**ETH**_zürich_

**Universität Zürich** UZH

# Gene identifier matching to join publicly available databases for the generation of a Mammalian Ortholog and Annotation Database with access from Galaxy-server

**Jochen Bick**[1], Mark Robinson[2], Susanne E. Ulbrich[1] and Stefan Bauersachs[1]

**[1]Animal Physiology, ETH Zurich;** [2]Institute of Molecular Life Sciences, University Zurich (UZH)
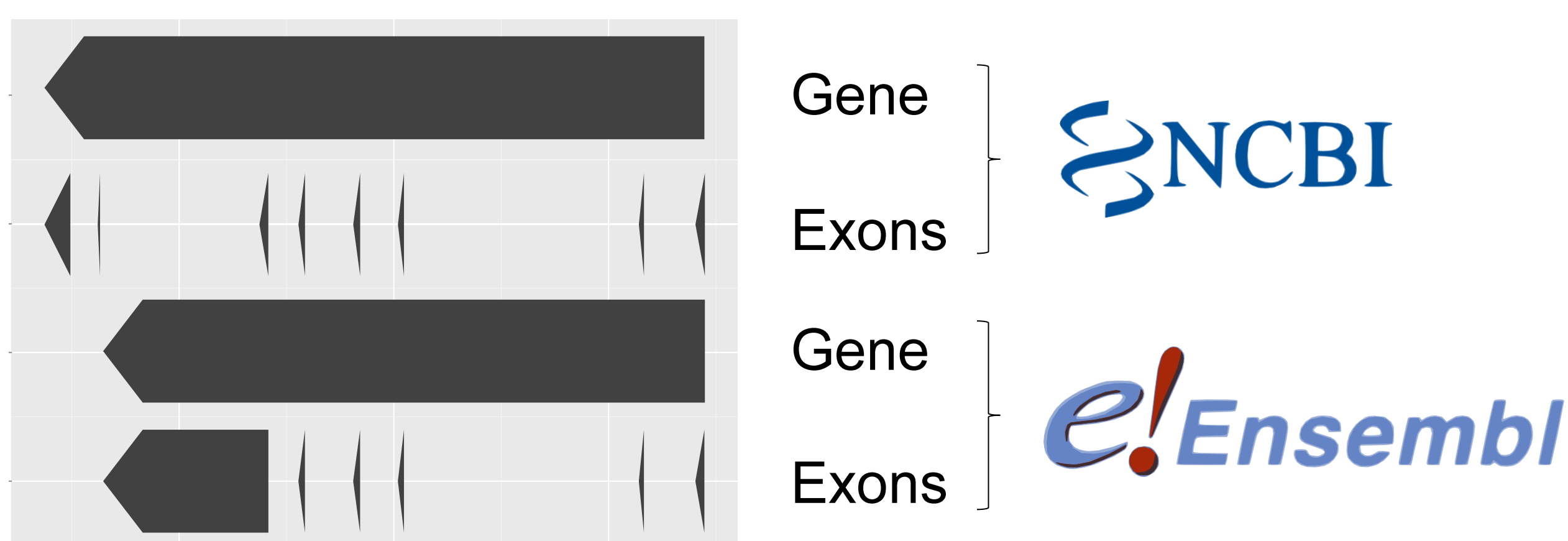
## Introduction

So far there is a number of well-organized databases that contain useful information regarding orthologous genes, e.g., EnsemblCompara ortholog database. The main problem when using information derived from different databases is to correctly assign different gene, transcript or protein identifiers. However, because NCBI annotation is for most species the most comprehensive, we need to map information from other databases to EntrezGene IDs. This is an important issue for the generation of a Mammalian Ortholog and Annotation Database (MOA-Db) which will be partially based on information from publicly available databases, which needs to be collected, analyzed, and connected. Since each public source database uses own unique identifiers, it is necessary to assign the corresponding database-specific identifiers. Existing lists that assign corresponding genes, e.g. between Ensembl and EntrezGene are incomplete and/or contain errors. R BioConductor packages were used to find overlapping gene and exon positions which were integrated as a lookup table into the MySQL database to handle the comparison of different database sources. Finally, this database will be integrated into our local Galaxy-server to give easy access to all our research groups and provide a useful interface with various options to parse information via SQL queries. The MOA-Db provides a basis for optimal across-species comparisons of transcriptome datasets from different mammalian species accessible within a Galaxy-server.
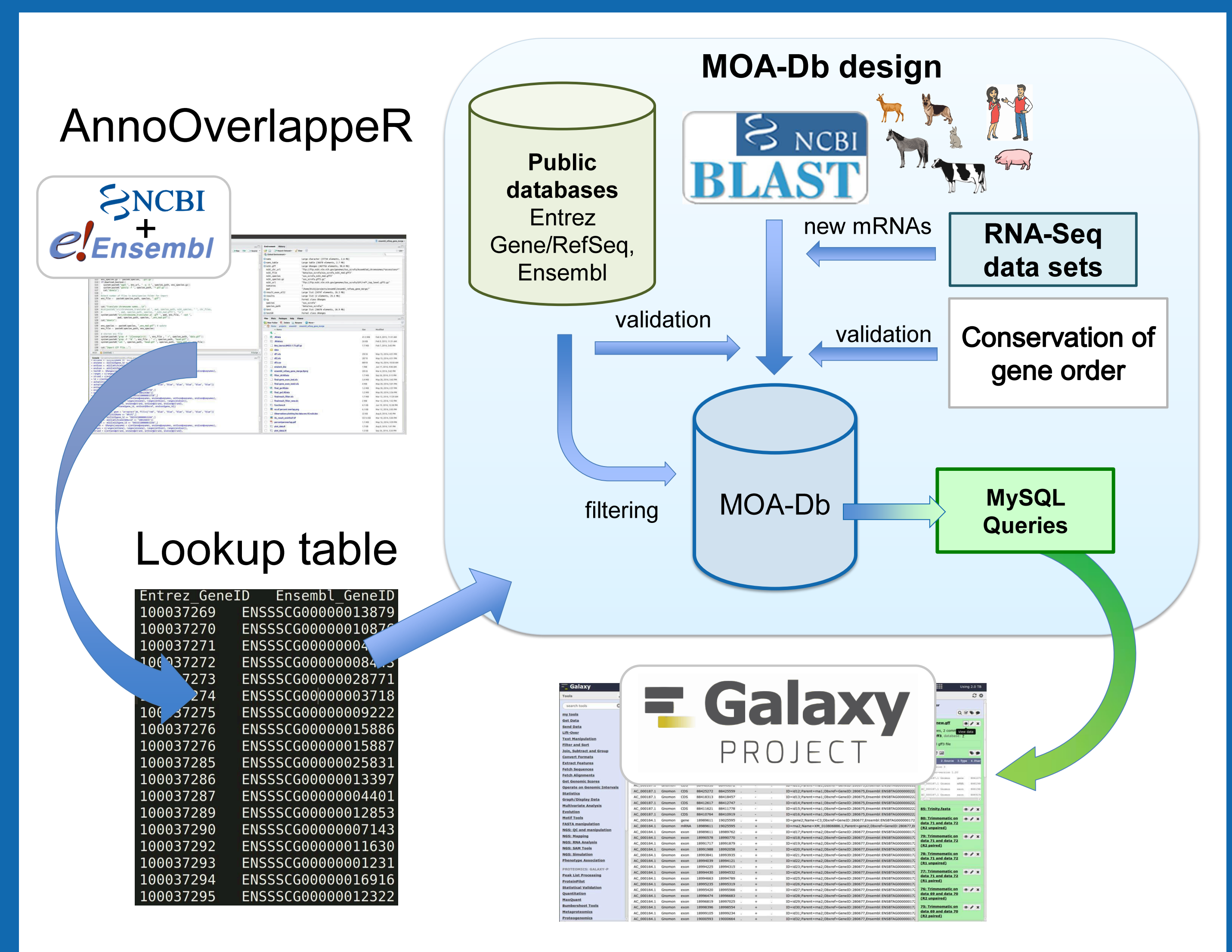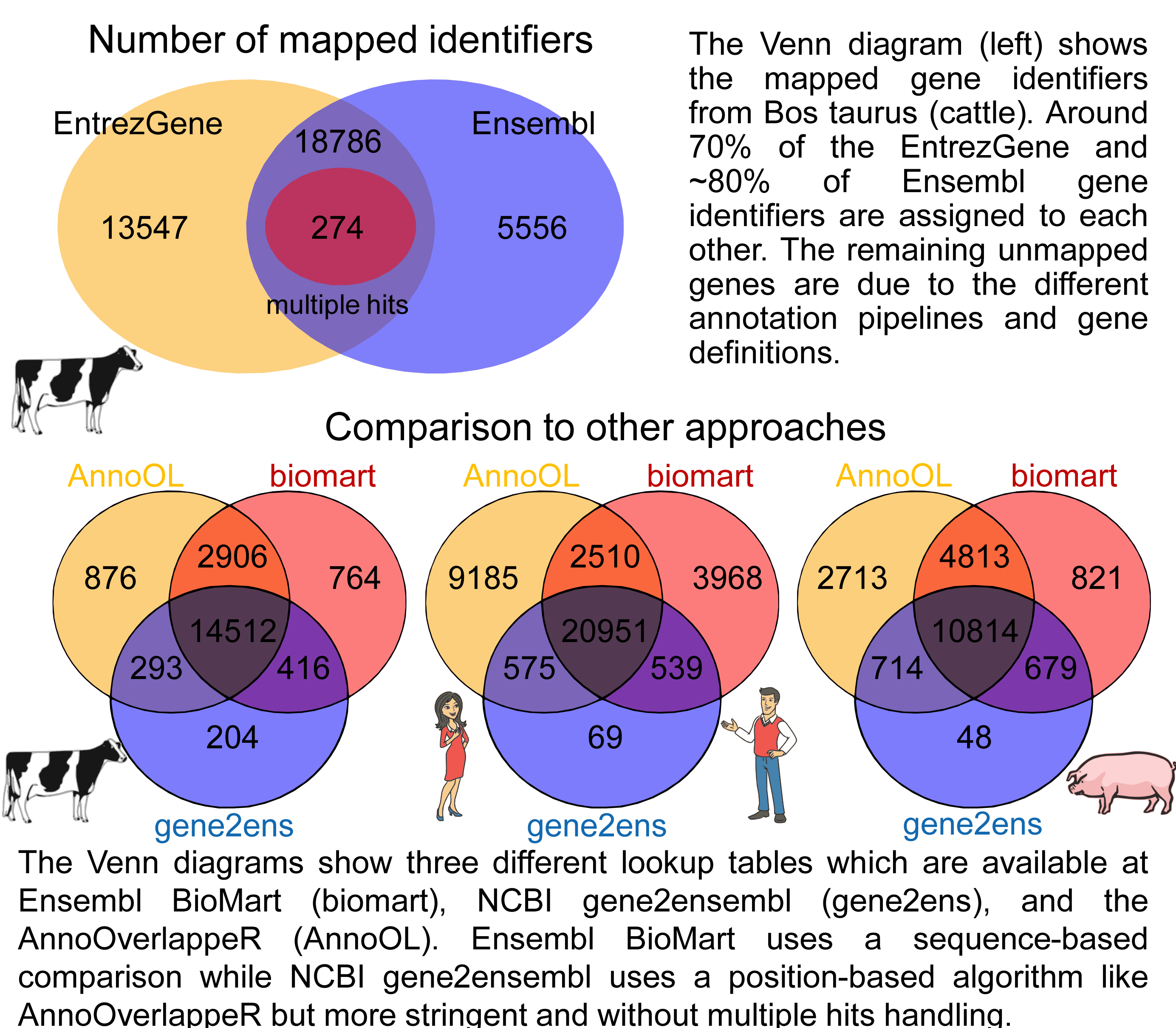
## AnnoOverlappeR

Three steps overlap algorithm:

1. Overlap on gene level
- Position-based overlap (chromosome, start, end, strand)
- Catch differently annotated genes (multiple hits)
2. Overlap on corresponding splice sites

- Cutoff > 50 % overlap
- If < 50% ➡ check exon overlap, if > 50% keep it
3. Handling of multiple hits
- One gene maps to multiple genes



Gene / Exons — **NCBI**

Gene / Exons — **e!Ensembl**

## Workflow



## Results

### Number of mapped identifiers



EntrezGene 13547 | 18786 | 274 | Ensembl 5556

multiple hits

The Venn diagram (left) shows the mapped gene identifiers from Bos taurus (cattle). Around 70% of the EntrezGene and ~80% of Ensembl gene identifiers are assigned to each other. The remaining unmapped genes are due to the different annotation pipelines and gene definitions.

### Comparison to other approaches



AnnoOL / biomart / gene2ens:
876 | 2906 | 764
14512
293 | 416
204

AnnoOL / biomart / gene2ens:
9185 | 2510 | 3968
20951
575 | 539
69

AnnoOL / biomart / gene2ens:
2713 | 4813 | 821
10814
714 | 679
48

The Venn diagrams show three different lookup tables which are available at Ensembl BioMart (biomart), NCBI gene2ensembl (gene2ens), and the AnnoOverlappeR (AnnoOL). Ensembl BioMart uses a sequence-based comparison while NCBI gene2ensembl uses a position-based algorithm like AnnoOverlappeR but more stringent and without multiple hits handling.
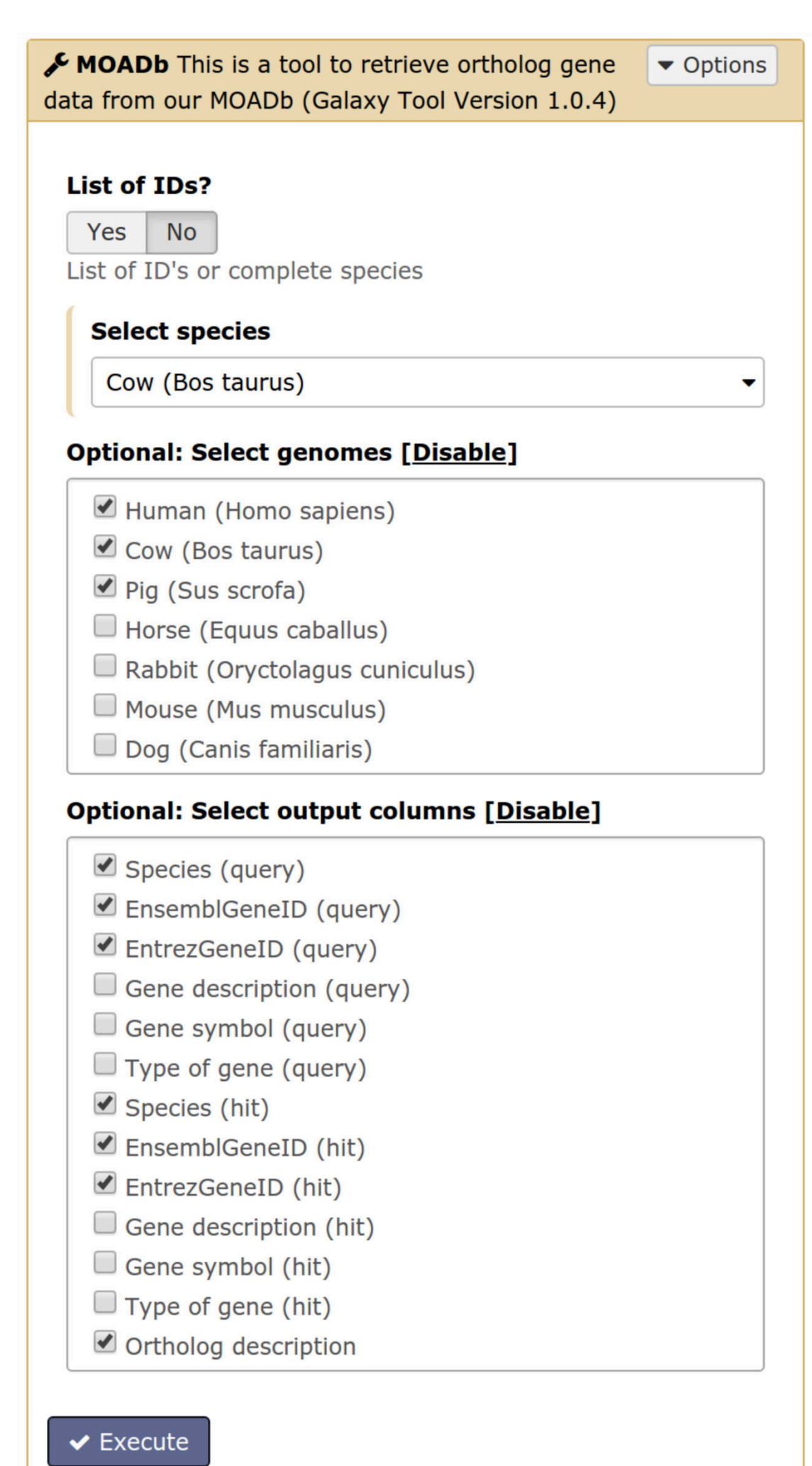
## Visualisation in Galaxy

The figure on the right shows a screenshot of the MOA-Db interface on our Galaxy server. For now you can chose between a list of identifiers (IDs) provided by your history or a complete species like selected in the example. 'Select genomes' will specify different mammals which will be matched inside the database. 'Select output columns' decides the columns to be shown in the output file.

Right now we support two different database IDs EntrezGene IDs and Ensembl gene IDs. Our database also contains data from seven different species human, cow, pig, horse, rabbit, mouse and dog.



The figure shows the output in Galaxy which is represented in "tabular" format. You can see all selected features in corresponding order as shown in the right figure.