

UniRule - Increasing Annotation Depth of Unreviewed Protein Entries in UniProtKB

The UniProt Knowledgebase

The mission of UniProt is to provide the scientific community with a comprehensive, high-quality and freely accessible resource of protein sequence and functional information. One of the databases provided by UniProt is the UniProt Knowledgebase (UniProtKB) which contains expertly reviewed (UniProtKB/Swiss-Prot) and unreviewed (UniProtKB/TrEMBL) records.

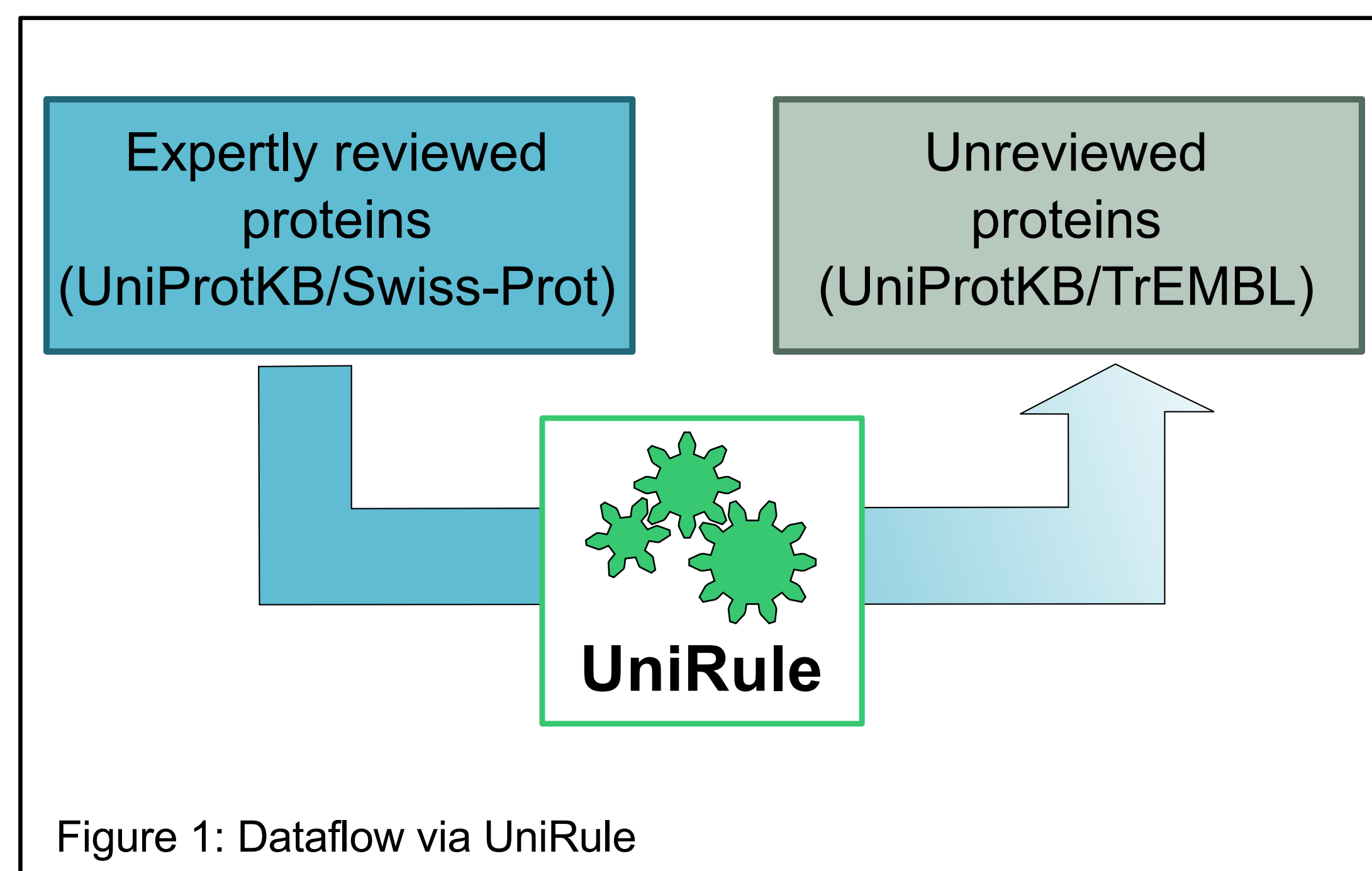


Figure 1: Dataflow via UniRule

Challenges of data growth and a solution

Currently, reviewed records constitute only about 1% of UniProtKB; expert curation is time-intensive and most published experimental data focuses on a rather limited range of model organisms. At the same time, the number of unreviewed records in UniProtKB (UniProtKB/TrEMBL) is growing continuously, yet for a large proportion of these records there is no experimental data available.

The UniRule system leverages expert curation for the automatic annotation of unreviewed UniProtKB records. It consists of manually created annotation rules that specify functional annotations and the conditions which must be satisfied for them to be applied, such as taxonomic scope, family membership as defined by InterPro and the presence of specific sequence features.

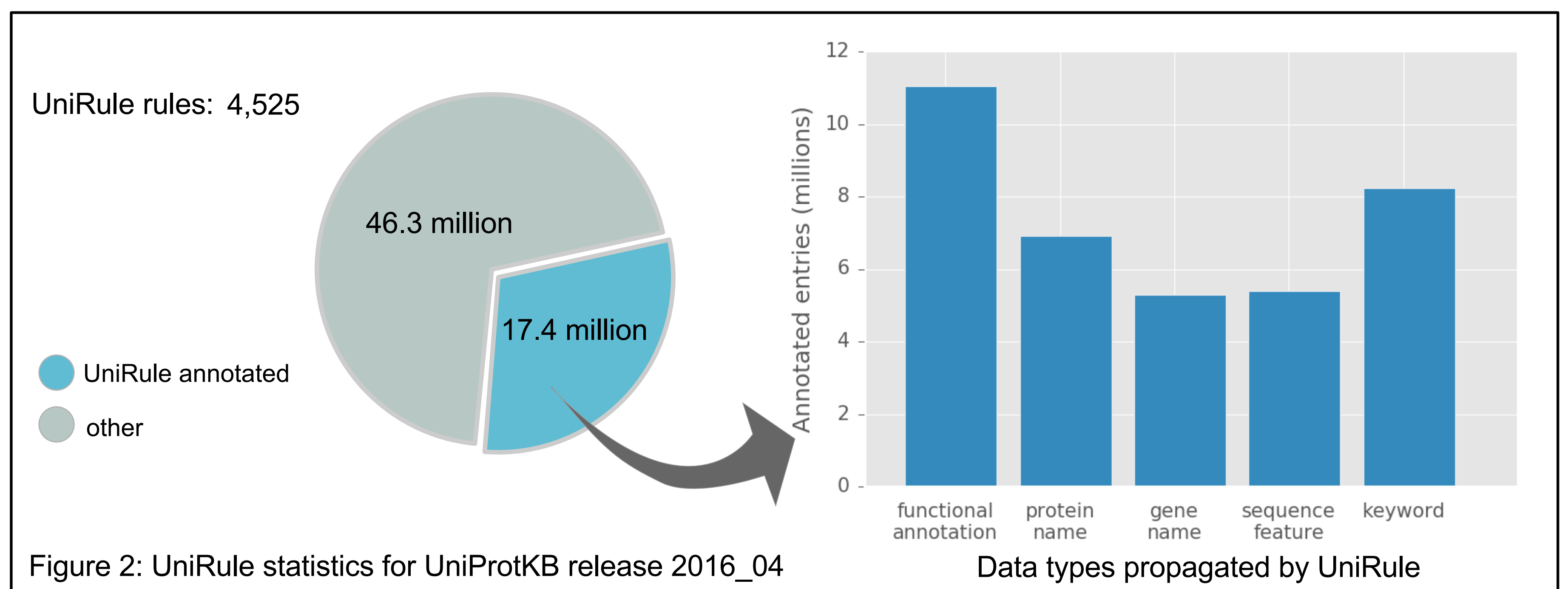


Figure 2: UniRule statistics for UniProtKB release 2016_04

Rule creation

Rule creation is part of an integrated workflow that starts with the expert curation of UniProtKB/Swiss-Prot records. Protein names, gene names, functional annotation, GO terms, keywords and sequence features are selected and combined with InterPro signatures and taxonomic restrictions to create a rule. Protein family signatures can be custom-built if necessary.

A. Flavin transferase that catalyzes the transfer of the FMN moiety of FAD and its covalent binding to the hydroxyl group of a threonine residue in a target flavoprotein. UniRule annotation

B. GO - Molecular function¹

- metal ion binding Source: UniProtKB-UniRule
- transferase activity Source: UniProtKB-UniRule

Figure 3: Evidence tags pointing to a UniRule rule. A, Annotation describing protein function. B, Associated Gene ontology terms

Rule maintenance

A range of statistics are used to evaluate the performance of rules (confidence, sensitivity, coverage). Statistics are recalculated after every change made to a rule and with every release. Rules performing suboptimally are flagged and are then reviewed by curators, thus ensuring the timely correction of any inaccuracies.

Rule application

Rules are reapplied with each UniProtKB release, keeping the propagated annotation up-to-date. On the UniProtKB website, information added by automatic annotation procedures is clearly highlighted as such using evidence tags (see Figure 3). The tags can also be used as search terms to specifically search for and/or filter out annotation added by a rule.

UniRule tool

The UniRule tool is a web interface used by curators for both the creation and maintenance of rules. It holds all the rules in the UniRule system and allows the editing of rules individually or in batch mode.

Rules can be built either step by step with the tool providing validation and assistance or they can be imported using an XML exchange format. The tool also exports to this format. Statistics on rule performance are calculated against all reviewed UniProtKB entries on the fly. This way, the effects of individual changes to a rule can be assessed. In case the effects are found to be detrimental to rule performance, a 'history' functionality allows previous versions to be reinstated.

People

R. Antunes¹, C. Arighi³, D. Baratin², A. Bridge², E. Coudert², B.A. Cuche², E. de Castro², J.S. Garavelli³, E. Hatton-Ellis¹, G. Keller², K. Laiho³, M. Martin¹, A. MacDougall¹, C. O'Donovan¹, I. Pedruzzi², K. Pichler¹, D. Poggioli¹, L. Pureza¹, N. Redaschi², A. Renaux¹, C. Rivoire², C.R. Vinayaka³, V. Volynkin¹, Q. Wang³, L.-S. Yeh³, H. Zellner¹

UniRule: UR000236406

Source ID: RU363002

View all proteins annotated by this rule Remove highlights

If a protein meets these conditions...¹

Common conditions

Matches Pfam signature PF02424
sequence length = 200 - 500
taxon = Bacteria

Special conditions

taxon = Proteobacteria, Chlamydiae, Spirochaetes

Predicted signal

... then these annotations are applied¹

Protein name¹ (then)

Recommended name:
FAD:protein FMN transferase (EC:2.7.1.180)

Sequence similarities¹

Belongs to the ApbE family.

Catalytic activity¹

FAD + [protein]-L-threonine = [protein]-FMN-L-threonine + AMP.

Cofactor¹

Mg²⁺

Subcellular location¹

Cell inner membrane

Function¹

Flavin transferase that catalyzes the transfer of the FMN moiety of FAD and its covalent binding to the hydroxyl group of a threonine residue in a target flavoprotein.

Chain¹

@CHAIN_NAME@¹ (to residues corresponding to positions @TO|+1@¹ - @CTER@¹)

Signal peptide¹

(to residues corresponding to positions @NTER@¹ - @TO@¹)

Keywords¹

FAD
Flavoprotein
Magnesium
Metal-binding
Transferase
Cell inner membrane
Lipoprotein
Signal

GO (Gene Ontology) terms¹

GO:0016740 transferase activity
GO:0046872 metal ion binding

Figure 4: Screenshot of a UniRule rule on uniprot.org

UniRule data in UniProtKB

UniRule-generated annotation in UniProtKB and all rules underlying the data can be explored. Clicking on conditions or annotations will highlight links between the two, thus explaining the rule's logic.

Figure 4 shows an example of a UniRule rule as it is displayed on uniprot.org including rule identifiers for reference and a quick link to a tabular view of all unreviewed entries touched by this rule. The rule defines the overall sequence space it covers using a Pfam signature, a taxonomic restriction to bacteria and constraints on the length of target sequences ('common conditions'). The bulk of the annotations is propagated if the common conditions are met (this is highlighted in the screenshot). Special conditions apply for propagation of the subcellular location comment. Furthermore, the rule calls a predictive algorithm (SignalP) for annotating potential signal peptides.