# Improved selection of canonical proteins for reference proteomes

Giuseppe Insana[1], Maria J. Martin[1], William R. Pearson[2]

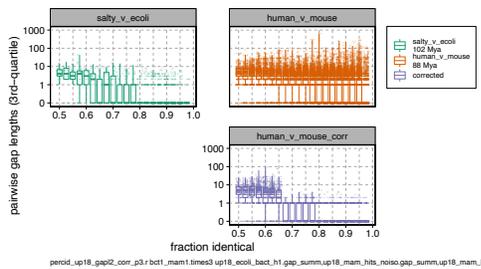(1) EMBL-EBI, Hinxton, UK  (2) U. of Virginia, Charlottesville, VA USA

## Introduction

The **Reference Proteomes** dataset seeks to provide complete proteomes for an evolutionarily diverse, less redundant, set of organisms.
As higher eukaryotes often encode multiple isoforms of a protein from a single gene, the Reference Proteomes pipeline selects a single representative ('**canonical**') sequence. UniProt identifies canonical isoforms using a *'Gene-Centric'* approach: proteins are grouped by gene-identifier and for each gene a single protein sequence is chosen.
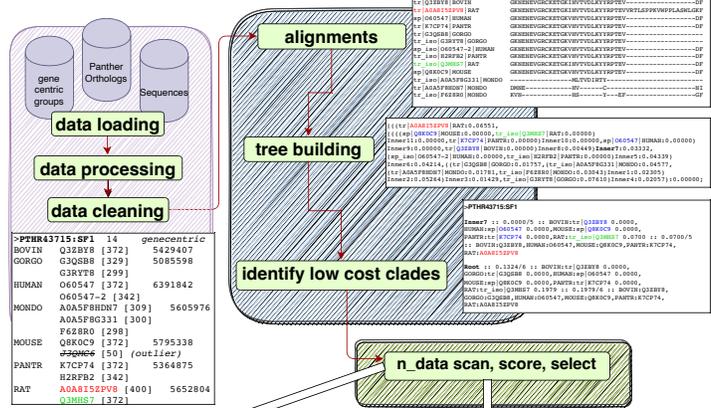
For unreviewed (UniProtKB/TrEMBL) protein sequences (and for some reviewed sequences), the longest sequence in the Gene-Centric group is chosen as canonical. This can create inconsistencies, selecting sequences with dramatically different lengths as canonical for orthologous genes. Biologically, it is unlikely that orthologous mammalian proteins differ greatly in length, but this happens about **10%** of the time for the **8 mammals in the Quest for Orthologs** set.

The **Ortho2tree** data pipeline examines **Gene-Centric** canonical and isoform sequences from sets of orthologous proteins (from **PantherDB**), and suggests replacements for canonicals that have lengths very different from closely related orthologs.
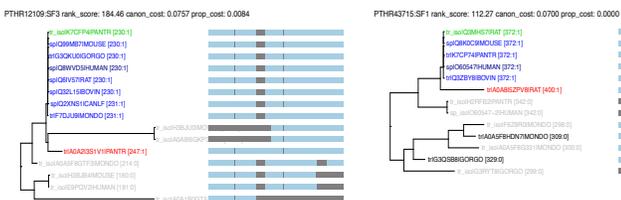
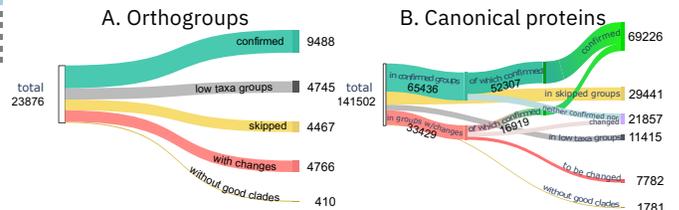### Closely related (>80% id) bacterial proteins have few gaps, unlike mammalian proteins



### Gaps in vertebrate ortholog protein alignments reflect isoform choice



Correction of chimp (A) or rat(B) canonical isoform selection using gap-based distance trees. Blue, deep-blue (MANE agreement), and red denote current canonical isoforms, red isoforms should not be canonical, they should be replaced by green isoforms.

### The ortho2tree pipeline combines gene-centric isoforms of orthologs to build length-consistent canonical orthologs



### Correction of QfO mammalian proteins



A. Orthogroups

B. Canonical proteins

## Methods & Results

Gap-distance trees are constructed from multiple sequence alignments of the sequences in each orthogroup. Clades with a low cost and a diverse set of taxa are identified from these trees. A weighting function is then used to seek the clades with more SwissProt entries and more canonicals, prioritising annotation from well characterized proteomes.

The pipeline suggests replacements of canonicals in the same Gene-Centric group and indicates when there is agreement between those suggestions and MANE (Matched Annotation from NCBI and EMBL-EBI) assignments.

For the ~140,000 proteins in ~24,000 Panther orthogroups from the eight mammalian proteomes (human, chimp, gorilla, mouse, rat, dog, cow and opossum), ortho2tree proposed 7782 canonical changes, while confirming 69,226 canonical assignments. When Ortho2tree suggested changing a HUMAN canonical with a MANE assignment, MANE suggested the same change 80% of the time; when no change to the HUMAN canonical was proposed, Ortho2tree and MANE agreement was 95%.

Ortho2tree can reduce canonical assignment errors among orthologous sequences that are more than 50% identical, such as sequences from vertebrates or higher plants.

### Complete results available online



**fasta.bioch.virginia.edu/ortho2tree/**