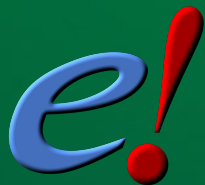


# Annotating Genomes the Ensembl Way



**Leanne Haggerty**

**Ensembl Genome Annotation Project Lead**

Biodiversity Genomics Academy 2023

Thursday 28th September 2023



**12:00 - 16:00**

**Annotating genomes the Ensembl way: General Concepts and Background**



**A Case study: Annotating Fungal Genomes - The Challenges!**



**14:00-16:00**

**Annotating genomes the Ensembl way: Hands-on - From RNAseq reads to gene models**



# What is Ensembl?

- A genomics platform for enabling and accelerating downstream science
- Most widely known for our genome browser
- Gene annotation, comparative genomics, variation and regulatory data
- Also provide data via FTP, REST/Perl API, MySQL dumps, BioMart

The screenshot shows the Ensembl genome browser homepage. At the top is a dark blue navigation bar with the Ensembl logo and links for BLAST/BLAT, VEP, Tools, BioMart, Downloads, Help & Docs, and Blog. A search bar on the right says "Search all species...". Below the navigation bar is a section with four tool links: "Tools" (with a sub-link "All tools"), "BioMart" (with a description: "Export custom datasets from Ensembl with this data-mining tool"), "BLAST/BLAT" (with a description: "Search our genomes for your DNA or protein sequence"), and "Variant Effect Predictor" (with a description: "Analyse your own variants and predict the functional consequences of known and unknown variants"). Below this is a search box with a dropdown menu set to "All species", a text input field, and a "Go" button. Below the search box is a section for "All genomes" with a dropdown menu set to "-- Select a species --". Below this is a section for "Pig breeds" with a pig icon and the text "Pig reference genome and 12 additional breeds". Below this is a link "View full list of all species". Below the "All genomes" section is a section for "Favourite genomes" with a pencil icon. Below the "Favourite genomes" section is a section with four tool links: "Compare genes across species", "Find SNPs and other variants for my gene" (with a DNA sequence snippet: "GTATACATTC", "CCTTAAAGTCTT", "CTCTTAAATGT", "GTACATTTTC"), "Gene expression in different tissues" (with a histology image), and "Retrieve gene sequence" (with a DNA sequence snippet: "GCTGAGCTCCGGGTTGG", "GGGCTTTGTGGGCGAGCT", "GGGCTTTGTGGGCGAGCT", "AGGAGGAGAGATTGTG", "GAGCTTTGAGAGGTTT", "CCGATCCAGGCTTGGG"). Below the "Retrieve gene sequence" section is a section for "Find a Data Display" and a section for "Use my own data in Ensembl" (with a diagram of a protein structure).

Ensembl is a genome browser for vertebrate genomes that supports research in comparative genomics, evolution, sequence variation and transcriptional regulation. Ensembl annotate genes, computes multiple alignments, predicts regulatory function and collects disease data. Ensembl tools include BLAST, BLAT, BioMart and the Variant Effect Predictor (VEP) for all supported species.

**Ensembl Release 109 (Feb 2023)**

- New gene sets for donkey and horse
- Updated SIFT and PolyPhen-2 missense variant pathogenicity
- New VEP plugins for UTR annotation
- New ATAC-seq tracks (peaks and signal) for fish species (Atlantic Salmon, European Seabass, Rainbow Trout and Turbot)

[More release news](#) on our blog

**Ensembl Rapid Release**

**New assemblies with gene and protein annotation every two weeks.**

Note: species that already exist on this site will continue to be updated with the full range of annotations.

[Go](#)

The Ensembl Rapid Release website provides annotation for recently produced, publicly available vertebrate and non-vertebrate genomes from biodiversity initiatives such as Darwin Tree of Life, the Vertebrate Genomes Project and the Earth BioGenome Project.

[Rapid Release news](#) on our blog

# Ensembl Division Sites

**EnsemblPlants** ~ 100 genomes

**EnsemblMetazoa** ~ 350 genomes

**EnsemblFungi** ~ 1 500 genomes

**EnsemblProtists** ~ 450 genomes

**EnsemblBacteria** ~ 30 000 genomes

**Archive sites**

Archive of release 45 of EnsemblBacteria: [pg45-bacteria.ensembl.org](http://pg45-bacteria.ensembl.org) (Sep 2019)

Archive of release 40 of EnsemblBacteria: [pg40-bacteria.ensembl.org](http://pg40-bacteria.ensembl.org) (July 2018)

Archive of release 37 of EnsemblBacteria: [pg37-bacteria.ensembl.org](http://pg37-bacteria.ensembl.org) (October 2017)

**Ensembl Bacteria**

Ensembl Bacteria is a browser for bacterial and archaeal genomes. These are taken from the databases of the [International Nucleotide Sequence Database Collaboration](http://www.ebi.ac.uk/seqdb/contributors/) (the European Nucleotide Archive at the EBI, GenBank at the NCBI, and the DNA Database of Japan).

**Non-redundant genomes**

As of release 35 (April 2017), we have only integrated new sequences that are non-redundant when compared to the existing data set, according to the criteria of the [UniProt Knowledgebase](http://www.ebi.ac.uk/seqdb/contributors/) (DOI: 10.1093/database/bax139). From early 2020, we will only be hosting non-redundant prokaryotic genomes. All existing data will be continue to be available via the archive sites.

**Data access**

Data can be visualised through the Ensembl genome browser and accessed programmatically via our Perl and RESTful APIs. Data is also accessible through public MySQL databases and our FTP site containing full data dumps in FASTA, EMBL, GTF, GFF3, JSON and RDF formats. A selection of over 100 key bacterial genomes have been included in the pan-taxonomic compars, and genes from all genomes have been classified into families using HAMAP and PANTHER ([more details](#)).

Ensembl Genomes is developed by [EMBL-EBI](http://www.ebi.ac.uk/seqdb/contributors/) and is powered by the [Ensembl](http://www.ebi.ac.uk/seqdb/contributors/) software system for the analysis and visualisation of genomic data. For details of our funding please [click here](#).

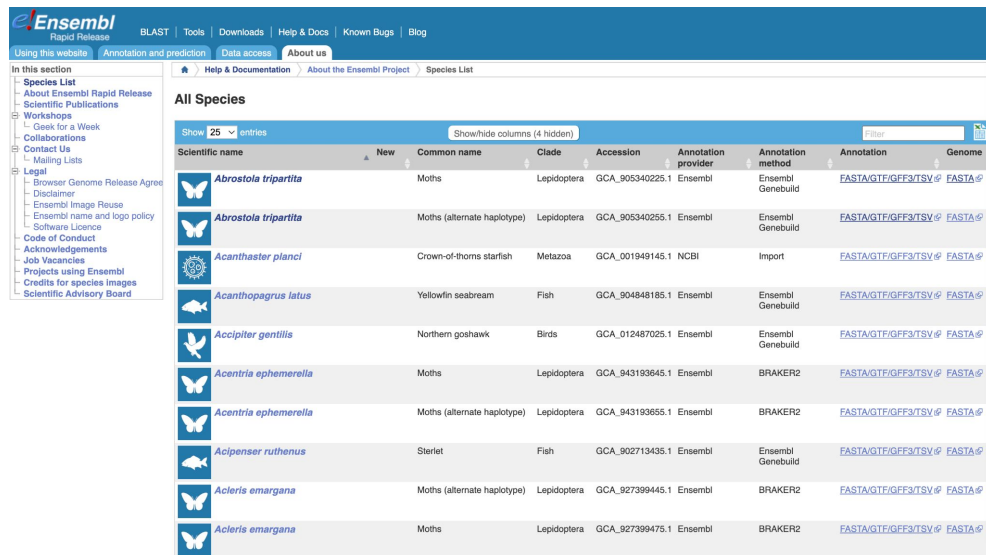
EMBL-EBI

# Rapid Release (rapid.ensembl.org)











- Runs on a two-week release cycle
- For deploying both verts/non-verts
- >1,700 genomes for >1000 different species since June 2020

- Current functionality

- Gene annotation
- Homologues
- Repeat tracks
- Protein feature annotation
- BLAST
- HAL multiple alignments for certain clades
- Variation data for human, fruit fly and pigeon pea



The screenshot shows the Ensembl Rapid Release website. The top navigation bar includes links for BLAST, Tools, Downloads, Help & Docs, Known Bugs, and Blog. Below this is a sub-navigation bar with links for Using this website, Annotation and prediction, Data access, and About us. The main content area is titled 'All Species' and displays a table of species. The table has columns for Scientific name, New, Common name, Clade, Accession, Annotation provider, Annotation method, Annotation, and Genome. The table lists 10 species, including *Abrostola tripartita*, *Acanthaster planci*, *Acanthopagrus latus*, *Accipiter gentilis*, *Acentria ephemerella*, *Acipenser ruthenus*, and *Acleris emargana*. Each row includes a small icon representing the species and links for FASTA/GTF/GFF3/TSV and FASTA files.

Scientific name	New	Common name	Clade	Accession	Annotation provider	Annotation method	Annotation	Genome
 <i>Abrostola tripartita</i>		Moths	Lepidoptera	GCA_905340225.1	Ensembl	Ensembl Genebuild	<a href="#">FASTA/GTF/GFF3/TSV</a>	<a href="#">FASTA</a>
 <i>Abrostola tripartita</i>		Moths (alternate haplotype)	Lepidoptera	GCA_905340255.1	Ensembl	Ensembl Genebuild	<a href="#">FASTA/GTF/GFF3/TSV</a>	<a href="#">FASTA</a>
 <i>Acanthaster planci</i>		Crown-of-thorns starfish	Metazoa	GCA_001949145.1	NCBI	Import	<a href="#">FASTA/GTF/GFF3/TSV</a>	<a href="#">FASTA</a>
 <i>Acanthopagrus latus</i>		Yellowfin seabream	Fish	GCA_904848185.1	Ensembl	Ensembl Genebuild	<a href="#">FASTA/GTF/GFF3/TSV</a>	<a href="#">FASTA</a>
 <i>Accipiter gentilis</i>		Northern goshawk	Birds	GCA_012487025.1	Ensembl	Ensembl Genebuild	<a href="#">FASTA/GTF/GFF3/TSV</a>	<a href="#">FASTA</a>
 <i>Acentria ephemerella</i>		Moths	Lepidoptera	GCA_943193645.1	Ensembl	BRAKER2	<a href="#">FASTA/GTF/GFF3/TSV</a>	<a href="#">FASTA</a>
 <i>Acentria ephemerella</i>		Moths (alternate haplotype)	Lepidoptera	GCA_943193655.1	Ensembl	BRAKER2	<a href="#">FASTA/GTF/GFF3/TSV</a>	<a href="#">FASTA</a>
 <i>Acipenser ruthenus</i>		Sterlet	Fish	GCA_902713435.1	Ensembl	Ensembl Genebuild	<a href="#">FASTA/GTF/GFF3/TSV</a>	<a href="#">FASTA</a>
 <i>Acleris emargana</i>		Moths (alternate haplotype)	Lepidoptera	GCA_927399445.1	Ensembl	BRAKER2	<a href="#">FASTA/GTF/GFF3/TSV</a>	<a href="#">FASTA</a>
 <i>Acleris emargana</i>		Moths	Lepidoptera	GCA_927399475.1	Ensembl	BRAKER2	<a href="#">FASTA/GTF/GFF3/TSV</a>	<a href="#">FASTA</a>

# Genome Annotation

# Genome Annotation

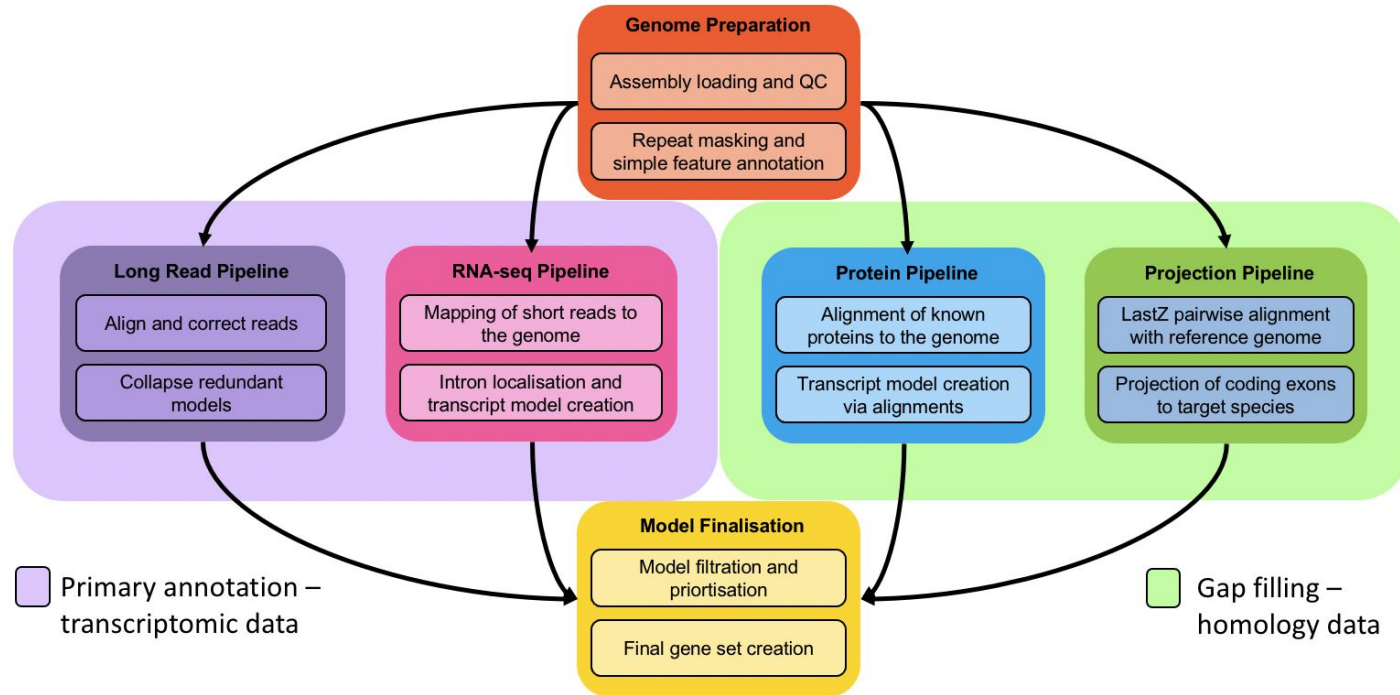
- The process of identifying and labelling information on a genome
- Coordinate-based: repeats, genes, transcripts, exons, variants, regulatory regions
- Knowledge-based: gene function, variant effect, repeat type
- Context-based: orthology/paralogy, synteny

# Annotation Approaches

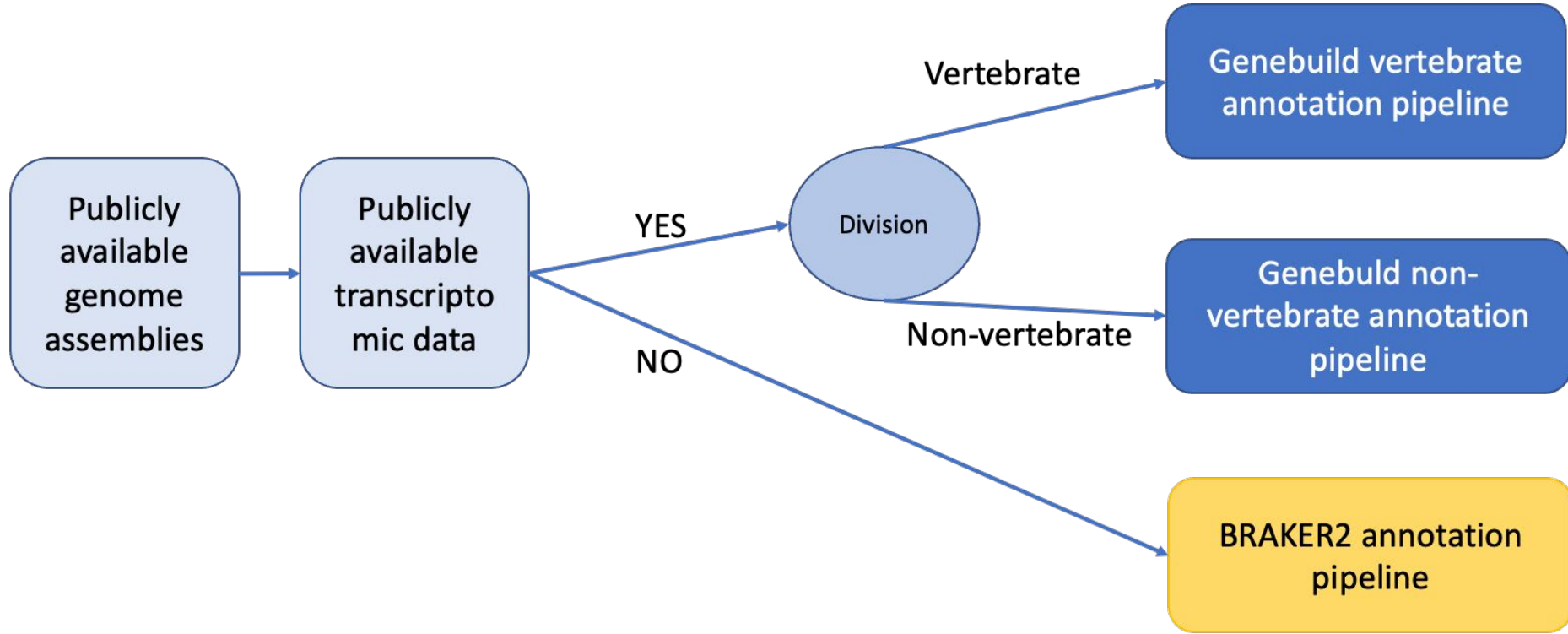
Approach	Main strengths	Main weaknesses
Ab initio	<ul style="list-style-type: none"><li>• Fast</li><li>• Only genome required as input</li></ul>	<ul style="list-style-type: none"><li>• Large numbers of false positives</li><li>• Large amounts of missing/spurious exons</li></ul>
Protein-genome alignments	<ul style="list-style-type: none"><li>• Fast</li><li>• Uses existing data</li><li>• Can give excellent coverage</li></ul>	<ul style="list-style-type: none"><li>• No UTR/lncRNAs</li><li>• Accuracy decreases significantly as evolutionary distance increases</li><li>• Highly dependent on the proteins available</li><li>• Doesn't capture novelty</li></ul>
Projection/liftover	<ul style="list-style-type: none"><li>• Can be highly accurate even at moderate distances</li><li>• Can leverage high quality reference annotations</li></ul>	<ul style="list-style-type: none"><li>• Expensive if using a full WGA</li><li>• Doesn't capture novelty</li><li>• Genes not covered in the WGA not projected</li></ul>
Transcriptomic	<ul style="list-style-type: none"><li>• Fast and accurate for finding structures</li><li>• Captures novel genes and transcripts</li><li>• Annotation of non-coding elements</li><li>• Allows for tissue/timepoint-specific tracks</li><li>• Can be used to help validate structures identified through other approaches</li></ul>	<ul style="list-style-type: none"><li>• Difficult to get high coverage of full gene set</li><li>• Sampling/Cost</li><li>• Fragmentation</li><li>• Short read data requires inference of structures</li></ul>



# Ensembl Annotation Pipelines



# Ensembl Annotation Pipelines



# Repeat Annotation

# Repeat Annotation

## Types of Repeats

- Low complexity regions
  - Poly-purine or poly-pyrimidine stretches
  - Regions of extremely high AT or GC content
- Transposable elements
  - Class I retrotransposons (“copy-and-paste”)
  - Class II DNA transposons (“cut-and-paste”)
- Satellite DNA
  - Short and long tandem repeats

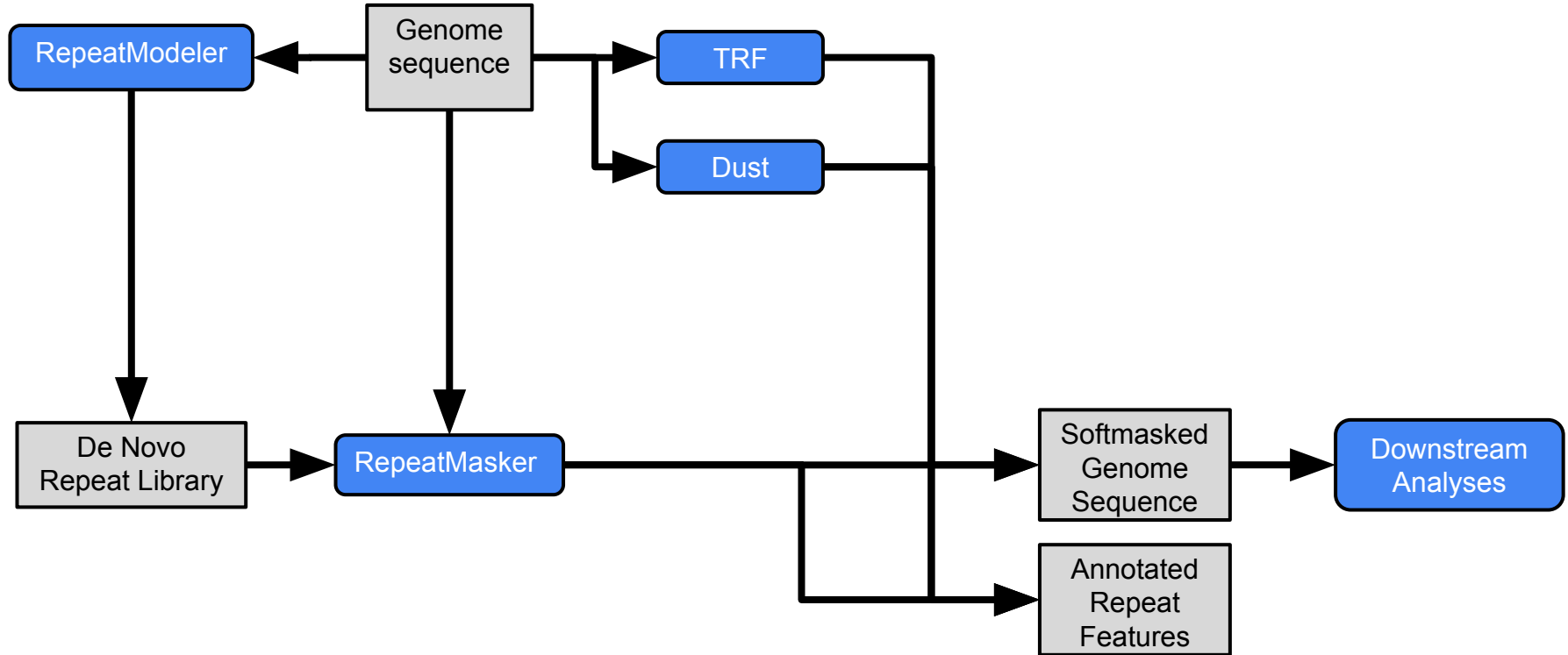
Favourite tracks	
Track order	
Search results	
[-] <b>Sequence and assembly</b>	<b>(2/29)</b>
└ Sequence	(2/4)
└ Simple features	(0/3)
└ Clones & misc. regions	(0/22)
<b>Genes and transcripts</b>	<b>(3/4)</b>
[-] <b>mRNA and protein alignments</b>	<b>(0/2)</b>
└ mRNA alignments	(0/1)
└ Protein alignments	(0/1)
[-] <b>Variation</b>	<b>(1/10)</b>
└ Sequence variants	(0/3)
└ Failed variants	(0/1)
└ Phenotype annotations	(0/2)
└ Structural variants	(1/4)
[-] <b>Regulation</b>	<b>(0/245)</b>
└ Regulatory Build	(0/1)

## Repeat regions

### Enable/disable all tracks

- ☒ All repeats
- ☐ Low complexity (Dust)
- ☐ Repeats (Mouse)
- ☐ LTRs (Repeats (Mouse))
- ☐ Other repeats (Repeats (Mouse))
- ☐ RNA repeats (Repeats (Mouse))
- ☐ Satellite repeats (Repeats (Mouse))
- ☐ Type I Transposons/LINE (Repeats (Mouse))
- ☐ Type I Transposons/SINE (Repeats (Mouse))
- ☐ Type II Transposons (Repeats (Mouse))
- ☐ Unknown (Repeats (Mouse))
- ☐ Tandem repeats (TRF)

# Repeat Annotation



# Repeat Annotation

- Repeat masking an important initial step in genome annotation
- Repeat annotation presents a lot of challenges
  - It is computationally very costly
  - Repeat libraries can sometimes contain gene families
  - Huge volume of software, but only a few long lasting/well supported ones
- Red (**RE**peat**D**etector) is extremely efficient tool if masking alone is the goal

# Gene Annotation

# Gene Annotation

- Can be broken into a number of approaches:
  - Transcriptomic – using long/short reads
  - Homology – Mapping/lifting data from other species
  - Ab initio – Using HMMs or other predictive methods
  - Hybrid – Using a combination of the above



# Gene Annotation

## *Transcriptomic*

# Gene Annotation - *Transcriptomic*

- Two major types of data – **short reads** (usually Illumina) and **long reads** (PacBio Iso-Seq, ONT)
- Short reads generally through two initial phases:
  - Align to genome to find exon positions
  - Reconstruction into potential transcript isoforms
- Long reads have one initial phase and one optional:
  - Align to the genome to find transcript structures
  - Error correcting the identified structures
- After the initial phases both data types have putative transcript models, then usually:
  - Collapse redundant/fragmented models
  - Identify longest ORF
  - Classify ORF as coding/non-coding

# Gene Annotation - *Transcriptomic*

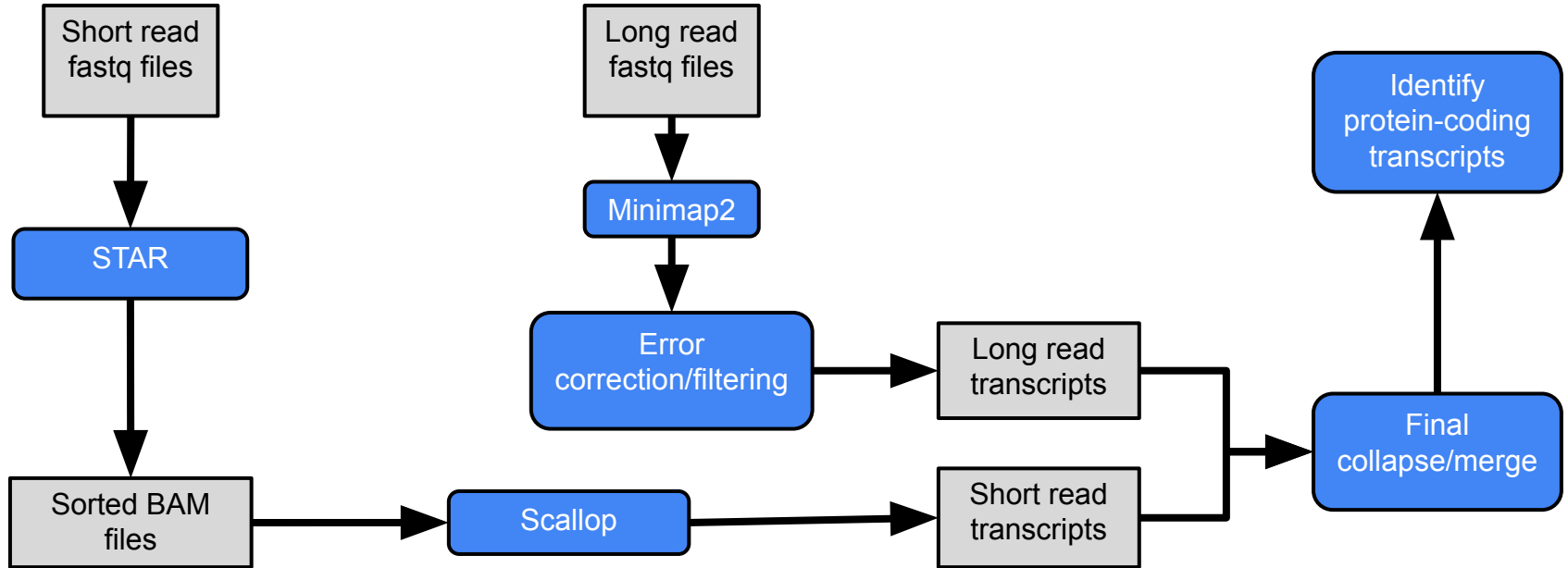
Approach	Main strengths	Main weaknesses
Short reads	<ul style="list-style-type: none"><li>• Better coverage</li><li>• Cheap</li><li>• High accuracy</li><li>• Can be easier to validate structures from other approaches</li></ul>	<ul style="list-style-type: none"><li>• Inferring transcripts, not observing them</li><li>• UTRs difficult to call</li></ul>
Long reads	<ul style="list-style-type: none"><li>• Full length transcript structures</li><li>• More accurate UTRs</li><li>• Can be more accurate for spliced lncRNAs</li></ul>	<ul style="list-style-type: none"><li>• Lots of single exon transcriptional noise</li><li>• High error rate for ONT data</li></ul>

# Gene Annotation - *Transcriptomic*

**Disclaimer: This applies to vertebrates!!!**

- **Minimal**
  - Short read data only
  - Highest value: brain, gonads, lung/gill, embryo
  - Lowest value: liver, muscle, blood
  - 100-150bp
  - 100 million reads+ per tissue
- **Ideal**
  - Short and long read data
  - 5+ tissues
  - Dev stages if possible
  - The more reads the better
  - Preference for consensus/cleaned over raw reads

# Gene Annotation - *Transcriptomic*



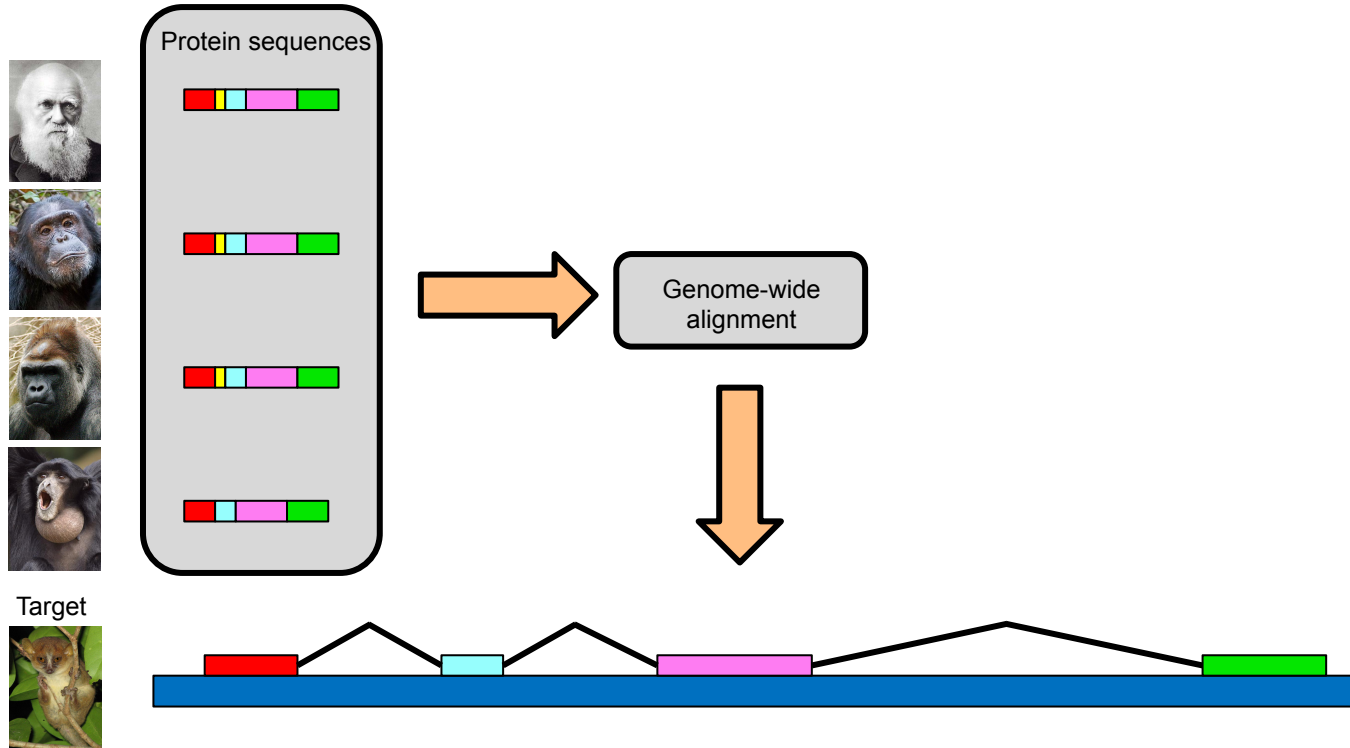
# Gene Annotation

## *Homology*

# Gene Annotation - *Homology*

- Cross species protein alignments
  - Protein sequences are more conserved than nucleotide sequences
  - There are vast numbers of available protein sequences, with varying levels of evidence

# Cross Species Protein Alignments



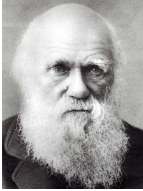


# Gene Annotation - *Homology*

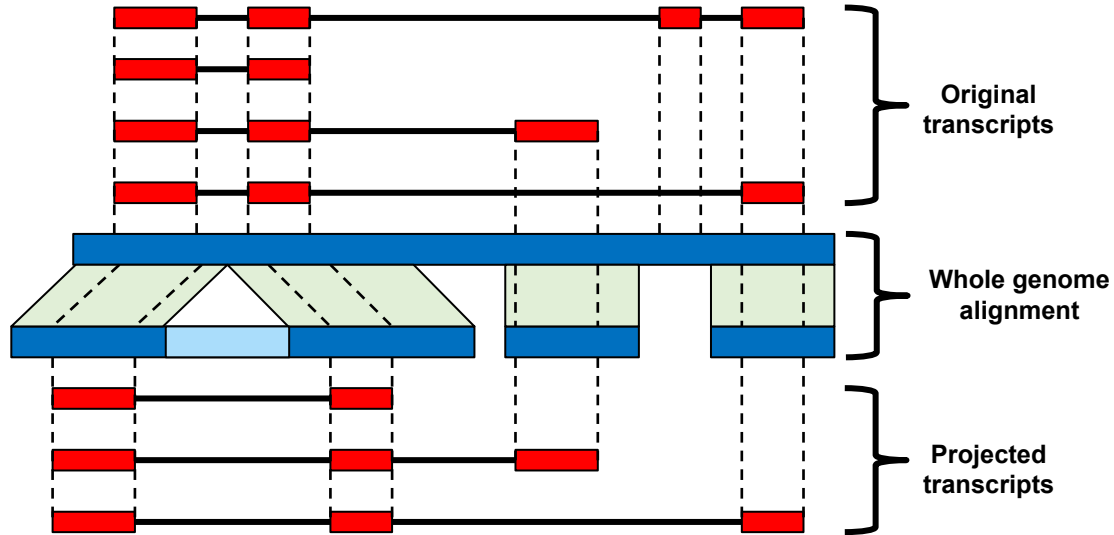
- Projection from a reference
  - Projection/liftover methods generally are based off a whole genome alignment
  - Generally involve high quality reference assembly and annotation being projected to a target assembly
  - Some methods focus on breeds/strains/haplotypes, and allow projection of coding and non-coding features
  - Others focus on only the coding exons to allow projection across greater evolutionary distance

# Projection From a Reference

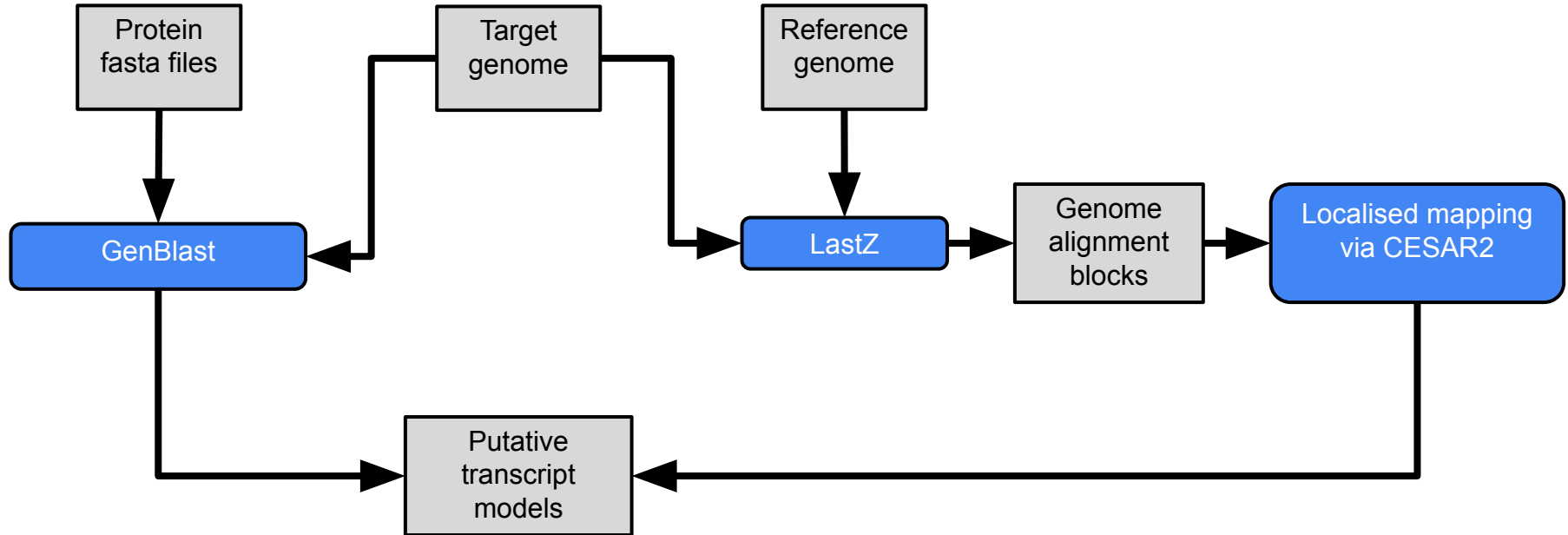
Species A



Species B



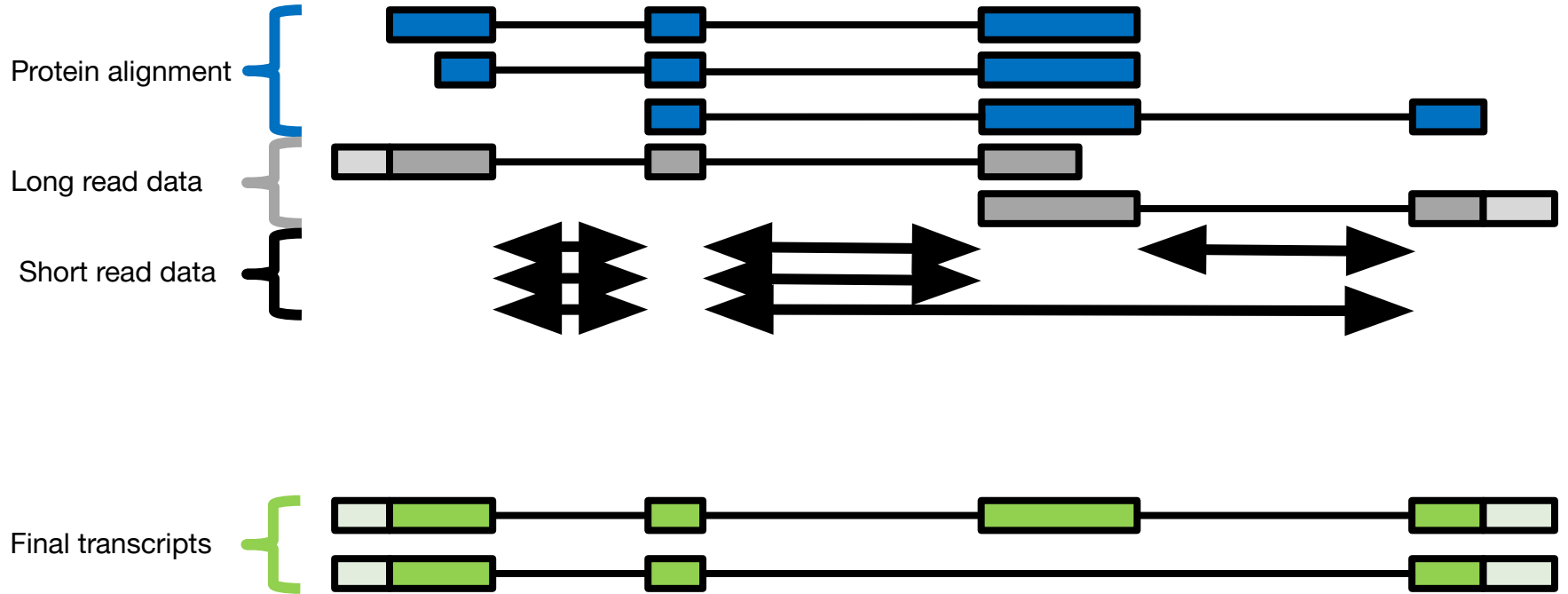
# Gene Annotation - *Homology*



# **Gene Annotation**

## ***Combined Evidence***

# Gene Annotation - *Combined Evidence*



# Gene Annotation - *Combined Evidence*

- Data rich:
  - More transcript isoforms
  - More genes captured per class, including timepoint/tissue specific ones
  - UTR features captured
  - Better exon capture, especially around start/end exons
  - Ability to measure confidence of the gene set
- Data poor:
  - Often focused on protein coding genes
  - Generally one isoform per gene
  - Often missing exons/misalignments of start/end exons
  - No UTRs
  - False positives

# **Gene Annotation**

## ***Assessing Quality***

# Gene Annotation - *Assessing Quality*

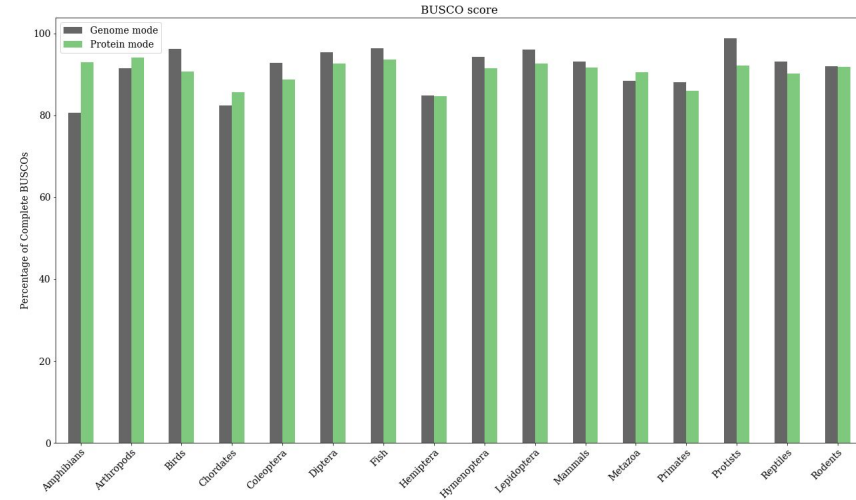
- Assessing the quality of a gene set can be difficult
- The closer to a well annotated reference, the easier the assessment
- Things to look for:
  - One-to-one orthologues (or reciprocal best BLAST hits) with references
  - Long/split/orphan gene counts
  - Average coding exons/genomic span/CDS length
  - BUSCO/ OMArk completeness for most appropriate taxonomic group



# Gene Annotation - *Assessing Quality*

## BUSCO

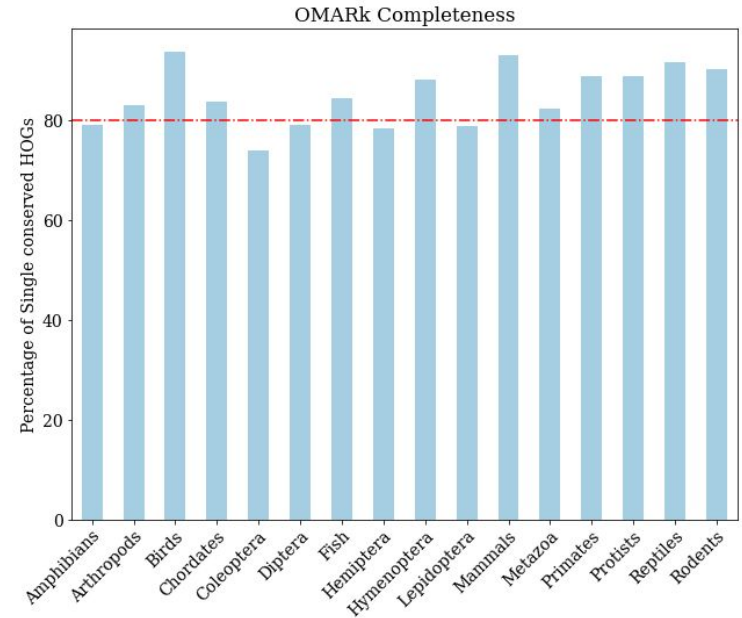
- A measure for quantitative assessment of genome assembly and annotation completeness based on evolutionarily informed expectations of gene content.
- Based on the concept of single-copy orthologs that should be highly conserved among the closely related species



# Gene Annotation - *Assessing Quality*

## OMArk

- Estimate the proteome completeness by comparison to conserved orthologous groups defined using Hierarchical Orthologous Groups (HOGs).
- Estimate the proportion of accurate and erroneous gene models in the proteome by comparing to the known gene families of the selected ancestral lineage
- Detect possible contamination from other species in the proteome.



# Summary

- **Repeat annotation** is usually the first step in genome annotation
- Many sources of repeat libraries and tools, but RepeatModeler and RepeatMasker combo is most popular approach
- Gene annotation can be done with a variety of approaches, with **transcriptomic data** being most valuable
- The **quality of annotation** will always be dependent on the quality of the input data used
- Important to take into account the quality for planning downstream analyses

## The Eukaryotic Annotation Team



**Fergal Martin**

Eukaryotic Annotation Team Leader

## The Genebuild Team



**Leanne Haggerty**

Ensembl Genome Annotation Project Lead



**Swati Sinha**

Senior Bioinformatician



**Francesca Floriana Tricomi**

Bioinformatician



**Jose Maria Gonzalez Perez-Silva**

Bioinformatician



**Vianey Paola Barrera Enriquez**

Bioinformatician

## The Comparative Genomics Team



**Thiago Augusto Lopes Genez**

Software Engineer



**Thomas Walsh**

Senior Bioinformatician



**Botond Sipos**

Bioinformatics Developer



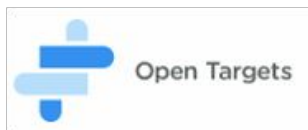
**Ivana Pilizota**

Bioinformatics Developer



**Simarpreet Kaur Bhurji**

Bioinformatician



Co-funded by  
the European  
Union

