# Ensembl gene annotation update (*e!90*)
# Sus scrofa, Sscrofa11.1

This document describes the annotation process of the high-coverage pig Sscrofa11.1 assembly, described in Figure 1. The first stage is Assembly Loading where databases are prepared and the assembly loaded into the database.
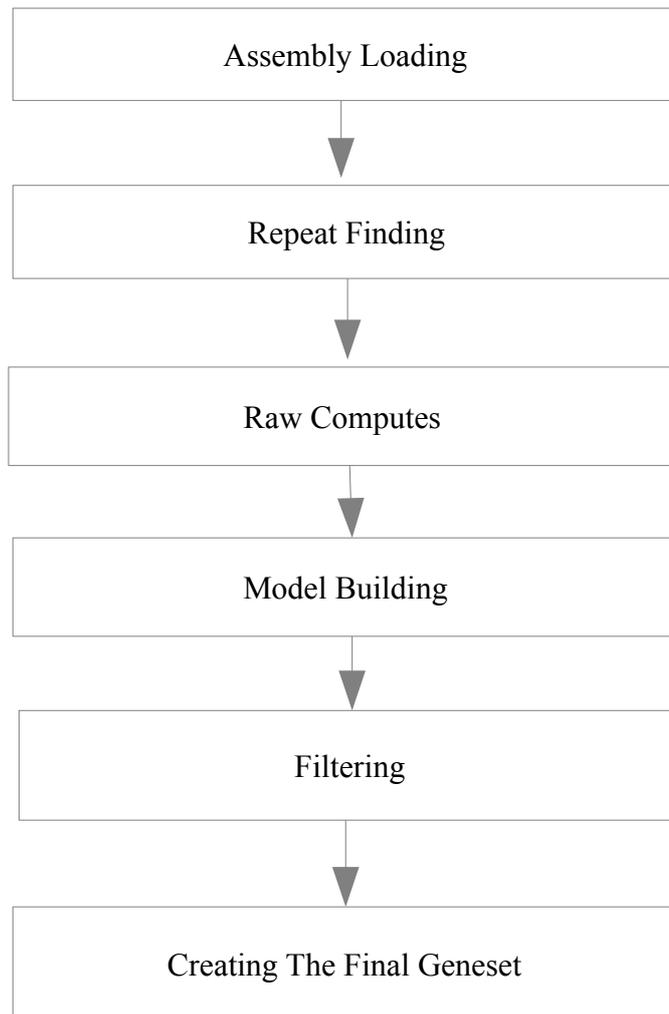


**Figure 1: The Gene Annotation Pipeline**

## Repeat Finding

After loading into a database the genomic sequence was screened for sequence patterns including repeats using RepeatMasker [1] (version 4.0.5 with parameters '-nolow -species "sus scrofa" -engine "crossmatch"', dustmasker [2] and TRF [3]. Both executions of RepeatMasker and dustmasker combined masked 45.04% of the assembly.

## Raw computes

Transcription start sites were predicted using Eponine–scan [4]. CpG islands [Micklem, G.] longer than 400 bases and tRNAs [5] were also predicted. The results of Eponine-scan, CpG, and tRNAscan are for display purposes only; they are not used in the gene annotation process.

Genscan [6] was run across repeat-masked sequence and the results were used as input for UniProt [7], UniGene [8] and Vertebrate RNA [9] alignments by BLAST+ [2]. Passing only Genscan results to BLAST is an effective way of reducing the search space and therefore the computational resources required. This resulted in 5,680,769 UniProt, 4,801,230 UniGene and 4,414,040 Vertebrate RNA sequences aligning to the genome.

## Model Generation

Various sources of transcript and protein data were 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 Table 1.

| Pipeline | Source | Number of Models |
|---|---|---|
| Species specific cDNAs | RefSeq, ENA | 45,589 |
| PacBio IsoSeq | Iowa State University | 326,217 |
| RNA-seq | Roslin Institute | 572,419 |
| Olfactory receptors | Human and mouse Ensembl release 89 | 1,212 |
| IG/TR genes | IMGT® | 1,803 |
| Protein-to-genome | Subset of UniProt vertebrate proteins | 509,769 |

**Table 1: Gene Model Generation Overview**

## cDNA Alignments

Pig cDNAs were downloaded from ENA and RefSeq, and aligned to the genome using Exonerate [10]. A minimal sequence length of 60bp was and a cut-off of 97% identity and 90% coverage were required for an alignment to be kept. The cDNAs are mainly used for display purposes, but can be used to add UTR to the protein coding transcript models if they have a matching set of introns.

| Species | Initial mRNA sequences | Sequences aligned |
|---------|------------------------|-------------------|
| Pig     | 45,571                 | 45,526            |

**Table 2: Species specific cDNAs aligned against Sscrofa11.1**

## PacBio IsoSeqs

PacBio IsoSeqs are transcriptomic long reads sequenced at a high coverage to allow correction of the technology. We downloaded the consensus sequences from SRA representing (PRJNA351265) nine tissue types; brain, diaphragm, hypothalamus, liver, longissimus muscle, pituitary, small intestine, spleen, thymus, after correction using Illumina short reads from the same tissue type. The sequences were aligned to the genome using Exonerate using a cut-off of 95% identity and 90% coverage.

All sets had 3' capping and were used for adding UTRs to homology-based protein-coding models. Both sets were used as lincRNA candidate for our lincRNA prediction pipeline.

Furthermore we collapsed the models created to produce a non-redundant set. We ran BLAST+ against a sub-set of UniProt to determine the coding potential of the set. We also checked the splice junctions with the RNA-seq data set. The models with a high coding potential and with full RNA-seq support were used as input for the gene model generation.

| Tissue sample | Initial IsoSeq sequences | Sequences aligned |
|---------------|--------------------------|-------------------|
| Liver         | 588,957                  | 491,796           |
| Thymus        | 567,700                  | 374,515           |

| | | |
|---|---|---|
| Hypothalamus | 414,021 | 256,930 |
| Brain | 398,629 | 354,494 |
| Longissimus muscle | 410,420 | 361,864 |
| Diaphragm | 459,911 | 391,813 |
| Spleen | 674,053 | 449,425 |
| Pituitary | 411,562 | 252,707 |
| Small intestine | 494,538 | 406,144 |

**Table 3: PacBio Isoseq sequences aligned against Sscrofa11.1**

## *Protein-to-genome Pipeline: Generating coding models using UniProt proteins*

Protein sequences were downloaded from UniProt and aligned to the genome in a splice aware manner using GenBlast [18]. The set of proteins aligned to the genome was a subset of UniProt proteins used to provide a broad, targeted coverage of the pig genome. The set consists of the following:

- Pig PE level 1, 2, 3

- Human PE level 1, 2, 3

- Mouse PE level 1, 2, 3

- Other mammals PE level 1, 2, 3

- Other vertebrates PE level 1, 2, 3

Note: PE level = protein existence level

A cut-off of 50 percent coverage and identity and an e-value of e-1 were used for GenBlast with the exon repair option turned on. The top 5 transcript models built by GenBlast for each protein passing the cut-offs were kept. This process produced 509,769 transcript models in total.

## *RNA-seq Pipeline*

RNA-seq data downloaded from ENA, PRJEB19386, was used in the annotation. This consisted of paired end, stranded data from twenty-eight tissue samples: alveolar macrophages, amygdala, brain stem, caecum,

cerebellum, colon, corpus callosum, duodenum, epididymis, frontal lobe, hippocampus, ileum, kidney cortex, left ventricle, mesenteric lymph node, medulla oblongata, occipital lobe, omentum, penis, pituitary gland, pituitary, pons, skeletal muscle, spleen, stomach, thalamus, tonsil, uterus. A merged file contain reads from all tissues was also created. The merged was less likely to suffer from model fragmentation due to read depth. The available reads were aligned to the genome using BWA. The Ensembl RNA-seq pipeline was used to process the BWA alignments and create further split read alignments using Exonerate.

The split reads and the processed BWA alignments were combined to produce 1,060,366 transcript models in total. The predicted open reading frames were compared to UniProt proteins using WU-BLAST. Models with poorly scoring or no BLAST alignments were split into a separate class and considered as potential lincRNAs.

## IG and TR genes

We downloaded all protein sequences from IMGT® [19] and aligned them against the genome using Exonerate using '--max-intron 50000' and only kept the models with 95% coverage and 80% identity. We generated 1,803 gene models.

## Olfactory receptor genes

We used the manually curated human and mouse set (Ensembl release 89) and pig olfactory receptor sequences [20]. We aligned the sequences against the genome with Exonerate and only kept the models with high similarity, 95% coverage and 95% identity. We generated 1,212 gene models.

## Selenocysteine proteins

We aligned known selenocysteine proteins against the genome using Exonerate. Then we checked 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. We generated 103 models.

# Filtering the Models

The filtering phase decided the subset of protein-coding transcript models, generated from the model-building pipelines, that comprise the final protein-coding gene set. Model are filtered based on information such as what pipeline they were generated using, how closely related the data are to the target species and how good the alignment coverage and precent identity to the original data are.

Models were filtered using the LayerAnnotation and GeneBuilder modules. The Apollo software [13] was used to visualise the results of filtering.

## *LayerAnnotation*

The LayerAnnotation module was used to define a hierarchy of input data sets, from most preferred to least preferred. The output of this pipeline included all transcript models from the highest ranked input set. Models from lower ranked input sets are 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 also separate 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 were not present in the vertebrate clade to avoid aligning the proteins multiple times.

**Layer 1:**

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

**Layer 2:**

- Pig cDNAs models with >= 90% coverage and 97% identity
- Pig IsoSeq models with protein support >= 80% coverage and identity and full RNA-seq support
- RNA-seq models with >= 95% coverage and identity

6

- Pig curated UniProt proteins from PE levels 1 & 2 with >= 80% coverage and identity and full RNA-seq support

- Pig curated UniProt proteins from PE levels 3 with >= 95% coverage and identity and full RNA-seq support

- All vertebrates curated UniProt proteins from PE levels 1 & 2 with >= 95% coverage and identity and full RNA-seq support

**Layer 3:**

- RNA-seq models with >= 80% coverage and identity

**Layer 4:**

- Pig curated UniProt proteins from PE levels 1 & 2 with >= 50% coverage and identity

- Pig IsoSeq models with protein support >= 80% coverage and identity

**Layer 5:**

- Pig curated UniProt proteins from PE levels 3 with >= 80% coverage and identity

- All vertebrates curated UniProt proteins from PE level 1 & 2 with >= 80% coverage and identity

**Layer 6:**

- RNA-seq models with >= 50% coverage and identity

- Pig IsoSeq models with protein support >= 50% coverage and identity

- Pig curated UniProt proteins from PE levels 3 with >= 50% coverage and identity

- All vertebrates curated UniProt proteins from PE level 1 & 2 with >= 50% coverage and identity

**Layer 7:**

- Pig UniProt proteins from PE levels 1 & 2 & 3 with >= 80% coverage and identity and full RNA-seq support

- All vertebrates UniProt proteins from PE levels 1 & 2 with >= 80% coverage and identity and full RNA-seq support

**Layer 8:**

- Pig UniProt proteins from PE levels 1 & 2 & 3 with >= 50% coverage and identity and full RNA-seq support

- All vertebrates UniProt proteins from PE levels 1 & 2 with >= 50% coverage and identity and full RNA-seq support

- Pig Isoseq models with protein support >= 50% coverage and identity which may

have a retained intron

**Layer 9:**

- Pig UniProt proteins from PE levels 1 & 2 & 3 with >= 80% coverage and identity

- All vertebrates UniProt proteins from PE levels 1 & 2 with >= 80% coverage and identity

**Layer 10:**

- Pig UniProt proteins from PE levels 1 & 2 & 3 with >= 50% coverage and identity

- All vertebrates UniProt proteins from PE levels 1 & 2 with >= 50% coverage and identity

## *Addition of UTR to coding models*

The set of coding models was extended into the untranslated regions (UTRs) using RNA-seq and cDNA and IsoSeqs sequences. The source of the UTRs was prioritised with UTR coming from cDNAs and IsoSeqs, then RNA-seq.

## *Generating multi-transcript genes*

The above steps generated a large set of potential transcript models, many of which overlapped 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.

At this stage the gene set comprised 23,025 genes with 46,511 transcripts.

## *Pseudogenes*

The Pseudogene module was run to identify pseudogenes from within the set of gene models. A total of 178 genes were labelled as pseudogenes or processed pseudogenes.

# Creating The Final Gene Set

## *Small ncRNAs*

Small structured non-coding genes were added using annotations taken from RFAM [14] and miRBase [15]. BLAST+ was run for these sequences and

models built using the Infernal software suite [17].

### *lincRNAs discovery*

Using the transcriptomic data set, we try to predict long intergenic non coding RNAs (lincRNAs). We used the RNA-seq data and the two IsoSeq sets which were filtered against the protein-coding gene set. The candidate lincRNAs should not overlap a protein-coding gene. The Pfam analysis of InterProScan is run against the filtered gene set. A potential lincRNA should not have a Pfam domain.

### *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.

As pig has been previously released in Ensembl a comparison was made to the previous gene set and as many stable identifiers as possible were mapped between the two annotations.

## Final Gene Set Summary

The final gene set consists of 22,439 protein coding genes, including 13 mitochondrial genes. These contain 45,898 transcripts. A total of 178 pseudogenes were identified. 2,320 small ncRNAs were added by the small ncRNA pipeline and 352 lincRNA were added by the lincRNA pipeline.

## *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 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 estimate

   o A higher coverage usually indicates a more complete assembly.

   o Using Sanger sequencing only, a coverage of at least 2x is preferred.

2. N50 of contigs and scaffolds

   o A longer N50 usually indicates a more complete genome assembly.

   o 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

   o A lower number toplevel sequences usually indicates a more complete genome assembly.

4. Alignment of cDNAs and ESTs to the genome

   o A higher number of alignments, using stringent thresholds, usually indicates a more complete genome assembly.

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

- Aken B et al.: **The Ensembl gene annotation system.** Database 2016. [PMID: 27337980]

- Potter SC, Clarke L, Curwen V, Keenan S, Mongin E, Searle SM, Stabenau A, Storey R, Clamp M: **The Ensembl analysis pipeline.** *Genome Res.* 2004, **14(5):**934-41. [PMID: 15123589]

- http://www.ensembl.org/info/genome/genebuild/index.html

- https://github.com/Ensembl/ensembl-doc/blob/master/pipeline_docs/the_genebuild_process.txt

## *References*

1  Smit, AFA, Hubley, R & Green, P: **RepeatMasker.** 1996-2010. www.repeatmasker.org

2  Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, Madden T.L: **BLAST+: architecture and applications.** *BMC Bioinformatics* 2008, **10**:421. [PMID: 20003500]

3  Benson G: **Tandem repeats finder: a program to analyze DNA sequences.** *Nucleic Acids Res.* 1999, **27(2):**573-580. [PMID: 9862982] http://tandem.bu.edu/trf/trf.html

4  Down TA, Hubbard TJ: **Computational detection and location of transcription start sites in mammalian genomic DNA.** *Genome Res.* 2002 **12(3):**458-461. http://www.sanger.ac.uk/resources/software/eponine/ [PMID: 11875034]

5  Lowe TM, Eddy SR: **tRNAscan-SE: a program for improved detection of transfer RNA genes in genomic sequence.** *Nucleic Acids Res.* 1997, **25(5):**955-64. [PMID: 9023104]

6  Burge C, Karlin S: **Prediction of complete gene structures in human genomic DNA.** *J Mol Biol.* 1997, **268(1):**78-94. [PMID: 9149143]

7  Goujon M, McWilliam H, Li W, Valentin F, Squizzato S, Paern J, Lopez R: **A new bioinformatics analysis tools framework at EMBL-EBI. Nucleic Acids Res.** 2010, **38 Suppl:**W695-699. http://www.uniprot.org/downloads [PMID: 20439314]

8    Sayers EW, Barrett T, Benson DA, Bolton E, Bryant SH, Canese K, Chetvernin V, Church DM, Dicuccio M, Federhen S, Feolo M, Geer LY, Helmberg W, Kapustin Y, Landsman D, Lipman DJ, Lu Z, Madden TL, Madej T, Maglott DR, Marchler-Bauer A, Miller V, Mizrachi I, Ostell J, Panchenko A, Pruitt KD, Schuler GD, Sequeira E, Sherry ST, Shumway M, Sirotkin K, Slotta D, Souvorov A, Starchenko G, Tatusova TA, Wagner L, Wang Y, John Wilbur W, Yaschenko E, Ye J: **Database resources of the National Center for Biotechnology Information.** *Nucleic Acids Res.* 2010, **38(Database issue):D5-16.** [PMID: 19910364]

9    http://www.ebi.ac.uk/ena/

10   Slater GS, Birney E: **Automated generation of heuristics for biological sequence comparison.** *BMC Bioinformatic*s 2005, **6:**31. [PMID: 15713233]

11   Birney E, Clamp M, Durbin R: **GeneWise and Genomewise.** *Genome Res.* 2004, **14(5):**988-995. [PMID: 15123596]

12   Eyras E, Caccamo M, Curwen V, Clamp M: **ESTGenes: alternative splicing from ESTs in Ensembl.** *Genome Res.* 2004 **14(5):**976-987. [PMID: 15123595]

13   Lewis SE, Searle SM, Harris N, Gibson M, Lyer V, Richter J, Wiel C, Bayraktaroglir L, Birney E, Crosby MA, Kaminker JS, Matthews BB, Prochnik SE, Smithy CD, Tupy JL, Rubin GM, Misra S, Mungall CJ, Clamp ME: **Apollo: a sequence annotation editor.** *Genome Biol.* 2002, **3(12):**RESEARCH0082. [PMID: 12537571]

14   Griffiths-Jones S., Bateman A., Marshall M., Khanna A., Eddy S.R: **Rfam: an RNA family database.** Nucleic Acids Research (2003) **31(1):**p439-441. [PMID: 12520045]

15   Griffiths-Jones S, Grocock RJ, van Dongen S, Bateman A, Enright AJ: **miRBase: microRNA sequences, targets and gene nomenclature.** NAR 2006 **34(Database Issue):**D140-D144 [PMID: 16381832]

16   Wilming L. G., Gilbert J. G. R., Howe K., Trevanion S., Hubbard T. and Harrow J. L: **The vertebrate genome annotation (Vega) database.** Nucleic Acid Res. 2008 Jan; Advance Access published on November 14, 2007; doi:10.1093/nar/gkm987 [PMID: 18003653]

17   Eddy, SR: **A memory-efficient dynamic programming algorithm for optimal alignment of a sequence to an RNA secondary structure**. BMC Bioinformatics 2002, 3:18. [PMID:12095421]

18   She R, Chu JS, Uyar B, Wang J, Wang K, and Chen N: **genBlastG: using BLAST searches to build homologous gene models.** Bioinformatics, 2011, [PMID: 21653517]

19  Lefranc MP, Giudicelli V, Duroux P, Jabado-Michaloud J, Folch G, Aouinti S, Carillon E, Duvergey H, Houles A, Paysan-Lafosse T, Hadi-Saljoqi S, Sasorith S, Lefranc G, Kossida S: **IMGT®, the international ImMunoGeneTics information system® 25 years on.** Nucleic Acids Res. 2012, **43:**D413-422; doi: 10.1093/nar/gku1056, [PMID: 25378316]

20  Nguyen DT, Lee K, Choi H, Choi MK, Le MT, Song N, Kim JH, Seo HG, Oh JW, Kim TH: **The complete swine olfactory subgenome: expansion of the olfactory gene repertoire in the pig genome.** BMC Genomics. 2012, **13:** 584-10.1186/1471-2164-13-584. [PMID: 23153364]