



AUTOMATISATION IN UNIPROTKB / SWISS-PROT ANNOTATION: NEW RULES AND TOOLS

[1] The HAMAP team, [1] The PROSITE team, [1] Edouard de Castro, [1] <u>Alan Bridge</u>, and the UniProt Consortium.
[1] Swiss-Prot Group, Swiss Institute of Bioinformatics, CMU, 1 Michel Servet, 1211 Geneva 4.

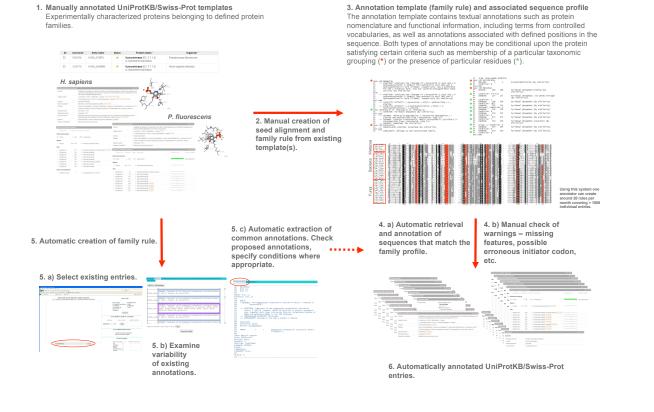
Introduction

The development of next generation sequencing technologies promises a massive increase in the rate of submission of new protein sequences to sequence databases such as the Universal Protein Resource Knowledge Base, UniProtKB. At UniProtKB/Swiss-Prot we propose to meet this challenge by continuing to expand and develop systems for the automatic propagation of existing annotation to newly submitted protein sequences. These developments will promote the standardization of ortholog annotation both across and within kingdoms and significantly enhance our ability to accurately annotate new protein sequences which are being produced at an ever increasing rate.

The HAMAP pipeline and automatic annotation of eukaryotic proteins

The existing prototype for automatic annotation in Swiss-Prot is the HAMAP project (High-quality Automated and Manual Annotation of microbial Proteomes http://www.expasy.org/sprot/hamap/index.html), which specifically targets microbial proteins. We have recently modified the existing HAMAP annotation pipeline to allow the automatic annotation of eukaryotic proteins. Steps 1 to 4 below illustrate a typical workflow for the construction of an **annotation template** or **family rule**. Such rules serve as the source of annotations for new entries that match an associated sequence profile. Steps 5a-c describe the application of a new tool for the automatic creation of such family rules from an existing set of manually annotated UniProtKB/Swiss-Prot entries.

We have identified eukaryotic proteins that are candidates for automatic annotation in two ways. First, during a recent annotation marathon targeting conserved proteins of the slime mold *Dictyostelium discoideum*, we manually annotated over 1000 proteins with orthologs in at least two other major eukaryotic phyla (*Homo sapiens/Mus musculus* and/or *Drosophila* melanogaster and/or *Sacharomyces cerevisiae/Schizosaccharomyces pombe)*. *D. discoideum* separated from metazoa and fungi prior to the metazoa:fungi split, and so annotation of *D. discoideum* orthologs provides an ideal opportunity for the identification of common annotations suitable for transfer to other eukaryotic orthologs. *Second*, we analyzed the coverage of eukaryotic proteins by existing HAMAP annotation rules. HAMAP includes 1551 family rules covering 214838 proteins (UniProtKB/Swiss-Prot Release 56.7 of 20-Jan-2009). Of these rules, 167 potentially match vertebrate proteins, while 228 match fungi and 147 match arthropods. The taxonomic scope of these existing HAMAP rules could conceivably be expanded to cover eukaryotes.



IniProt is mainly supported by the National Institutes of Health (NIH) grant 2 001 HG02712-04. Additional support for the EBI workernent in UniProt comes from the European Commission (EC)/S FELCS grant (02028/III) and from the NIH grant R01HG02273-01. Swiss-Prot activities at the SIB are supported by the Swiss Federal Government through the Federal Offic ducation and Science and the European Commission contract FELCS (021902RII3). PIR activities are also support for the Hill grants and contract HHSV86260400061C, NCI-caBIG, and 1R01GM080646-01. and the National Science Foundation NSF) grant IIS-0430743. PROSITE is supported by the Swiss National Science Foundation, grant n. 315280-116864 and the nterPo rant n. 121037 from the European Linon.

Contact help@uniprot.org www.uniprot.org

See also: www.isb-sib.ch