Databases and ontologies

# OMAMO: orthology-based model organism selection

**Alina Nicheperovich** [1,], **Adrian M. Altenhoff** [2,4], **Christophe Dessimoz** [3,4,5,6,*] **and Sina Majidian** [3,4,*]

[1] Department of Structural and Molecular Biology, University College London, London WC1E, UK, [2] Department of Computer Science, ETH, 8092 Zurich, Switzerland, [3] Department of Computational Biology, University of Lausanne, 1015 Lausanne, Switzerland, [4] SIB Swiss Institute of Bioinformatics, 1015 Lausanne, Switzerland, [5], Department of Computer Science, University College London, London WC1E 6BT, UK and [6] Department of Genetics, Evolution and Environment, University College London, London WC1E, UK

*To whom correspondence should be addressed.

Associate Editor: XXXXXXX

Received on XXXXX; revised on XXXXX; accepted on XXXXX

## Abstract

**Summary:** The conservation of pathways and genes across species has allowed scientists to use non-human model organisms to gain a deeper understanding of human biology. However, the use of traditional model systems such as mice, rats, and zebrafish is costly, time-consuming and increasingly raises ethical concerns, which highlights the need to search for less complex model organisms. Existing tools only focus on the few well-studied model systems, most of which are higher animals. To address these issues, we have developed **O**rthologous **Ma**trix and **M**odel **O**rganisms, a software and a website that provide the user with the best simple organism for research into a biological process of interest based on orthologous relationships between the human and the species. The outputs provided by the database were supported by a systematic literature review.

**Availability and implementation:** https://omabrowser.org/omamo/, https://github.com/DessimozLab/omamo

**Contact:** christophe.dessimoz@unil.ch and sina.majidian@unil.ch

**Supplementary information:** Supplementary data are available at *Bioinformatics* online.

## 1 Introduction

Model organisms are non-human species used in human biomedical research to study development, gene regulation, and other cellular processes because they are relatively fast-growing, inexpensive, and easy to manipulate. Most importantly, their use has been possible due to the evolutionary conservation of biological processes (Wangler *et al.*, 2017). Fast-moving progress in comparative genomics has allowed scientists to identify these evolutionary relationships by inferring human orthologs, genes that have diverged due to speciation (Fitch, 1970). Since orthologous genes tend to be functionally conserved and have common gene expression patterns, they are a better basis for model organism selection than other subtypes of homologs, which tend to functionally diverge faster (Altenhoff *et al.*, 2012; Zheng-Bradley *et al.*, 2010).

Currently used model organisms range from bacteria to complex mammals. The scientific community, however, aims to reduce the use of animals in research due to ethical implications, opting to use less complex organisms where possible. Currently available databases include MARVVEL (Wang *et al.*, 2019), the Alliance of Genome Resources portal (Alliance of Genome Resources Consortium, 2000), and MORPHIN (Hwang *et al.*, 2014). They focus on five to nine 'traditional' model organisms, most of which are higher organisms like mouse, rat and zebrafish. Moreover, their scope is restricted to human disease-related research. The only unicellular organisms considered in these databases are fission and budding yeast, whilst abundance of unicellular species in nature and their unique features make it difficult to find other non-complex model organisms for a biological process of interest.

To address the challenges above, we created an orthology-based database tool OMAMO alongside a user-friendly website that helps to select the best non-complex model organism for a biological process. Because the majority of species in the database have not been considered as model systems in the past, OMAMO has the potential to extend the set of organisms used in human biomedical research.

**1**

picture(0,0)(-35,0)(1,0)30 (0,35)(0,-1)30 picture                    picture(0,0)(35,0)(-1,0)30 (0,35)(0,-1)30 picture

## 2 Methods

OMAMO takes advantage of the OMA database of orthologous genes. For a given biological process, the output presents a list of potential model organisms ranked based on their orthologous relationships with human.

For each species, pyOMA library was used to extract human orthologs (Altenhoff *et al.*, 2021). For each ortholog, pyOMA was used to retrieve Gene Ontology (GO) terms, which provide information about the gene product and can represent one of the following three aspects: molecular function, cellular component and biological process (Gene Ontology Consortium, 2021). Some GO terms are general (e.g. 'cell division'), whilst others are more specific ('G2/M transition of mitotic cycle'). To quantify specificity of a GO term, we used information content (IC) calculated as *-log(p)* where *p* is its empirical frequency in the UniProt database (Pesquita, 2017), hence more specific GO terms have a higher IC value. The IC values were used to calculate functional similarity for each orthologous pair (Supplementary Section 1).

Orthologous pairs with functional similarity of $< 0.05$ were discarded. This aims to reduce the number of orthologs that only share general GO terms in the output. Consequently, gene pairs from a given species were grouped according to the biological process GO term they share. To maintain sufficient specificity in functional similarity considered, only GO terms with the IC value of $\geq 5$ were kept. Finally, for each biological process GO term, species were ranked based on a scoring system, which takes into account the number of orthologs relevant to the biological process and average functional similarity across the genes.

We developed a freely accessible website for OMA (Fig.1), with the source code publically available. Out of the 50 species currently present in OMAMO, 31 are unicellular eukaryotes and the rest are bacteria. We suggest at least one model organism for 4620 out of 28,923 available biological GO terms (Gene Ontology Consortium, 2021). Since OMAMO is integrated in the OMA, it will be updated alongside the browser,meaning that the set of organisms will continue to grow and the database will include the latest GO annotations.

To validate our results, we referred to experimental evidence through a systematic literature search on PubMed (Supplementary Section 2). The top five review articles on three of the most well-studied organisms in OMA (*D. discoideum, N. crassa, S. pombe*) published in 2010-2021 were selected from the search output. Out of all biological processes which have been studied in one of the three organisms, the species of interest was in the top 5 model organism candidates in 42.6% of respective searches in OMAMO (Supplementary Section 2). This indicates that our algorithm is well supported by experimental data found in the literature.

## 3 Discussion

OMAMO is a freely-available database which aims to help scientists exploit alternative model species for human biomedical research. With the limited number of presently used model systems, the scientific community can now benefit from using other organisms, some of which could become model systems for processes that have previously only been studied in animals, leading to reduction in their use in experimental research. Moreover, this is the first database that provides such a wide range of potential model organisms. Due to the lack of literature on using species presented in OMAMO, the validation of results proved to be challenging. The following step for output validation would be to utilise proposed model species as model systems in wet-lab experiments. In the future, we plan to greatly expand the set of species and improve the scoring system by considering sequence similarity, conservation of protein structure and reproduction time. Additionally, we hope to provide unicellular model organisms for studying species other than human, for example animals for veterinary science research.



**Fig. 1.** Website interface. (A) The main browser page of OMAMO. The user can search a GO term ('DNA repair') or a GO ID (0006281). (B) The output page gives a list of species ranked based on the score, but the user has the option to sort the output based on the total number of orthologs or the average functional similarity by clicking on the up-down sorting icon. The user can view orthologs by clicking on the '+' button.

## References

Alliance of Genome Resources Consortium (2020) Alliance of Genome Resources Portal: unified model organism research platform. *Nucleic Acids Res.*, **48**, D650–D658.

Altenhoff,A.M. *et al* (2021) OMA orthology in 2021: website overhaul, conserved isoforms, ancestral gene order and more. *Nucleic Acids Res.*, **49**, D373–D379.

Altenhoff,A.M *et al.* (2012) Resolving the ortholog conjecture: orthologs tend to be weakly, but significantly, more similar in function than paralogs. *PLoS Comput. Biol.*, **8**, e1002514.

Fitch,W.M. (1970) Distinguishing homologous from analogous proteins. *Syst. Zool.*, **19**, 99-113.

Gene Ontology Consortium (2021) The Gene Ontology resource: enriching a GOld mine. *Nucleic Acids Res.*, **49**, D325–D334.

Hwang,S. *et al.* (2014) MORPHIN: a web tool for human disease research by projecting model organism biology onto a human integrated gene network. *Nucleic Acids Res.*, **42**, W147–53.

Pesquita,C. (2017) Semantic Similarity in the Gene Ontology. *Methods Mol. Biol.*, **1446**, 161–173.

Wang,J. *et al.* (2019) Navigating MARRVEL, a Web-Based Tool that Integrates Human Genomics and Model Organism Genetics Information. *J. Vis. Exp.*

Wangler,M.F. *et al.* (2017) Model Organisms Facilitate Rare Disease Diagnosis and Therapeutic Research. *Genetics*, **207**, 9-27.

Zheng-Bradley,X. *et al.* (2010) Large scale comparison of global gene expression patterns in human and mouse. *Genome Biol.*, **11**, R124.

picture(0,0)(-35,0)(1,0)30 (0,-35)(0,1)30 picture                    picture(0,0)(35,0)(-1,0)30 (0,-35)(0,1)30 picture