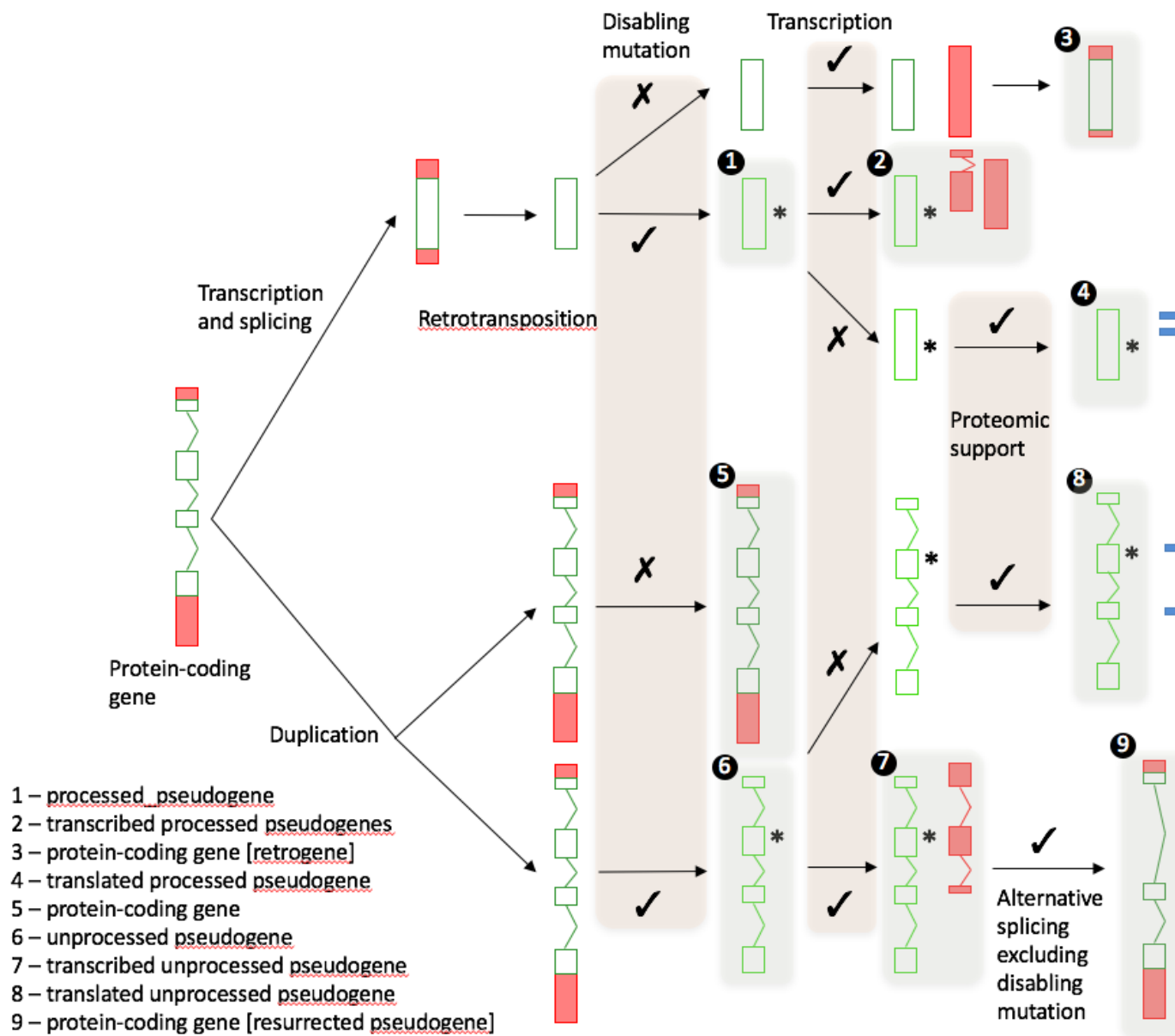
HAVANA annotation guidelines: protein-coding genes



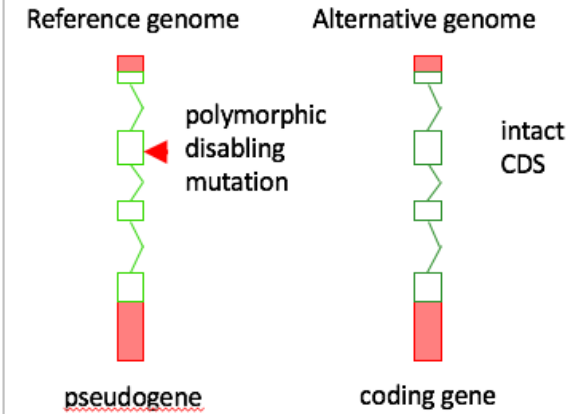
Improvements of lincRNA annotation: understanding functionality



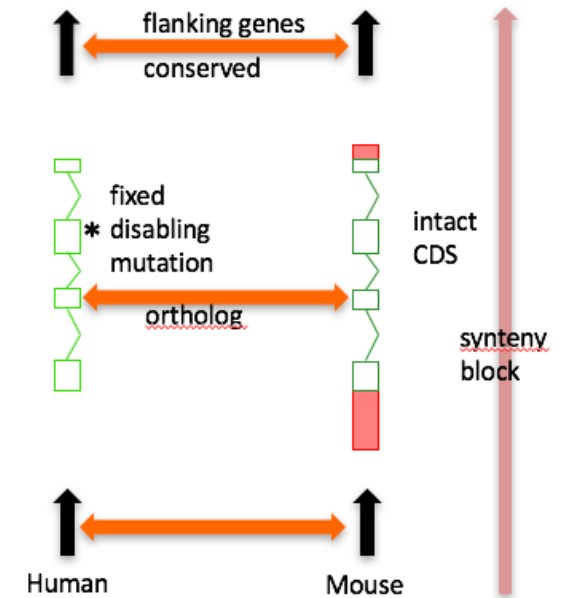
HAVANA annotation guidelines: pseudogenes



Polymorphic_pseudogene

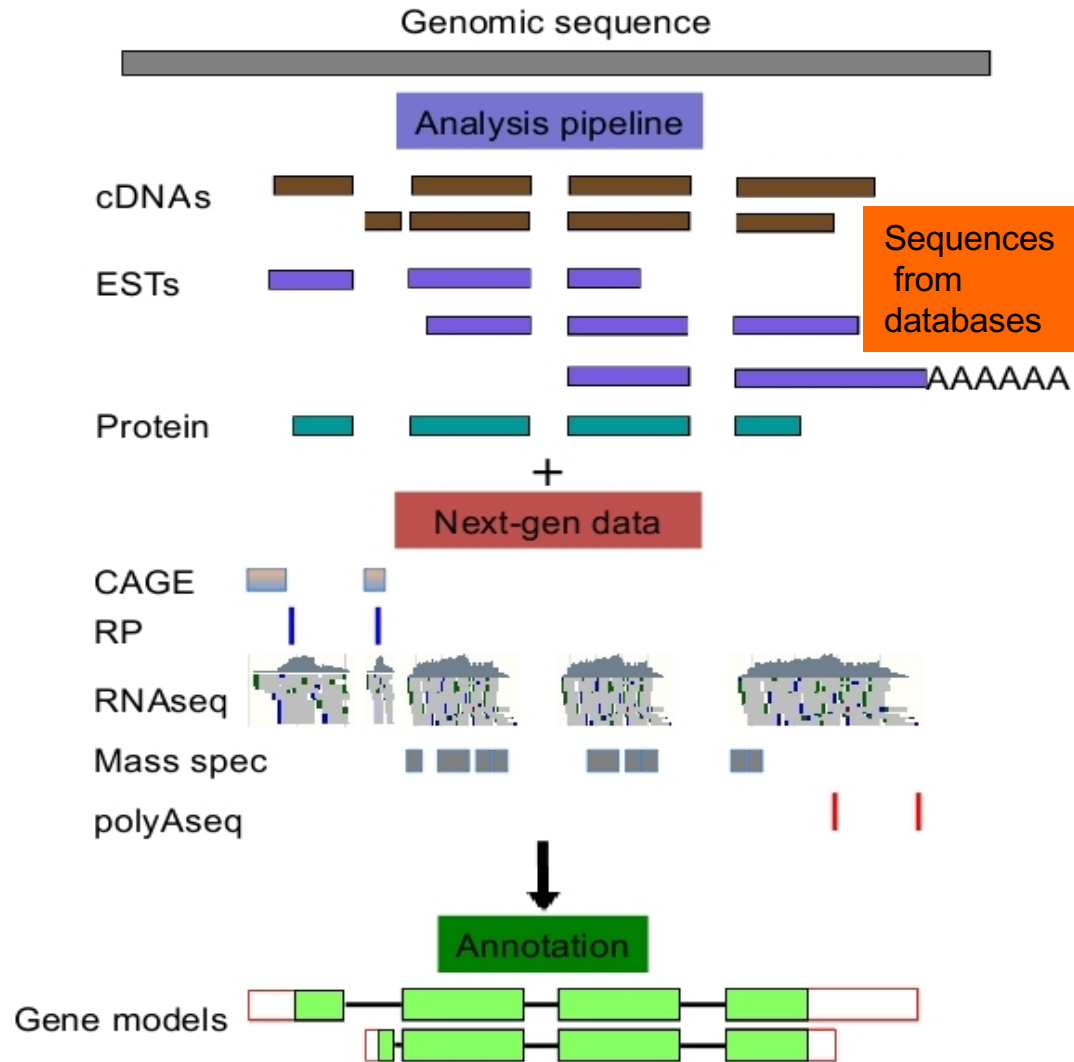


Unitary_pseudogene



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

GENCODE Biotypes

Version 27 (January 2017 freeze, GRCh38) - Ensembl 90, 91

General stats

Total No of Genes	58288
Protein-coding genes	19836
Long non-coding RNA genes	15778
Small non-coding RNA genes	7569
Pseudogenes	14694
- processed pseudogenes:	10704
- unprocessed pseudogenes:	3469
- unitary pseudogenes:	206
- polymorphic pseudogenes:	63
- pseudogenes:	18
Immunoglobulin/T-cell receptor gene segments	
- protein coding segments:	410
- pseudogenes:	234
Total No of Transcripts	200401
Protein-coding transcripts	80930
- full length protein-coding:	55406
- partial length protein-coding:	25524
Nonsense mediated decay transcripts	14208
Long non-coding RNA loci transcripts	27908
Total No of distinct translations	60297
Genes that have more than one distinct translations	13580

biotype	genes	transcripts
3prime_overlapping_ncRNA	31	35
antisense_RNA	5521	11050
bidirectional_promoter_lncRNA	19	40
IG_C_gene	14	23
IG_C_pseudogene	9	9
IG_D_gene	37	37
IG_J_gene	18	18
IG_J_pseudogene	3	3
IG_pseudogene	1	1
IG_V_gene	144	144
IG_V_pseudogene	188	188
lincRNA	7499	13348
macro_lincRNA	1	1
miRNA	1881	1881
misc_RNA	2213	2227
Mt_rRNA	2	2
Mt_tRNA	22	22
non_coding	3	3
non_stop_decay	0	84
nonsense_mediated_decay	0	14208
polymorphic_pseudogene	63	89
processed_pseudogene	10240	10243
processed_transcript	544	28230
protein_coding	19836	80930
pseudogene	18	37
retained_intron	0	27239
ribozyme	8	8
rRNA	544	544
scaRNA	49	49
scRNA	1	1
sense_intronic	905	963
sense_overlapping	189	339
snoRNA	943	955
snRNA	1900	1900
sRNA	5	5
TEC	1066	1165
TR_C_gene	6	6
TR_D_gene	4	4
TR_J_gene	79	79
TR_J_pseudogene	4	4
TR_V_gene	108	108
TR_V_pseudogene	30	30
transcribed_processed_pseudogene	462	462
transcribed_unitary_pseudogene	111	113
transcribed_unprocessed_pseudogene	830	836
translated_processed_pseudogene	2	2
unitary_pseudogene	95	95
unprocessed_pseudogene	2639	2640

IWGSC RefSeq v1.0 Gene stats:

~4000 manually curated genes

107,886 high confidence genes

Duplicated genes (Inparalogs): 27% of the high confidence gene set
Lineage specific duplications

~162,000 low confidence genes: unsure if pseudogene or gene or genomic error

~300,000 pseudogenes

~4 million transposable elements (85% genome)

850 RNAseq sets available

Please speak to us about manual annotation

Guidelines can be found at:

[ftp://ftp.sanger.ac.uk/pub/project/havana/Guidelines/
Guidelines_March_2016.pdf](ftp://ftp.sanger.ac.uk/pub/project/havana/Guidelines/Guidelines_March_2016.pdf)

Please contact us at:

gene-annotation@ebi.ac.uk

GENCODE Acknowledgements

Ensembl-HAVANA:

Adam Frankish
If Barnes
Andrew Berry
Alex Bignell
Sarah Donaldson
Matt Hardy
Toby Hunt
Jane Loveland
Jonathan Mudge
Gaurab Mukherjee
Marie-Marthe Suner
Mark Thomas

TGMI:

Joannella Morales
Ruth Bennett
Claire Davidson
Mike Kay

Annotrack/GENCODE:

Jose Manuel Gonzalez

Ensembl:

Paul Flicek
Bronwen Aken
Fiona Cunningham
Thibaut Hourlier
Carlos García Girón
Fergal Martin

GENCODE Consortium

Roderic Guigo, CRG

Julien Legarde

Barbara Uszczyński

Rory Johnson

Manolis Kellis, MIT

Irwin Jungreis

Michael Tress, CNIO

Alex Reymond, UNIL

Anne-Maude Ferreira

Mark Gerstein, Yale

Cristina Sisu

Fabio Navara

Benedict Paten, UCSC

Mark Diekhans

Tim Hubbard, KCL

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

Ensembl Acknowledgements

The Entire Ensembl Team

Daniel R. Zerbino¹, Premanand Achuthan¹, Wasiu Akanni¹, M. Ridwan Amode¹, Daniel Barrell^{1,2}, Jyothish Bhai¹, Konstantinos Billis¹, Carla Cummins¹, Astrid Gall¹, Carlos García Giroñ¹, Laurent Gil¹, Leo Gordon¹, Leanne Haggerty¹, Erin Haskell¹, Thibaut Hourlier¹, Osagie G. Izuogu¹, Sophie H. Janacek¹, Thomas Juettemann¹, Jimmy Kiang To¹, Matthew R. Laird¹, Ilias Lavidas¹, Zhicheng Liu¹, Jane E. Loveland¹, Thomas Maurel¹, William McLaren¹, Benjamin Moore¹, Jonathan Mudge¹, Daniel N. Murphy¹, Victoria Newman¹, Michael Nuhn¹, Denye Ogeh¹, Chuang Kee Ong¹, Anne Parker¹, Mateus Patricio¹, Harpreet Singh Riat¹, Helen Schuilenburg¹, Dan Sheppard¹, Helen Sparrow¹, Kieron Taylor¹, Anja Thormann¹, Alessandro Vullo¹, Brandon Walts¹, Amonida Zadissa¹, Adam Frankish¹, Sarah E. Hunt¹, Myrto Kostadima¹, Nicholas Langridge¹, Fergal J. Martin¹, Matthieu Muffato¹, Emily Perry¹, Magali Ruffier¹, Dan M. Staines¹, Stephen J. Trevanion¹, Bronwen L. Aken¹, Fiona Cunningham¹, Andrew Yates¹ and Paul Flicek^{1,3}

¹European Molecular Biology Laboratory, European Bioinformatics Institute, Wellcome Genome Campus, Hinxton, Cambridge CB10 1SD, UK, ²Eagle Genomics Ltd., Wellcome Genome Campus, Hinxton, Cambridge CB10 1DR, UK and ³Wellcome Trust Sanger Institute, Wellcome Genome Campus, Hinxton, Cambridge CB10 1SA, UK

Funding

