# A computational architecture designed for genome annotation: oak genome sequencing project as a use case

Nicolas Francillonne, Tina Alaeitabar, Françoise Alfama-Depauw, Loïc Couderc, Claire Guerche, Thomas Letellier, Mikaël Loaec, Isabelle Luyten, Célia Michotey, Jean-Marc Aury, et al.

# A computational architecture designed for genome annotation:
## Oak genome sequencing project as a use case

N, FRANCILLONNE [1], T, ALAEITABAR [1], F, ALFAMA [1], L, COUDERC [1], C, GUERCHE [1], T, LETELLIER [1], M, LOAEC [1], I, LUYTEN [1], C, MICHOTEY [1], J, AURY [2], H, QUESNEVILLE [1], C. PLOMION [3], J. AMSELEM [1] and the Oak genome sequencing consortium

nicolas.francillonne@versailles.inra.fr
joelle.amselem@versailles.inra.fr

1 INRA, UR1164 URGI - Research Unit in Genomics-Info, INRA de Versailles, Route de Saint-Cyr, Versailles, 78026, France
2 CEA, Institut de Génomique, Genoscope, BP5706, 91057 Evry, France
3 INRA, UMR1202 – Biodiversity, Genes and Communities - 69, route d'Arcachon, F-33610 Cestas, France

The ANR Genoak project aims to study the two key evolutionary processes that explain the remarkable diversity found within the oak genus. We performed an automated structural annotation (transposable elements (TEs) and genes) and functional annotation of predicted genes using robust pipelines i/ REPET for TEs ii/ Eugene for gene prediction iii/ FunAnnotPipe (in-house pipeline) mainly based on InterproScan for functional annotation. Further objectives were to: i/ integrate the whole genome with all the features annotated into a Genome Browser, ii/ provide an interface for gene prediction curation/validation, and iii/ provide an information system pointing towards accessibility and interoperability.
We are setting up a fast and flexible genome browser WebApollo_oak allowing the edition of genes structure. This tool based on Jbrowse is used to visualize, identify and curate gene predictions. For functional annotation, we set up a data warehouse QuercusRoburMine based on Intermine technology. Its user friendly interface gives access to functional information which allow cross queries between different data sources. These combined tools constitute powerful resources for whole genome annotation. Next step will be to automatically update all functional information in QuercusRoburMine for gene curated through WebApollo_oak. We will present some case studies to illustrate questions raised by the users.

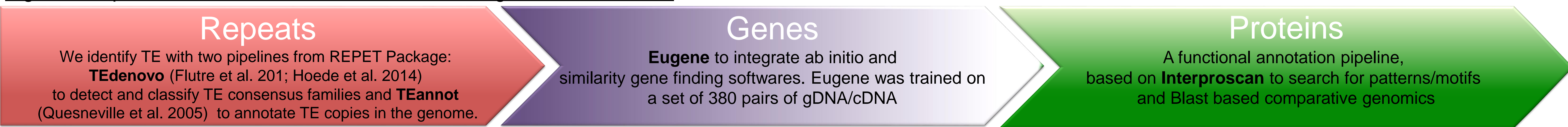## Figure 1: Pipelines used for structural and functional genome annotation



**Repeats** — We identify TE with two pipelines from REPET Package: **TEdenovo** (Flutre et al. 201; Hoede et al. 2014) to detect and classify TE consensus families and **TEannot** (Quesneville et al. 2005) to annotate TE copies in the genome.

**Genes** — **Eugene** to integrate ab initio and similarity gene finding softwares. Eugene was trained on a set of 380 pairs of gDNA/cDNA

**Proteins** — A functional annotation pipeline, based on **Interproscan** to search for patterns/motifs and Blast based comparative genomics

### Figure 2a: Several ways to query data



1 - Quick search
2 – Protein identifier
3 - Genomic region
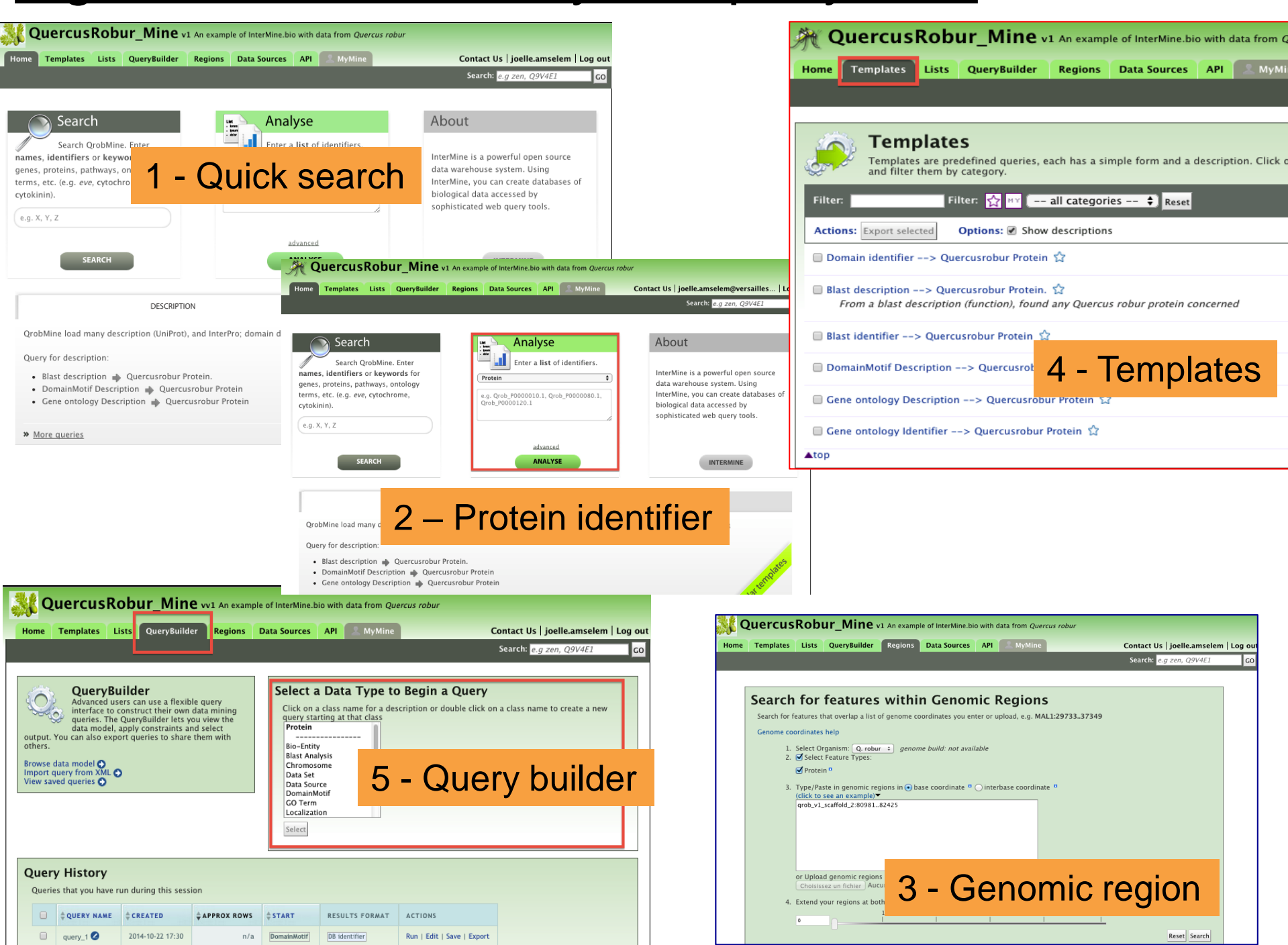4 - Templates
5 - Query builder

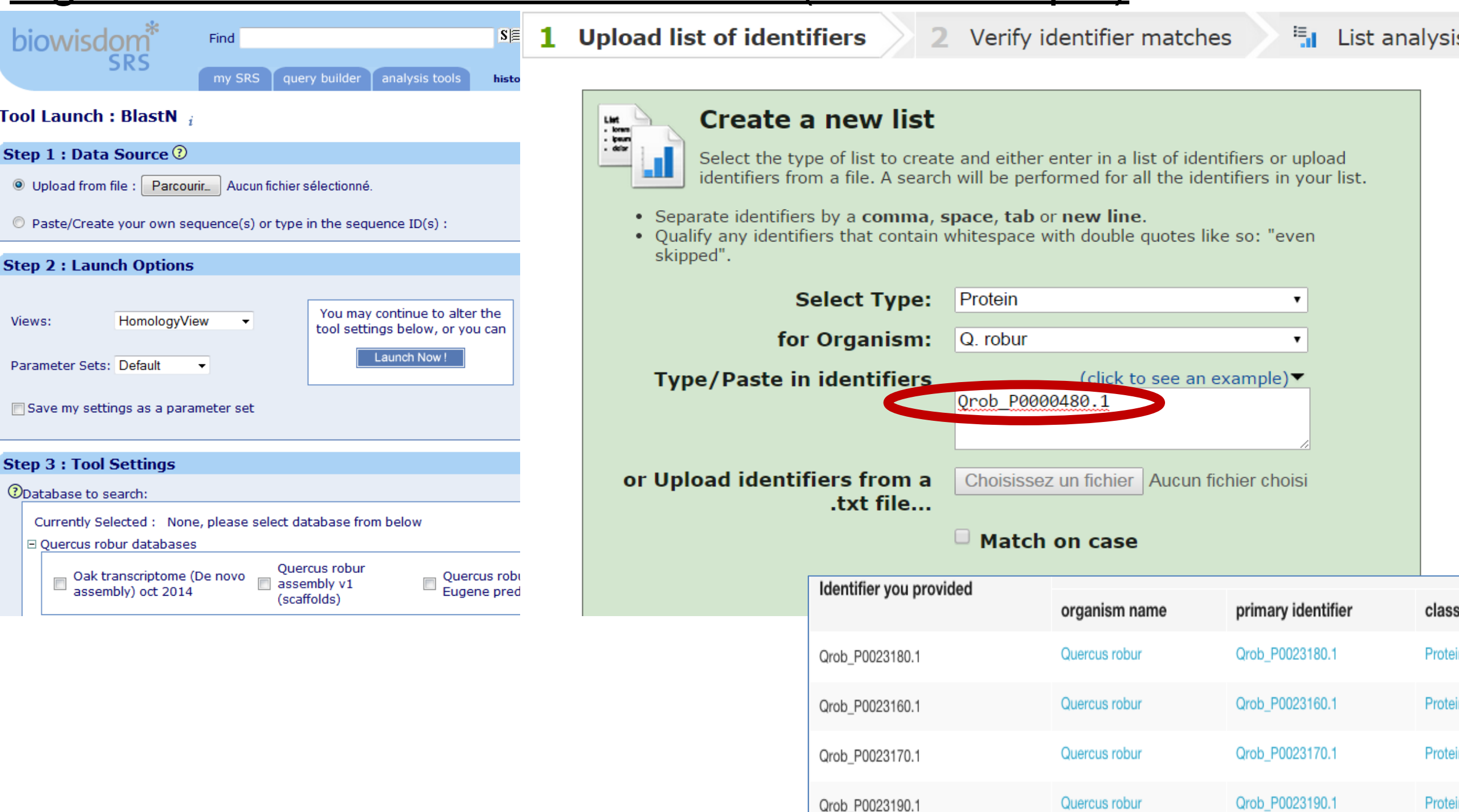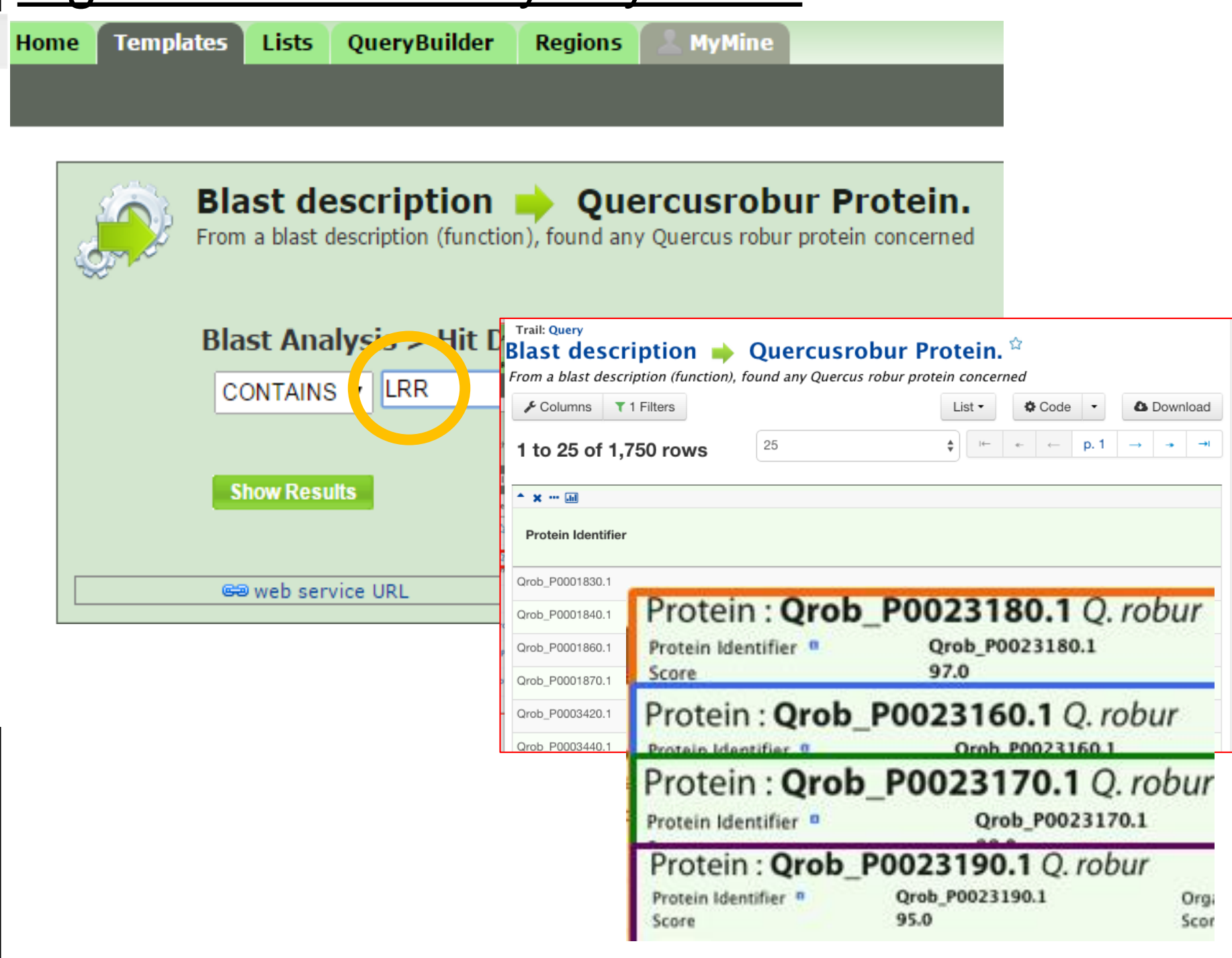### Figure 2b: Search a list of identifiers (blast example)



### Figure 2c: Search by keywords



## Architecture

To curate an annotation, we need to combine different sources of evidence altogether. Those information provided by structural and functional genome annotation pipeline (figure 1) are available through two powerful resources. Intermine is used to gather and give access to functional annotation using several ways of querying data (figure 2a). With a list of identifiers, obtained by blast comparison against the newly predicted set of genes (figure 2b) or using keywords based search (figure 2c), we can retrieve all functional annotation displayed through a gene card (figure 2d). Those cards linked functional annotation and structural annotation by featuring a WebApollo (figure 3). WebApollo add several tools to the jbrowse (figure 3a), which allow gene curation like merging predicted transcripts into new gene (figure 3b). This architecture evolution would be to have an automated and continuous update of all the manual annotation made through this system.
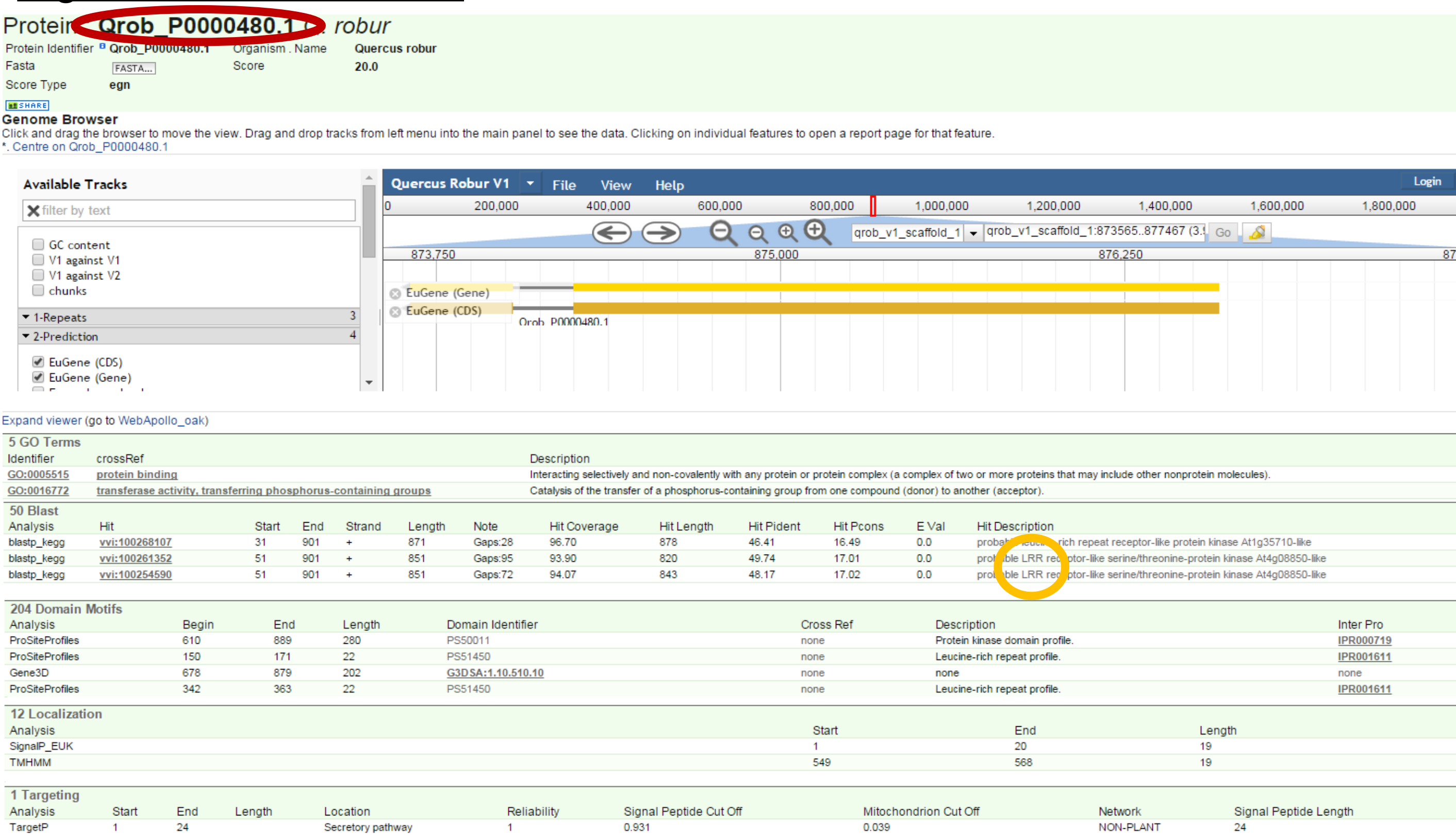
### Figure 2d: Protein card
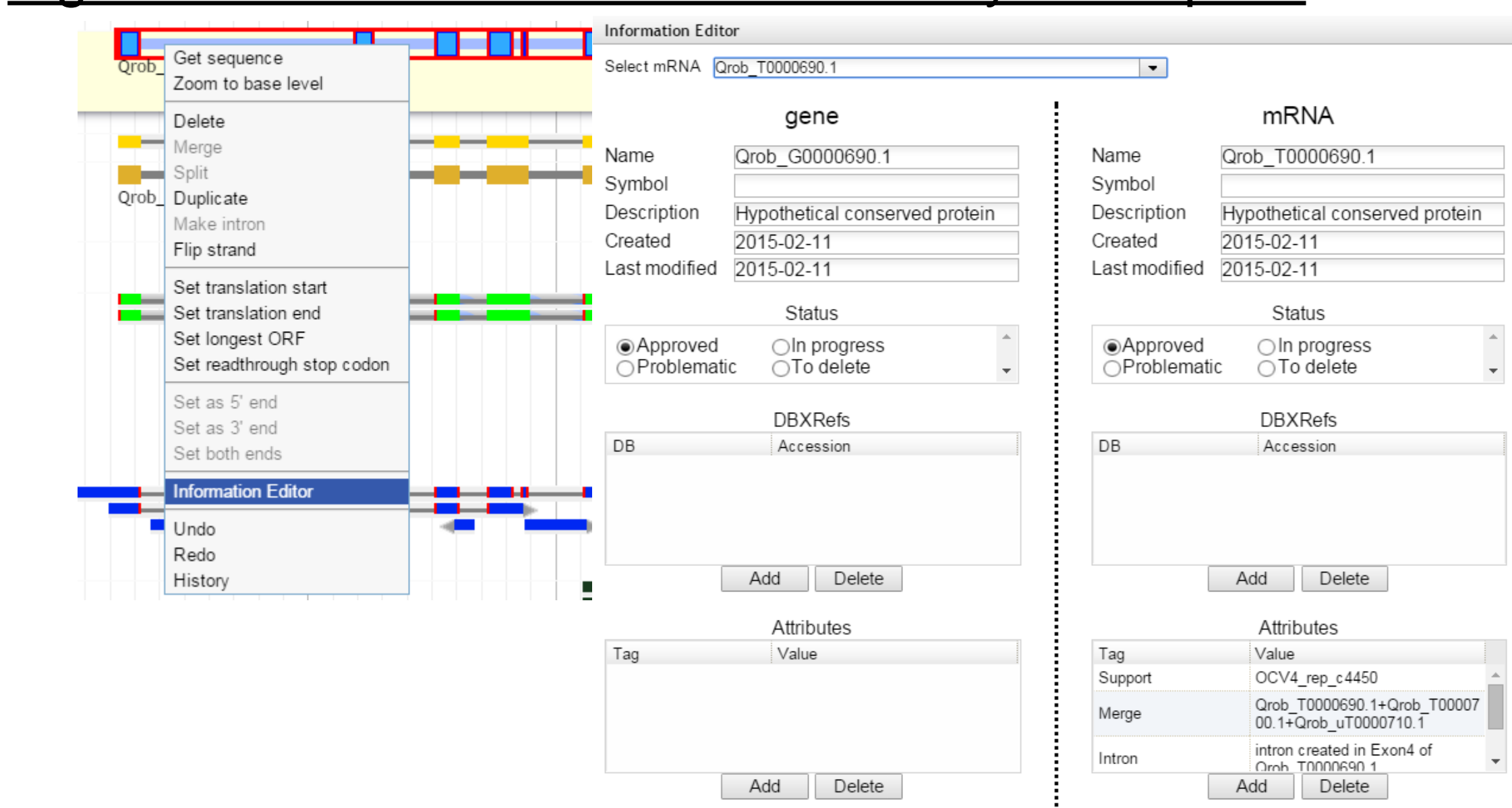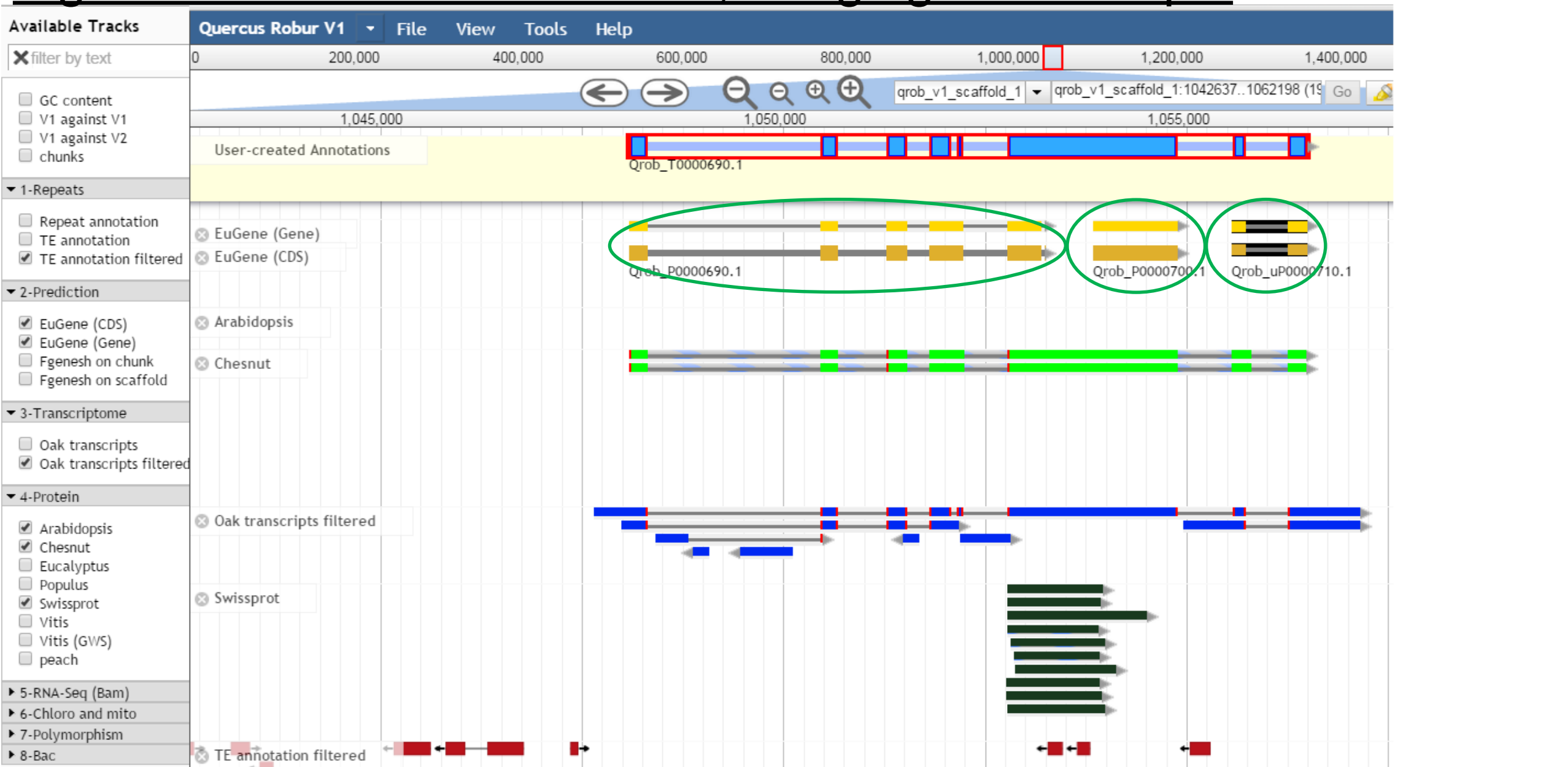


### Figure 3a: Visualization and curation by WebApollo



### Figure 3b: Curation used-case, merging 3 transcripts