

Ensembl Gene Annotation (*e!98*)

Pig (*Sus scrofa*)

Assembly: Sscrofa11.1, GCA_000003025.6

Table of Contents

<u>SECTION 1: GENOME PREPARATION</u>	<u>4</u>
REPEAT FINDING	4
LOW COMPLEXITY FEATURES, AB INITIO PREDICTIONS AND BLAST ANALYSES	4
<u>SECTION 2: PROTEIN-CODING MODEL GENERATION</u>	<u>5</u>
PROJECTION MAPPING PIPELINE	5
PROTEIN-TO-GENOME PIPELINE	6
RNA-SEQ PIPELINE	6
<u>SECTION 3: FILTERING THE PROTEIN-CODING MODELS</u>	<u>8</u>
PRIORITISING MODELS AT EACH LOCUS	8
ADDITION OF UTR TO CODING MODELS	8
GENERATING MULTI-TRANSCRIPT GENES	9
PSEUDOGENES	9
IMMUNOGLOBULIN AND T-CELL RECEPTOR GENES	9
<u>SECTION 4: CREATING THE FINAL GENE SET</u>	<u>10</u>
SMALL ncRNAs	10
CROSS-REFERENCING	10
STABLE IDENTIFIERS	10
<u>SECTION 5: FINAL GENE SET SUMMARY</u>	<u>11</u>
<u>SECTION 6: APPENDIX - FURTHER INFORMATION</u>	<u>12</u>

ASSEMBLY INFORMATION	13
STATISTICS OF INTEREST	13
LAYERS IN DETAIL	14
MORE INFORMATION	16
<u>REFERENCES</u>	<u>17</u>

This document describes the annotation process of an assembly. The first stage is Assembly Loading where databases are prepared and the assembly loaded into the database.

Section 1: Genome Preparation

The genome phase of the Ensembl gene annotation pipeline (Fig. 1) involves loading an assembly into the Ensembl core database schema and then running a series of analyses on the loaded assembly to identify an initial set of genomic features.

The most important aspect of this phase is identifying repeat features (primarily through RepeatMasker) as soft masking of the genome is used extensively later in the annotation process.

Repeat Finding

After the genomic sequence has been loaded into a database, it is screened for sequence patterns including repeats using RepeatMasker [1] (version 4.0.5 with parameters, *-nolow-engine "crossmatch"*), Dust [2] and TRF [3]. For the pig annotation, the Repbase mammals library was used with RepeatMasker. In addition to the Repbase library, a custom repeat library was used with RepeatMasker. This custom library was created using RepeatModeler [1].

Low complexity features, ab initio predictions and BLAST analyses

Transcription start sites are predicted using Eponine-scan [4]. CpG islands longer than 400 bases and tRNAs are also predicted. The results of Eponine-scan, CpG, and tRNAscan [5] are for display purposes only; they are not used in the gene annotation process. Genscan [6] is run across repeat-masked sequence to identify ab initio gene predictions. The results of the Genscan analyses are also used as input for UniProt [7], UniGene [8] and Vertebrate RNA alignments by NCBI-BLAST [9]. Passing only Genscan results to BLAST is an effective way of reducing the search space and therefore the computational resources required.

Genscan predictions are for display purposes only and are not used in the model generation phase.

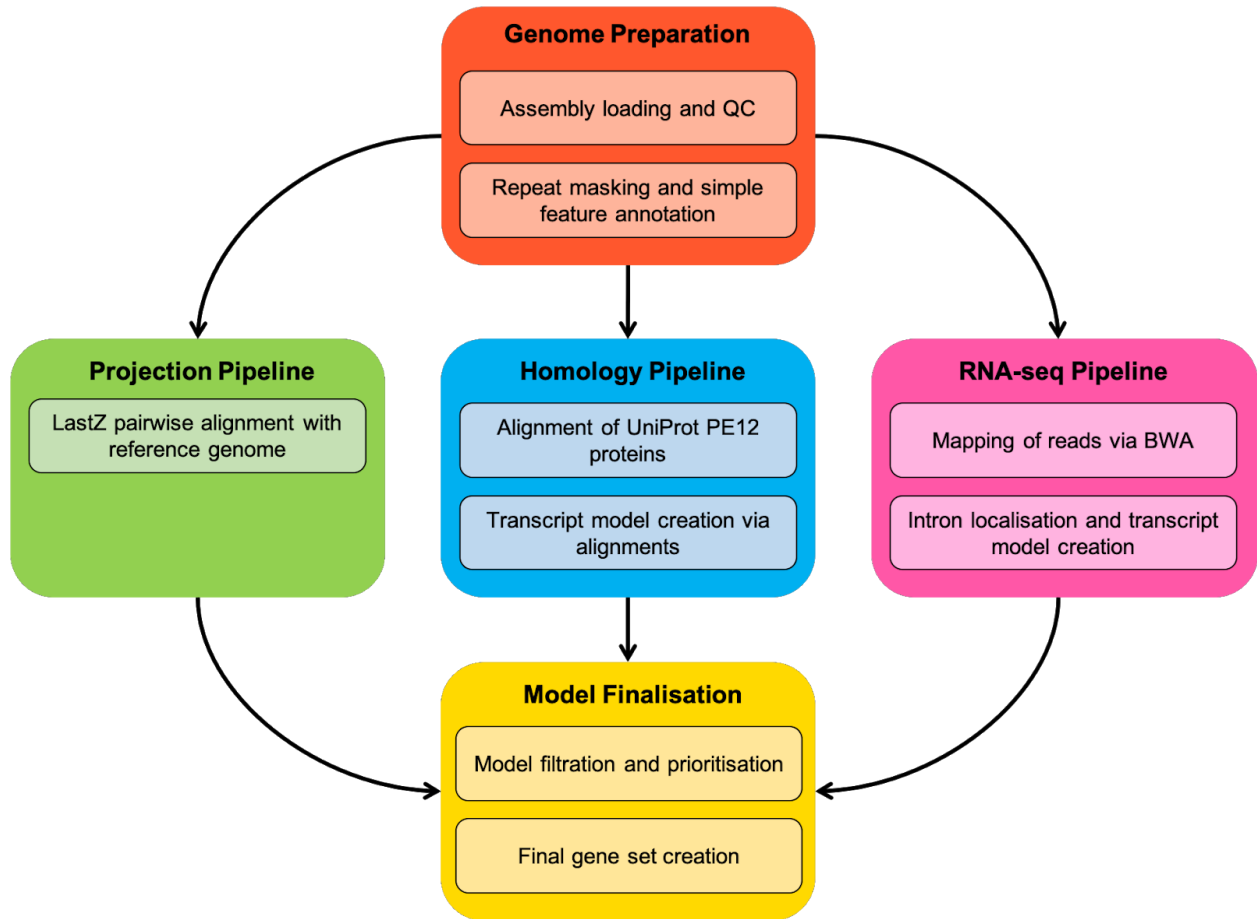


Fig. 1: Flowchart of the protein-coding annotation pipeline. Small ncRNAs, Ig genes, TR genes, and pseudogenes are computed using separate pipelines.

Section 2: Protein-Coding Model Generation

Various sources of transcript and protein data are investigated and used to generate gene models using a variety of techniques. The data and techniques employed to generate models

are outlined here. The numbers of gene models generated are described in gene summary.

Species specific cDNA and protein alignments

cDNAs are downloaded from ENA (www.ebi.ac.uk/ena) and RefSeq [10], and aligned to the genome using Exonerate [11]. Only known mRNAs are used (NMs). The cDNAs can be used to add UTR to the protein coding transcript models if they have a matching set of introns. Proteins are downloaded from UniProt and filtered based on protein existence (PE) at protein level and transcript level. The proteins are aligned to the genome using PMATCH to reduce the search space, then with genewise, which is a splice-aware aligner, to generate spliced models.

Projection mapping pipeline

For all species we generate a whole genome alignment against a suitable reference assembly using LastZ [12]. Syntenic regions identified using this alignment are then used to map protein-coding annotation from the most recent closely-related species or the Ensembl/GENCODE [13] gene set. For the pig annotation, we used the human assembly GRCh38.p12 as a reference to map protein-coding annotation. For each protein-coding gene in human, we projected the coding exons within the canonical transcript to pig. In case of exonic overlap on the projected sequence, the longest exon took precedence. If the mapping did not succeed, we selected the next successful projection of the transcript having the longest translation.

Protein-to-genome pipeline

Protein sequences were downloaded from UniProt and aligned to the genome in a splice aware manner using GenBlast [12]. The set of proteins aligned to the genome is a subset of UniProt proteins used to provide a broad, targeted coverage of the pig proteome. The set consists of the following:

- Self SwissProt/TrEMBL PE 1 & 2
- Human SwissProt/TrEMBL PE 1 & 2
- Mouse SwissProt/TrEMBL PE 1 & 2
- Other mammals SwissProt/TrEMBL PE 1 & 2
- Vertebrates SwissProt/TrEMBL PE 1 & 2

Note: PE level = protein existence level

For the pig annotation, logical thresholds for coverage (50%), percent identity (30%) and e-value (E-1) were used for GenBlast with the exon repair option turned on. The top 10 transcript models built by GenBlast for each protein passing the cut-offs are kept.

RNA-seq pipeline

RNA-seq data was downloaded from ENA (<https://www.ebi.ac.uk/ena/>) and used in the annotation. A merged file containing reads from all tissues/samples was created. The merged data is less likely to suffer from model fragmentation due to read depth. The available reads were aligned to the genome using BWA [13], with a tolerance of 50% mismatch to allow for intron identification via split read alignment. Initial models generated from the BWA alignments were further refined using exonerate. Protein-coding models were identified from BLAST alignments of the longest ORF against the UniProt vertebrate PE 1 & 2 data set. create a gene track for each tissue/sample viewable in the Ensembl browser and accessible through the Ensembl API.

Immunoglobulin and T-cell Receptor genes

Translations of different human IG gene segments were downloaded from the IMGT database [14] and aligned to the genome using GenBlast. For the pig annotation, logical thresholds for coverage (80%), percent identity (70%) and e-value (E-1) were used for GenBlast with the exon

repair option turned on. The top 10 transcript models built by GenBlast for each protein passing the cut-offs are kept.

Selenocysteine proteins

We aligned known selenocysteine proteins against the genome using Exonerate, checking that the generated model had a selenocysteine in the same positions as the known protein. We only kept models with at least 90% coverage and 95% identity.

Iso-Seq Pipeline

PacBio Iso-Seqs are transcriptomic long reads sequenced at high coverage. We downloaded from ENA the consensus sequences with ids: SRR5012257, SRR5012258, SRR5012866, SRR5012867, SRR5012868, SRR5012869, SRR5060320, SRR5120058, SRR5120059, SRR5250920, SRR5250921, SRR5275317, SRR5275318, SRR5275320, SRR5275321 and SRR5275382.

Samples were aligned to the genome using Minimap2[20], and then overlapping models were collapsed. Protein-coding models are identified via BLAST alignment of the longest ORF against the UniProt vertebrate PE 1 & 2 data set. As a result, our PacBio pipeline built 618,679 models for pig.

Section 3: Filtering the Protein-Coding Models

The filtering phase decides the subset of protein-coding transcript models, generated from the model-building pipelines, that comprise the final protein-coding gene set. Models were filtered based on information such as: what pipeline was used to generate them, how closely related to the target species and how good the alignment coverage and percent identity are to the original data.

Prioritising models at each locus

The LayerAnnotation module was used to define a hierarchy of input data sets, from most preferred to least preferred. The output of this pipeline includes all transcript models from the highest ranked input set. Models from lower ranked input sets were included only if their exons do not overlap a model from an input set higher in the hierarchy. Note that models cannot exist in more than one layer. For UniProt proteins, models were separated into clades, to help selection during the layering process. Each UniProt protein was in one clade only, for example mammal proteins were present in the mammal clade and not in the vertebrate clade to avoid aligning proteins multiple times.

When selecting the model or models kept at each position, we prioritise based on the highest layer with available evidence. In general, the highest layers contain the set of evidence containing the most trustworthy evidence in terms of both alignment/mapping quality, and also in terms of relevance to the species being annotated. So, for example, when a primate is being annotated, well aligned evidence from either the species itself or other closely related vertebrates would be chosen over evidence from more distant species. Regardless of what species is being annotated, well-aligned human proteins are usually included in the top layer as human is the current most complete vertebrate annotation. For further details on the exact layering used please refer to section 6.

Addition of UTR to coding models

The set of coding models were extended into the untranslated regions (UTRs) using RNA-seq data and alignments of species-specific RefSeq cDNA sequences. The criteria for adding UTR from cDNA or RNA-seq alignments to protein models lacking UTR (such as the projection models or the protein-to-genome alignment models) is that the intron coordinates from the model missing UTR exactly match a subset of the coordinates from the UTR donor model.

Generating multi-transcript genes

The above steps generated a large set of potential transcript models, many of which overlap one another. Redundant transcript models were collapsed and the remaining unique set of transcript models were clustered into multi-transcript genes where each transcript in a gene has at least one coding exon that overlaps a coding exon from another transcript within the same gene.

Pseudogenes

Pseudogenes were annotated by looking for genes with evidence of frame-shifting or lying in repeat heavy regions. Single exon retrotransposed pseudogenes were identified by searching for a multi-exon equivalent elsewhere in the genome. Identified pseudogenes and processed pseudogenes are included in the core db, please check Final Gene set Summary (Fig. 2).

Immunoglobulin and T-cell Receptor genes

Translations of different human IG gene segments were downloaded from the IMGT database [14] and aligned to the genome using GenBlast. For the pig annotation, logical thresholds for coverage (80%), percent identity (80%) and e-value ($E-1$) were used for GenBlast with the exon repair option turned on. The top 10 transcript models built by GenBlast for each protein passing the cut-offs are kept.

Section 4: Creating the Final Gene Set

Small ncRNAs

Small non-coding (sncRNA) genes were added using annotations taken from RFAM [15] and miRBase [16]. For miRNAs, NCBI-BLAST was run for these sequences to identify homologs in the genome sequence and models were evaluated for expected stem-loop structures using RNAfold [17]. Additional machine learning based filters were applied to exclude predictions with sub-optimal alignments to the genome and non-conforming secondary structures. For other sncRNAs, models were built using the Infernal software suite [18].

Cross-referencing

Before public release the transcripts and translations were given external references (cross-references to external databases). Translations were searched for signatures of interest and labelled where appropriate.

Stable Identifiers

Stable identifiers were assigned to each gene, transcript, exon and translation. When annotating a species for the first time, these identifiers are auto-generated. In all subsequent annotations for a species, the stable identifiers are propagated based on comparison of the new gene set to the previous gene set.

Section 5: Final Gene Set Summary

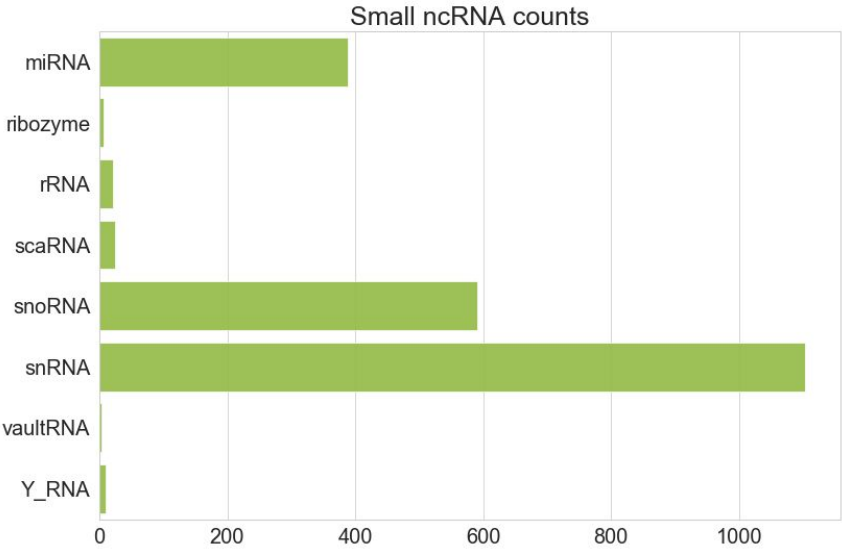
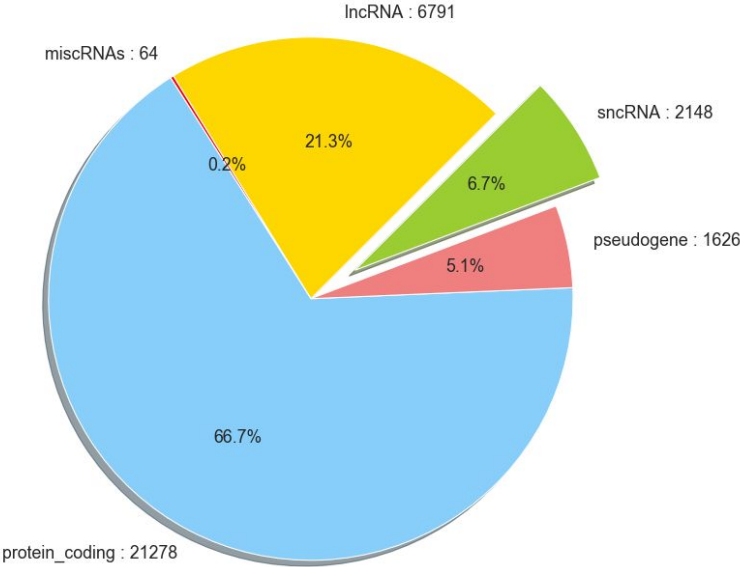


Figure 2: Counts of the major gene classes in Fig.

Section 6: Appendix - Further information

The Ensembl gene set is generated automatically, meaning that gene models are annotated using the Ensembl gene annotation pipeline. The main focus of this pipeline is to generate a conservative set of protein-coding gene models, although non-coding genes and pseudogenes may also be annotated. Every gene model produced by the Ensembl gene annotation pipeline is supported by biological sequence evidence (see the “Supporting evidence” link on the left-hand menu of a Gene page or Transcript page); *ab initio* models are not included in our gene set. *Ab initio* predictions and the full set of cDNA and EST alignments to the genome are available on our website.

The quality of a gene set is dependent on the quality of the genome assembly. Genome assembly can be assessed in a number of ways, including:

1. Coverage estimates
 - A higher coverage usually indicates a more complete assembly.
 - Using Sanger sequencing only, a coverage of at least 2x is preferred.
2. N50 of contigs and scaffolds
 - A longer N50 usually indicates a more complete genome assembly.
 - Bearing in mind that an average human gene may be 10-15 kb in length, contigs shorter than this length will be unlikely to hold full-length gene models.
3. Number of contigs and scaffolds
 - A lower number top level sequences usually indicates a more complete genome assembly.
4. Alignment of cDNAs and ESTs to the genome
 - A higher number of alignments, using stringent thresholds, usually indicates a more complete genome assembly.

Assembly Information

Table 1: Assembly information

Species	Common name	Assembly	Genbank accession	Date released
<i>Sus scrofa</i>	pig	Sscrofa11.1	GCA_000003025.6	2016-06

Statistics of Interest

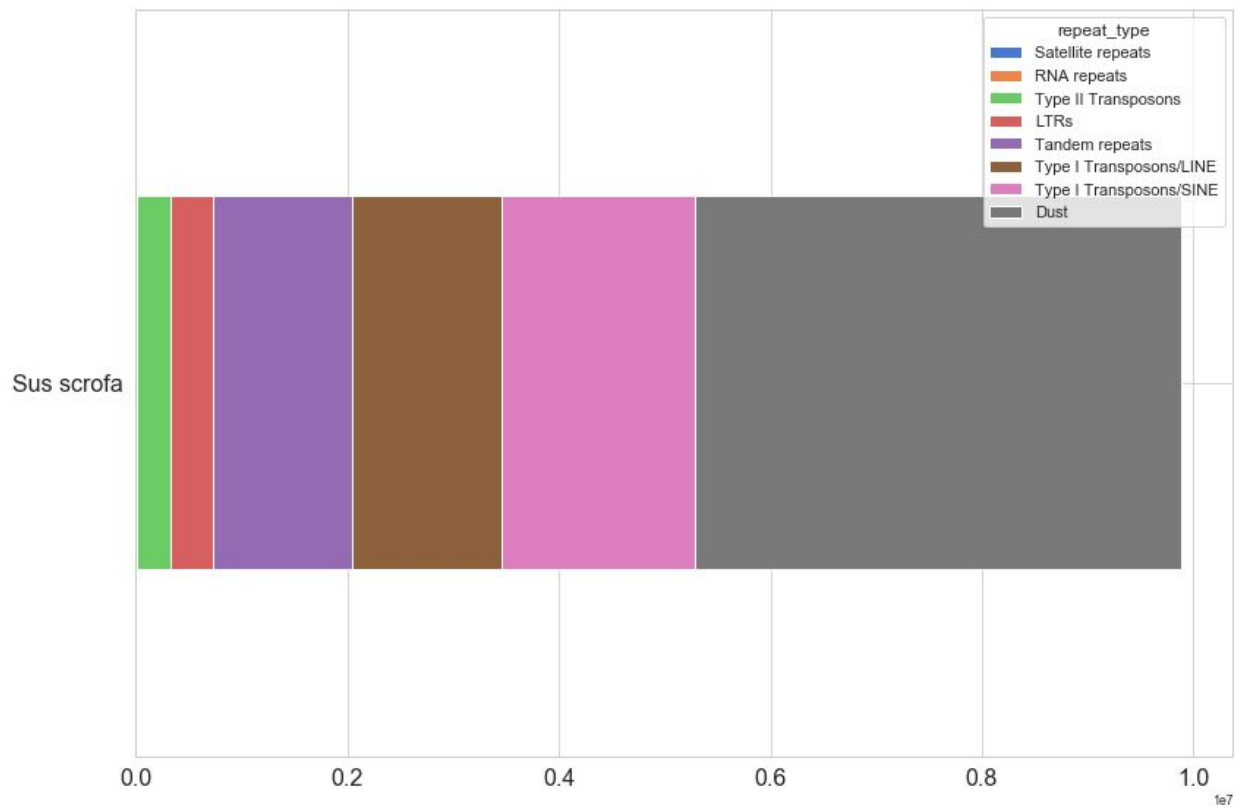


Figure 3: Number of repeat features identified in pig using "rebase_mammals" repeat library

Table 2: Transcript sizes

	min	median	max
biotype			
IG_C_gene	6548	6548	6548
IG_V_gene	272	321.5	8161
Mt_rRNA	959	1264.5	1570
Mt_tRNA	58	68	74
TR_J_gene	59	66.5	77
TR_V_gene	284	290	4303
Y_RNA	80	97	112
lncRNA	358	6822	577662
miRNA	53	81	116
misc_RNA	112	298	306
processed_pseudogene	302	1970	3752
protein_coding	203	29950.5	2268979
pseudogene	72	3409	607687
rRNA	118	118	1869
ribozyme	187	327	425
scaRNA	81	135	433
snRNA	59	106	220
snoRNA	53	110	328
vaultRNA	92	94	95

Table 3: Transcript sizes (spliced)

	min	median	max
biotype			
IG_C_gene	936	936	936
IG_V_gene	272	296	341
Mt_rRNA	959	1264.5	1570
Mt_tRNA	58	68	74
TR_J_gene	59	66.5	77
TR_V_gene	284	290	589
Y_RNA	80	97	112
lncRNA	66	2407	9221
miRNA	53	81	116
misc_RNA	112	298	306
processed_pseudogene	302	1970	3752
protein_coding	40	3069	24374
pseudogene	66	2988	15782
rRNA	118	118	1869
ribozyme	187	327	425
scaRNA	81	135	433
snRNA	59	106	220
snoRNA	53	110	328
vaultRNA	92	94	95

Layers in detail

Layer 1:

- Pig seleno-proteins
- Pig olfactory receptors with $\geq 90\%$ coverage and 97% identity
- All vertebrates seleno-proteins with full RNA-seq support
- IG and TR genes

Layer 2:

- Pig cDNAs models with $\geq 90\%$ coverage and 97% identity
- Pig IsoSeq models with protein support $\geq 80\%$ coverage and identity and full RNA-seq support
- RNA-seq models with $\geq 95\%$ coverage and identity
- Pig curated UniProt proteins from PE levels 1 & 2 with $\geq 80\%$ coverage and identity and full RNA-seq support
- Pig curated UniProt proteins from PE levels 3 with $\geq 95\%$ coverage and identity and full RNA-seq support
- All vertebrates curated UniProt proteins from PE levels 1 & 2 with $\geq 95\%$ coverage and identity and full RNA-seq support

Layer 3:

- RNA-seq models with $\geq 80\%$ coverage and identity

Layer 4:

- Pig curated UniProt proteins from PE levels 1 & 2 with $\geq 50\%$ coverage and identity
- Pig IsoSeq models with protein support $\geq 80\%$ coverage and identity

Layer 5:

- Pig curated UniProt proteins from PE levels 3 with $\geq 80\%$ coverage and identity
- All vertebrates curated UniProt proteins from PE level 1 & 2 with $\geq 80\%$ coverage and identity

Layer 6:

- RNA-seq models with $\geq 50\%$ coverage and identity
- Pig IsoSeq models with protein support $\geq 50\%$ coverage and identity
- Pig curated UniProt proteins from PE levels 3 with $\geq 50\%$ coverage and identity
- All vertebrates curated UniProt proteins from PE level 1 & 2 with $\geq 50\%$ coverage and identity

Layer 7:

- Pig UniProt proteins from PE levels 1 & 2 & 3 with $\geq 80\%$ coverage and identity and full RNA-seq support
- All vertebrates UniProt proteins from PE levels 1 & 2 with $\geq 80\%$ coverage and identity and full RNA-seq support

Layer 8:

- Pig UniProt proteins from PE levels 1 & 2 & 3 with $\geq 50\%$ coverage and identity and full RNA-seq support
- All vertebrates UniProt proteins from PE levels 1 & 2 with $\geq 50\%$ coverage and identity and full RNA-seq support
- Pig Isoseq models with protein support $\geq 50\%$ coverage and identity

which may have a retained intron

Layer 9:

- Pig UniProt proteins from PE levels 1 & 2 & 3 with $\geq 80\%$ coverage and identity
- All vertebrates UniProt proteins from PE levels 1 & 2 with $\geq 80\%$ coverage and identity

Layer 10:

- Pig UniProt proteins from PE levels 1 & 2 & 3 with $\geq 50\%$ coverage and identity
- All vertebrates UniProt proteins from PE levels 1 & 2 with $\geq 50\%$ coverage and identity

[More information](#)

More information on the Ensembl automatic gene annotation process can be found at:

- Publication : Aken B et al.: The Ensembl gene annotation system. Database 2016.
- Web: [Link to Ensembl gene annotation documentation](#)

References

1. Smit, A., R. Hubley, and P. Green, <http://www.repeatmasker.org/>. RepeatMasker Open, 1996. **3**: p. 1996-2004.
2. Kuzio, J., R. Tatusov, and D. Lipman, *Dust*. Unpublished but briefly described in: Morgulis A, Gertz EM, Schäffer AA, Agarwala R. A Fast and Symmetric DUST Implementation to Mask Low-Complexity DNA Sequences. *Journal of Computational Biology*, 2006. **13**(5): p. 1028-1040.
3. Benson, G., *Tandem repeats finder: a program to analyze DNA sequences*. *Nucleic acids research*, 1999. **27**(2): p. 573.
4. Down, T.A. and T.J. Hubbard, *Computational detection and location of transcription start sites in mammalian genomic DNA*. *Genome research*, 2002. **12**(3): p. 458-461.
5. Lowe, T.M. and S.R. Eddy, *tRNAscan-SE: a program for improved detection of transfer RNA genes in genomic sequence*. *Nucleic acids research*, 1997. **25**(5): p. 955-964.
6. Burge, C. and S. Karlin, *Prediction of complete gene structures in human genomic DNA*. *Journal of molecular biology*, 1997. **268**(1): p. 78-94.
7. Consortium, U., *UniProt: the universal protein knowledgebase*. *Nucleic acids research*, 2017. **45**(D1): p. D158-D169.
8. Pontius, J.U., L. Wagner, and G.D. Schuler, *21. UniGene: A unified view of the transcriptome*. *The NCBI Handbook*. Bethesda, MD: National Library of Medicine (US), NCBI, 2003.
9. Altschul, S.F., et al., *Basic local alignment search tool*. *Journal of molecular biology*, 1990. **215**(3): p. 403-410.
10. O'Leary, N.A., et al., *Reference sequence (RefSeq) database at NCBI: current status, taxonomic expansion, and functional annotation*. *Nucleic acids research*, 2015. **44**(D1): p. D733-D745.
11. Slater, G.S.C. and E. Birney, *Automated generation of heuristics for biological sequence comparison*. *BMC bioinformatics*, 2005. **6**(1): p. 31.

12. She, R., et al., *genBlastG: using BLAST searches to build homologous gene models*. *Bioinformatics*, 2011. **27**(15): p. 2141-2143.
13. Li, H. and R. Durbin, *Fast and accurate short read alignment with Burrows–Wheeler transform*. *Bioinformatics*, 2009. **25**(14): p. 1754-1760.
14. Lefranc, M.-P., et al., *IMGT®, the international ImMunoGeneTics information system® 25 years on*. *Nucleic acids research*, 2014. **43**(D1): p. D413-D422.
15. Griffiths-Jones, S., et al., *Rfam: an RNA family database*. *Nucleic acids research*, 2003. **31**(1): p. 439-441.
16. Griffiths-Jones, S., et al., *miRBase: microRNA sequences, targets and gene nomenclature*. *Nucleic acids research*, 2006. **34**(suppl_1): p. D140-D144.
17. R. Lorenz, S.H. Bernhart, C. Hoener zu Siederdisen, H. Tafer, C. Flamm, P.F. Stadler and I.L. Hofacker (2011), "ViennaRNA Package 2.0", *Algorithms for Molecular Biology*: 6:26
18. Nawrocki, E.P. and S.R. Eddy, *Infernal 1.1: 100-fold faster RNA homology searches*. *Bioinformatics*, 2013. **29**(22): p. 2933-2935.
19. Finn, R.D., et al., *The Pfam protein families database: towards a more sustainable future*. *Nucleic acids research*, 2016. **44**(D1): p. D279-D285.
20. Heng Li, *Minimap2: pairwise alignment for nucleotide sequences*, *Bioinformatics*, Volume 34, Issue 18, 15 September 2018, Pages 3094–3100, <https://doi.org/10.1093/bioinformatics/bty191>