



HAL
open science

OrthoInspector 3.0: open portal for comparative genomics

Yannis Nevers, Arnaud Kress, Audrey Defosset, Raymond Ripp, Benjamin Linard, Julie D Thompson, Olivier Poch, Odile Lecompte

► To cite this version:

Yannis Nevers, Arnaud Kress, Audrey Defosset, Raymond Ripp, Benjamin Linard, et al.. OrthoInspector 3.0: open portal for comparative genomics. *Nucleic Acids Research*, 2019, 47 (D1), pp.D411-D418. 10.1093/nar/gky1068 . hal-01984857

HAL Id: hal-01984857

<https://hal.science/hal-01984857>

Submitted on 17 Jan 2019

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial 4.0 International License

OrthoInspector 3.0: open portal for comparative genomics

Yannis Nevers¹, Arnaud Kress¹, Audrey Defosset¹, Raymond Ripp¹, Benjamin Linard^{2,3,4}, Julie D. Thompson¹, Olivier Poch¹ and Odile Lecompte^{1,*}

¹ Department of Computer Science, ICube, UMR 7357, University of Strasbourg, CNRS, Fédération de Médecine Translationnelle de Strasbourg, Strasbourg, France

² LIRMM, Univ Montpellier, CNRS, Montpellier, France

³ ISEM, Univ Montpellier, CNRS, IRD, EPHE, CIRAD, INRAP, Montpellier, France

⁴ AGAP, Univ Montpellier, CIRAD, INRA, Montpellier Supagro, Montpellier, France

* To whom correspondence should be addressed. Tel: +33 (0)3 68 85 32 96 ; Email: odile.lecompte@unistra.fr

ABSTRACT

OrthoInspector is one of the leading software suites for orthology relations inference. In this paper, we describe a major redesign of the OrthoInspector online resource along with a significant increase in the number of species: 4753 organisms are now covered across the three domains of life, making OrthoInspector the most exhaustive orthology resource to date in terms of covered cellular species. The new website integrates original data exploration and visualization tools in an ergonomic interface. Distributions of protein orthologs are represented by heatmaps summarizing their evolutionary histories, and proteins with similar profiles can be directly accessed. Two novel tools have been implemented for comparative genomics: a phylogenetic profile search that can be used to find proteins with a specific presence-absence profile and investigate their functions and, inversely, a GO profiling tool aimed at deciphering evolutionary histories of molecular functions, processes or cell components. In addition to the re-designed website, the OrthoInspector resource now provides a REST interface for programmatic access. OrthoInspector 3.0 is available at <http://lbgi.fr/orthoinspectorv3>.

INTRODUCTION

Genes descending from a common ancestor, or homologs, are commonly divided into two classes: orthologs, that are derived from a speciation event, and paralogs, that are derived from a duplication event (1). According to the ortholog conjecture (2), which has been debated recently but still holds (3, 4), orthologs generally conserve the same function in distinct species while paralogs can evolve different or specialized functions. Furthermore, a discrimination between outparalogs and inparalogs is needed when studying evolutionary and functional relationships between proteins (5). Outparalogs are produced by a duplication event anterior to a given speciation event, while inparalogs result from a 'recent' duplication, posterior to a speciation event. Thus, inparalogs in one species are assumed to be relatively close to each other and are considered co-orthologs to their counterparts in another species deriving from the considered speciation event.

These notions are key principles in current biology and inferring the true orthologs or co-orthologs of proteins is crucial for comparative genomics and molecular biology. For example, it is essential in the transfer of data from experimental studies between species, thus making it possible to study human health in model organisms. It is also the keystone of phylogenetic profiling, an approach that exploits the presence and absence of protein orthologs across multiple species (6). The method is based on the principle that two proteins that interact or are involved in the same biological process tend to be conserved and lost together (7). Applications of phylogenetic profiling include protein-protein interaction inference and genotype-phenotype correlation as genes associated with a certain phenotypic trait tend to have a profile correlated with that trait's phylogenetic distribution (8).

1
2
3 More than thirty resources have been developed to address the challenges of orthologous relation
4 inference and community efforts have been directed towards standardization and benchmarking of these
5 resources, in the form of the Quest for Orthologs consortium (9). OrthoInspector (10, 11) was shown to be
6 one of the three most balanced methods of orthology inference in terms of precision and recall in a
7 standardized benchmarking test (12) and performed well in other comparative studies (13). The previous
8 release of OrthoInspector (11) provided two precomputed databases (Prokaryotes and Eukaryotes) that
9 could be queried from its website, however since the last release the number of available annotated
10 genomes has significantly increased and standards for web interfaces have evolved.

11
12 Here, we present the third release of OrthoInspector that includes a number of important developments.
13 First, we report a major increase in the number of species represented in the OrthoInspector precomputed
14 databases across the three domains of cellular life, including both in-domain and cross-domain relations,
15 making the OrthoInspector databases the most exhaustive orthology resource to date in terms of covered
16 species. Second, to manage the massive increase of data, the OrthoInspector website has been entirely
17 redesigned to provide a streamlined and intuitive experience for users, including a summary visualization
18 of ortholog distributions and novel tools allowing powerful comparative genomics analyses.

19 RESULTS

20 Improved coverage of the tree of life

21 *Proteome selection*

22
23 When designing the OrthoInspector databases, we focused on providing a broad coverage of the tree of
24 life, with a selection of organisms that are representative of the taxonomic diversity. In order to meet this
25 goal, we used the Uniprot Reference Proteomes (14), which result from an effort to efficiently sample the
26 tree of life and limit redundancy. Incomplete genomes, mispredicted or fragmentary protein sequences
27 constitute an important source of errors in orthology inference. Therefore, we used a combination of filters
28 (see supplementary materials and methods) to exclude proteomes with abnormally small proteome size, a
29 high proportion of small proteins or of proteins that do not start with a methionine.

30
31 Starting from the 5443 Reference Proteomes, the quality filtering step resulted in the exclusion of 690
32 proteomes (13%). The percentages of excluded proteomes were similar across domains: 119 out of 830
33 eukaryotes (14%), 537 out of the 4400 Bacteria (12%) and 34 out of 213 Archaea (16%). In one case, we
34 privileged the coverage of the tree of life over quality measures and kept the proteome of *Lokiarchaeum*
35 *sp. GC14_75* owing to the general interest for representatives of the Asgard group in comparative
36 genomics (15, 16).

37 *Database architecture*

38
39 The OrthoInspector 3.0 databases cover 4753 organisms (+144% compared with the previous release):
40 3863 bacteria (+146%), 711 (+174%) eukaryotes and 179 archaea (+49%) (Figure 1). This is, to our
41 knowledge, the widest coverage available for an orthology inference resource in terms of cellular species.

42
43 The database architecture is designed to cover the essential use cases for orthology data. It relies on
44 three main databases, one for each domain of life. Each database provides all the orthologous relations
45 between proteins of each species within the domain. This exhaustive coverage of each domain is suitable
46 for fine grained studies, as it provides a good resolution at low taxonomic levels.

47
48 We designed a fourth database to provide orthologous relationships across a wider evolutionary spectrum
49 and specifically, to cover the three domains simultaneously. To facilitate handling and interpretation of
50 these cross-domain comparisons, we defined a subset of significant species that we will refer to as “model
51 species” (see supplementary Table 1). We selected these species according to their importance in the
52 biological field (e.g. model species such as *Mus musculus* or *Caenorhabditis elegans*) and/or to ensure a
53 good taxonomic sampling (Figure 1). This selection corresponds to 317 species: 144 eukaryotes, 142
54 bacteria and 31 archaea.

OrthoInspector can thus be used to find intra-domain orthologs in a large number of species and to find inter-domain orthologs in fewer, well-studied, species. Users interested in orthology relationships between non-model species from different domains can find them by transitivity, by first finding orthologs in close ‘model species’. This original implementation involving the co-existence of databases with different levels of granularity implies that orthologs can be found in all our available species without the huge computational burden a ‘full’ inference would require.

Complete information about the database content is available in supplementary Table 1 and in the database tab on the website.

A new information design

To cope with the massive increase in the number of species available in the OrthoInspector databases and the corresponding increase in the number of orthology relationships, we implemented a new website interface providing a smooth navigation in the new datasets.

Access to protein entries

The OrthoInspector website offers two main ways to access the data: by protein identifier and by sequence similarity searches.

The protein identifier search is accessible from the main page, or anywhere on the site using the navigation bar. The user should define the appropriate database by selecting the domain of life of the query protein. Typing in the search bar triggers autocompletion and dynamically proposes a list of clickable protein entries available in the selected OrthoInspector database. The identifier search currently supports both Uniprot identifiers and Uniprot access numbers.

A sequence similarity search is also available from the OrthoInspector webpage or by selecting ‘BLAST search’ on the database tab. This launches a BLASTp (17, 18) search against all protein sequences in the OrthoInspector databases. The result is a formatted BLAST output of the 50 best hits along with their corresponding local alignments and links to the corresponding protein pages in OrthoInspector.

Protein page

The data in OrthoInspector can be explored from protein pages. The protein page header gives a quick summary of the protein (gene name, description, organism). All Gene Ontology (19) terms associated with this protein are displayed in an extendable panel when available, as well as the protein sequence and a schematic view of InterPro (20) domains found in the protein. The protein page is the core section of the website architecture and provides access to orthology relations, taxonomic distribution and proteins with similar distribution (detailed below).

Orthology data

Orthologous relationships are presented in the ‘Orthologs and taxonomic distribution’ section of the protein page. A menu allows users to choose display options, depending on their needs:

- **Domain’s model organisms:** only orthologs found in the “model organisms” of Eukaryotes, Bacteria or Archaea are shown in this tab. This view is used to find orthologs in popular species and avoids overwhelming the user with superfluous information. The page shown by default should meet the requirements of most users and thus serves as a suitable entry point.
- **Whole domain:** orthologs in all species of the in-domain databases are shown in this tab. This exhaustive view is suitable for an in-depth exploration of intra-domain relationships.
- **Three domains:** orthologs in “model organisms” of the three domains of life are shown in this tab. This view, which provides orthologs across all domains of life, is relevant for broader comparative genomics studies. This tab is only available for proteins belonging to “model organisms”.

1
2
3 All ortholog relations are shown in a table giving basic information: the type of relations (one-to-one, one-
4 to-many, many-to-one, many-to-many), identifiers of all inparalogs (for many-to-*) and orthologs with links
5 to their respective protein pages on the OrthoInspector and Uniprot web sites, the species name (linking to
6 the NCBI taxonomy) and a summary of the species taxonomy. Additional information about orthologs
7 (protein description and length) can be shown by customizing the output using the columns output button,
8 in the top right corner.
9

10 By default, the table is ordered according to the taxonomic distance of the target species from the query
11 species, as inferred from the NCBI taxonomy. Thus, except in the case of unusual evolutionary events, the
12 first orthologs shown will be more closely related to the query protein. In the case of proteins with a large
13 number of orthologs, a search bar allows the user to search specific results by identifier, species name,
14 species taxid or even a specific clade name. For example, if a user is interested in orthologs of a human
15 protein in representatives of the carnivore clade only, typing 'carnivora' on the search bar will achieve this.
16

17 *Data export*

18 From the protein page, multiple export options are available. Exports of the table itself are available in
19 numerous formats (Excel, CSV, XML...) via the top right corner 'Export' button. User can also retrieve all
20 sequences involved in selected relations (all inparalogs and orthologs) in FASTA format, which could
21 serve as a starting point for further analyses.
22 OrthoInspector also offers the possibility to directly generate a multiple sequence alignment of the query
23 protein and all its orthologs in selected species (and inparalogs, if any) using the latest version of the
24 alignment workflow PipeAlign 2.0 (<http://www.lbgi.fr/pipealign>) (Kress, in prep).
25

26 Finally, on each protein page, the selected orthologous relations can be downloaded in the standardized
27 OrthoXML format, as defined by the Quest for Orthologs consortium (21).
28

29 *Taxonomic distribution summary*

30
31 The orthologs table contains, as seen above, all information about orthology relations. However, making
32 sense of such tables can be a daunting task, especially for proteins involved in many orthology relations.
33 To facilitate knowledge extraction, the OrthoInspector protein page provides a summary view of the
34 ortholog distribution at three levels of granularity: the domain's model organisms, the whole domain and all
35 three domains.
36

37 This information appears in a banner above the orthologs table after complete loading and is displayed as
38 a heatmap (see Figure 2) on a single row. Each tile of the heatmap corresponds to a major clade (Figure
39 1) of the selected domain, defined either from the NCBI taxonomy (22) or in some cases from the
40 consensus in the literature (for example, 'Excavata' appears in the cross-domains banner and is widely
41 accepted by the community despite not existing as such in the NCBI taxonomy). For each clade, the
42 corresponding tile is colored in green if orthologs are found in all its representatives and red if no orthologs
43 are found, with intermediary states between these two colors if orthologs are found in a subset of
44 representatives. The number of species in which orthologs are found and the total number of species
45 belonging to the clade represented in the OrthoInspector database are both displayed when hovering over
46 the tiles.
47

48 The clades on the heatmap are ordered according to the taxonomy: clades close to each other are side by
49 side on the heatmap. The heatmap provides users with preliminary information about the evolutionary
50 history (emergence and losses in major clades) of their protein family at a glance.
51

52 The clades displayed in this view depend on the granularity level selected by the user. In the cross-
53 domain view, only high-level clades are indicated ('First level' in Figure 1, Figure 2a), for instance
54 Opisthokonta. The domain of each clade is clearly indicated in the banner, by an indicator above the
55 heatmap and by a color code. Some of the high-level clades are detailed in the 'domain's model
56 organisms' and 'whole domain' views. For instance, Opisthokonta appear as Fungi, Choanoflagellida,
57
58
59
60

1
2
3 Metazoa and Other Opisthokonta ('Second level' in Figure 1, Figure 2b). Additionally, major clades
4 referencing many species can be further divided by clicking on the tile to display subtaxa and show a more
5 nuanced version of the distribution (see 'Third level' in Figure 1, Figure 2c). For instance, 15 phyla or
6 subphyla can be visualized for the Metazoa kingdom (156 species including 47 'model' species). These
7 clickable tiles are identified by a blue frame.

8 9 *Inparalogs distribution*

10 Information about presence and absence of orthologs is fundamental when studying the evolutionary
11 histories of proteins, but can miss some evolutionary events, notably duplication events. To address this
12 issue, the taxonomic summary banner also provides a 'See inparalogs' button, that shows all inparalogs of
13 the query protein relative to the considered clade. They are represented by ticks under the heatmap tiles
14 that provide information about the timing of each duplication during the gene's evolutionary history (Figure
15 2a). For example, an inparalog of a human protein found in relation to all species except Opisthokonts
16 may indicate a duplication of the ancestral gene in the Opisthokonta common ancestor.

17
18 Finally, the summary section also includes the list of species in which no orthologs were found.

19 20 **Phylogenetic profiling tools**

21 The presence and absence of orthologs summarized in the above section can be represented as detailed
22 binary profiles, the phylogenetic profiles.

23 24 *Searching for proteins with similar evolutionary histories*

25
26 The OrtholInspector protein page can be used to find other proteins of the same species with similar
27 phylogenetic profiles. This information is available under the 'Proteins with similar distribution' section on
28 the Protein page. The data available in these sections are based on the Jaccard distance between all
29 phylogenetic profiles of proteins in the same species (see supplementary materials and methods). The
30 identifiers of proteins exhibiting a phylogenetic profile distance below 0.4 are shown, along with a short
31 description of their functions and the exact value of the distance. For clarity, only the five closest proteins
32 are shown; additional proteins can be visualized by clicking 'See more'.

33
34 Distances are available both from a domain centric point of view (calculated on profiles limited to species
35 of the same domain) or from a cross-domain point of view. While the domain specific section is available
36 for all species in OrtholInspector, the cross-domain section is only available for 'model species'. Distances
37 between intra-domain and cross-domain profiles may differ significantly only for proteins that are present
38 in multiple domains.

39
40 Ciliary proteins are a good example of proteins whose phylogenetic profiles are clearly correlated to their
41 function, since the cilium has a very specific evolutionary history in Eukaryotes including multiple
42 independent losses (8). The cilium critically depends on molecular complexes to function properly, notably
43 the intraflagellar transport (IFT) complexes (23). We searched a core protein of the IFT-A complex, IFT122
44 (IF122_HUMAN) on the OrtholInspector website. In the 'Proteins with similar distribution in Eukaryota'
45 section, we found a list of 33 proteins, showing a significant enrichment in the GO term 'cilium' (P-
46 value: 4.93×10^{-43}). This list includes 4 out of the 5 other components of the IFT-A complex and 8 out of the
47 16 components of IFT-B, most of them with a distance <0.3 .

48
49 As illustrated by this example, these sections provide an original perspective when studying the function of
50 proteins and can be used to obtain a list of other proteins with potentially similar functions and possible
51 interaction partners.

52 53 *Searching proteins with a known profile*

54 Genes associated with a given phylogenetic trait tend to share the same distribution. The distribution of a
55 trait can thus be exploited to identify associated genes. OrtholInspector offers an original tool for
56
57
58
59
60

1
2
3 phylogenetic profiling, i.e. to search for proteins with orthologs present in a defined set of species or
4 clades and absent in others. This tool is available from the home page and under the 'Access/Search by
5 profile' tab. Users should select their query species on the dropdown menu and then interact with a
6 dynamic representation of the NCBI taxonomic tree to define the profile. Clicking once on a clade imposes
7 the presence of orthologs in at least one species of the clade, double clicking imposes the absence in all
8 species, a third click removes the constraints. Once the constraints are set, the database is queried to find
9 all proteins meeting the user's requirements (Figure 3).

10
11 The resulting proteins are displayed as panels in the result windows with their distribution summary (see
12 above) to facilitate identification of distribution subcategories within the results. Each protein panel also
13 contains a short description of the protein along with the associated Gene Ontology terms. For a functional
14 analysis of the complete protein list, a button can be clicked to run a GO term enrichment analysis using
15 the Panther webservice (24). The full list of proteins obtained can be exported using the 'Download list'
16 button, for further analysis.

17
18 Using this tool, we performed a phylogenetic profile search on the cross-domain database. Our objective
19 was to identify Eukaryotic Signature Proteins (ESP) that were also present in Asgard Archaea, a clade
20 whose discovery sparked interest due to its unexpected similarity to Eukaryotes (15, 16). We searched for
21 orthologs of *Homo sapiens* proteins present in Archaea of the Asgard group but absent in other Archaea
22 and in Bacteria (Figure 3a). This operation resulted in a total of 69 proteins with the required distributions
23 (Figure 3b). The list shows a strong enrichment in proteins with GTPase activity (P-value: 4.97×10^{-28}) and
24 vesicle-mediated transport (P-value: 5.12×10^{-36}), in agreement with previous studies (16). We also
25 retrieved actin-cytoskeleton proteins and ubiquitin-associated proteins, two iconic examples of ESP
26 previously reported in the Asgard group (16). As shown here, the phylogenetic profile search rapidly
27 provides both a list of genes associated with specific distributions and the tools required to extract
28 functional knowledge.

29 *Identifying profiles linked to a functional category*

30
31 OrthoInspector provides an original tool to explore the evolutionary history of a biological function, process
32 or component. This module, available on the home page or via the 'Access/GO profile' tab, provides the
33 distribution of all proteins of a species associated with a given GO term. After selecting the database,
34 species, and GO term of interest, the user retrieves the list of matching proteins, in the format described
35 above with the summary of the distribution of each associated protein. In this way, users can derive the
36 distribution associated with their function of interest and explore the different evolutionary histories of
37 proteins involved in the same biological system.

38 **Data and software accessibility**

39
40 This database update is complemented by the release of the new version 3.0 of the OrthoInspector
41 software suite, developed in Java, and available for download on the website in the download section
42 (http://www.lbgi.fr/orthoinspectorv3/download_Package). This release does not involve changes to the
43 main algorithm (10) but provides several software improvements.

44 *Software improvement*

45
46 Several code modifications were performed to optimize the management of the massive quantities of data.
47 This implies type changes to handle larger datasets, code optimization by reducing loop redundancy, the
48 use of more optimized data structures (library Fastutil, (25)) and more efficient database access from the
49 software (fewer SQL queries). This version of OrthoInspector runs faster than the precedent for large
50 computations and can still be parallelized when installing a large database.

51 *Improved accessibility*

52
53 Following feedback from users, the new OrthoInspector version provides an easier accessibility for small
54 datasets. Until now, fully supported database systems included MySQL and PostgreSQL, which require
55
56
57
58
59
60

1
2
3 prior experience of SQL management systems. This version comes with full support for SQLite database,
4 which eliminates most of the preliminary steps for computing a local database since no database server
5 configuration is required. We recommend the use of the easily accessible SQLite database option when
6 installing small local databases and, for performance reasons, the use of PostgreSQL and MySQL
7 systems for larger databases (several hundred of species). Updated tutorials for the installation procedure
8 are available on the website.
9

10 *Precomputed databases*

11 All four precomputed databases (Eukaryotes, Bacteria, Archaea, Cross-Domain) can be accessed *via* the
12 website interface. Due to the data volume (up to multiple terabytes in a single database), the database
13 dump is not available for direct download but could be made available on demand.
14

15 *Quest for Ortholog consortium reference proteome*

16
17 The Quest for Ortholog (QFO) consortium is part of an ongoing effort from the community pushing for
18 standardization in orthology inferences. The QFO consortium published a list of 78 reference proteomes
19 representing high quality proteomes and recommend using it for benchmarking purposes. The
20 precomputed orthology relationship made using this benchmark are available on
21 <http://www.lbgi.fr/orthoinspectorv3/QFO>.
22

23 *Webservices*

24 In addition to the web interface, a programming access is a major requirement for modern databases, as it
25 allows more flexible use of data. In this release, we introduce a Representational State Transfer (REST)
26 API providing access to most data available from the website, through the Swagger framework
27 (<https://swagger.io>). The documentation is available on the website
28 (<http://www.lbgi.fr/orthoinspectorv3/API>) where all endpoints and their parameters are described. All
29 queries can be executed with custom parameters directly from the documentation page.
30

31 **CONCLUSIONS AND FUTURE DIRECTIONS**

32
33 With this new release of OrtholInspector, we provide improvement in two main areas: proteome coverage
34 and information design.
35

36 The new databases boast a massive increase in the number of species across the three domains of life
37 and provide the most comprehensive ortholog relations resource in terms of species coverage.
38 Nevertheless, this increase did not involve simply adding a substantial number of species. Special
39 attention was paid to both quality of proteomes and taxonomic coverage. With the increasing rate of
40 genome sequencing, our scheduled strategy to ensure scalability will include regular updates of the
41 current proteome content and the addition of new species while maintaining our standard of proteome
42 quality. This will come with an updating procedure directly added to the software suite to allow any user to
43 easily update their local databases with the latest data.
44

45 In terms of accessibility, the installation process of local databases using the software suite has been
46 simplified and more importantly, the web interface of the OrtholInspector precomputed databases has
47 been significantly reorganized. The new design offers improved access to orthologous data in the three
48 domains of life. In addition, we believe that the implementation of original and user-friendly comparative
49 genomics tools will be useful for anyone interested in comparative genomics and evolutionary studies of
50 protein families. The next step for OrtholInspector will be the automated definition and analysis of
51 orthologous families among 'model species' by exploiting our experience in multiple sequence alignment
52 construction (26, 27) (Kress, in prep). This will allow the exploration of protein evolution through the three
53 life domains at different levels of resolution, from presence/absence of orthologs to subtler patterns of
54 differential conservation at domain or block levels.
55

56 **SUPPLEMENTARY DATA**

57
58
59
60

Additional information about proteome selection, database computations and retrieval of data displayed on the website are available in Supplementary material and methods.

Supplementary table 1. Species available in OrthoInspector. List of all species available in the OrthoInspector databases, along with their domain and their classification under our 3 levels of taxonomic granularity. Model species are indicated.

Supplementary table 2. Species for which a Uniprot Reference proteome was available and filters applied to include and exclude them in the final release.

ACKNOWLEDGEMENTS

We thank the Bio-statistics, Informatics and Complex System platform (BICS) and BISTRO bioinformatics platforms for informatics support and the European Grid Infrastructure for cloud computing facilities. We also thank our users for their feedback that helped to improve our suite and website.

FUNDING

This work was supported by the Agence Nationale de la Recherche [BIPBIP: ANR-10-BINF-03-02, ReNaBi-IFB: ANR-11-INBS-0013, Labex Agro: ANR-10-LABX-0001-01 to B.L., Labex CeMEB: ANR-10-LABX-0004 to B.L., Labex NUMEV: ANR-10-LABX-20 to B.L.]; and Institute funds from the Centre National de la Recherche Scientifique and the Université de Strasbourg. Funding for open access charge: Centre National de la Recherche Scientifique.

REFERENCES

1. Fitch,W.M. (1970) Distinguishing homologous from analogous proteins. *Syst. Zool.*, **19**, 99–113.
2. Tatusov,R.L., Koonin,E.V. and Lipman,D.J. (1997) A Genomic Perspective on Protein Families. *Science*, **278**, 631–637.
3. Nehrt,N.L., Clark,W.T., Radivojac,P. and Hahn,M.W. (2011) Testing the Ortholog Conjecture with Comparative Functional Genomic Data from Mammals. *PLOS Comput. Biol.*, **7**, e1002073.
4. Altenhoff,A.M., Studer,R.A., Robinson-Rechavi,M. and Dessimoz,C. (2012) Resolving the ortholog conjecture: orthologs tend to be weakly, but significantly, more similar in function than paralogs. *PLoS Comput. Biol.*, **8**, e1002514.
5. Koonin,E.V. (2005) Orthologs, paralogs, and evolutionary genomics. *Annu. Rev. Genet.*, **39**, 309–338.
6. Pellegrini,M., Marcotte,E.M., Thompson,M.J., Eisenberg,D. and Yeates,T.O. (1999) Assigning protein functions by comparative genome analysis: protein phylogenetic profiles. *Proc. Natl. Acad. Sci. U. S. A.*, **96**, 4285–4288.
7. Pellegrini,M. (2012) Using phylogenetic profiles to predict functional relationships. *Methods Mol. Biol. Clifton NJ*, **804**, 167–177.
8. Nevers,Y., Prasad,M.K., Poidevin,L., Chennen,K., Allot,A., Kress,A., Ripp,R., Thompson,J.D., Dollfus,H., Poch,O., *et al.* (2017) Insights into Ciliary Genes and Evolution from Multi-Level Phylogenetic Profiling. *Mol. Biol. Evol.*, **34**, 2016–2034.

- 1
2
3 9. Forslund,K., Pereira,C., Capella-Gutierrez,S., Silva,D., Sousa,A., Altenhoff,A., Huerta-Cepas,J.,
4 Muffato,M., Patricio,M., Vandepoele,K., *et al.* (2018) Gearing up to handle the mosaic nature of
5 life in the quest for orthologs. *Bioinformatics*, **34**, 323–329.
6
- 7 10. Linard,B., Thompson,J.D., Poch,O. and Lecompte,O. (2011) OrthoInspector: comprehensive orthology
8 analysis and visual exploration. *BMC Bioinformatics*, **12**, 11.
9
- 10 11. Linard,B., Allot,A., Schneider,R., Morel,C., Ripp,R., Bigler,M., Thompson,J.D., Poch,O. and
11 Lecompte,O. (2015) OrthoInspector 2.0: Software and database updates. *Bioinforma. Oxf. Engl.*,
12 **31**, 447–448.
13
- 14 12. Altenhoff,A.M., Boeckmann,B., Capella-Gutierrez,S., Dalquen,D.A., DeLuca,T., Forslund,K., Huerta-
15 Cepas,J., Linard,B., Pereira,C., Prysycz,L.P., *et al.* (2016) Standardized benchmarking in the quest
16 for orthologs. *Nat. Methods*, **13**, 425–430.
17
- 18 13. Liebeskind,B.J., McWhite,C.D. and Marcotte,E.M. (2016) Towards Consensus Gene Ages. *Genome*
19 *Biol. Evol.*, **8**, 1812–1823.
20
- 21 14. UniProt: the universal protein knowledgebase (2017) *Nucleic Acids Res.*, **45**, D158–D169.
22
- 23 15. Spang,A., Saw,J.H., Jørgensen,S.L., Zaremba-Niedzwiedzka,K., Martijn,J., Lind,A.E., van Eijk,R.,
24 Schleper,C., Guy,L. and Ettema,T.J.G. (2015) Complex archaea that bridge the gap between
25 prokaryotes and eukaryotes. *Nature*, **521**, 173–179.
26
- 27 16. Zaremba-Niedzwiedzka,K., Caceres,E.F., Saw,J.H., Bäckström,D., Juzokaite,L., Vancaester,E.,
28 Seitz,K.W., Anantharaman,K., Starnawski,P., Kjeldsen,K.U., *et al.* (2017) Asgard archaea
29 illuminate the origin of eukaryotic cellular complexity. *Nature*, **541**, 353–358.
30
- 31 17. Altschul,S.F., Gish,W., Miller,W., Myers,E.W. and Lipman,D.J. (1990) Basic local alignment search tool.
32 *J. Mol. Biol.*, **215**, 403–410.
33
- 34 18. Camacho,C., Coulouris,G., Avagyan,V., Ma,N., Papadopoulos,J., Bealer,K. and Madden,T.L. (2009)
35 BLAST+: architecture and applications. *BMC Bioinformatics*, **10**, 421.
36
- 37 19. Gene Ontology Consortium (2015) Gene Ontology Consortium: going forward. *Nucleic Acids Res.*, **43**,
38 D1049-1056.
39
- 40 20. Finn,R.D., Attwood,T.K., Babbitt,P.C., Bateman,A., Bork,P., Bridge,A.J., Chang,H.-Y., Dosztányi,Z., El-
41 Gebali,S., Fraser,M., *et al.* (2017) InterPro in 2017—beyond protein family and domain
42 annotations. *Nucleic Acids Res.*, **45**, D190–D199.
43
- 44 21. Dessimoz,C., Gabaldón,T., Roos,D.S., Sonnhammer,E.L.L., Herrero,J., Altenhoff,A., Apweiler,R.,
45 Ashburner,M., Blake,J., Boeckmann,B., *et al.* (2012) Toward community standards in the quest
46 for orthologs. *Bioinformatics*, **28**, 900–904.
47
- 48 22. Sayers,E.W., Barrett,T., Benson,D.A., Bryant,S.H., Canese,K., Chetvernin,V., Church,D.M., DiCuccio,M.,
49 Edgar,R., Federhen,S., *et al.* (2009) Database resources of the National Center for Biotechnology
50 Information. *Nucleic Acids Res.*, **37**, D5-15.
51
52
53
54
55
56
57
58
59
60

- 1
2
3 23. Lechtreck,K.F. (2015) IFT-Cargo Interactions and Protein Transport in Cilia. *Trends Biochem. Sci.*, **40**,
4 765–778.
5
6 24. Mi,H., Poudel,S., Muruganujan,A., Casagrande,J.T. and Thomas,P.D. (2016) PANTHER version 10:
7 expanded protein families and functions, and analysis tools. *Nucleic Acids Res.*, **44**, D336-342.
8
9 25. Boldi,P., Marino,A., Santini,M. and Vigna,S. (2016) BUBiNG: Massive Crawling for the Masses.
10 *ArXiv160106919 Cs*.
11
12 26. Vanhoutreuve,R., Kress,A., Legrand,B., Gass,H., Poch,O. and Thompson,J.D. (2016) LEON-BIS: multiple
13 alignment evaluation of sequence neighbours using a Bayesian inference system. *BMC*
14 *Bioinformatics*, **17**, 271.
15
16 27. Kress,A., Lecompte,O., Poch,O. and Thompson,J.D. (2018) PROBE: analysis and visualization of
17 protein block-level evolution. *Bioinforma. Oxf. Engl.*, 10.1093/bioinformatics/bty367.
18
19
20
21

22 TABLES AND FIGURE LEGENDS

23
24 Figure 1. Taxonomic distribution of species represented in OrthoInspector. The domain trees are
25 distributed on three 'levels'. The first level corresponds to the cross-domain taxonomic distribution
26 heatmap shown when browsing the cross-domain database, the second level is shown on the heatmap for
27 domain specific databases and the third level is the 'focus view' available for certain clades (See figure 2).
28 The size of a node is proportional to the number of species in the corresponding clades according to
29 indicated scales. The number of species and model species in first-level clades are displayed in black and
30 pink respectively.
31

32 Figure 2. Taxonomic distribution heatmaps. Each labelled tile corresponds to a clade and is colored
33 according to the proportion of species in the clade with at least one ortholog. Colors range from red (no
34 species) to green (all species). a. Heatmap corresponding to the cross-domain database. The domain of
35 life of the clades is shown by an additional label and a color code. Inparalogs distribution is indicated by a
36 tick under each clade. b. Heatmap corresponding to the eukaryotic database. The box framed by a thin
37 blue outline can be expanded to 'focus view'. c. Heatmap corresponding to the 'focus view' of Metazoa.
38

39 Figure 3. Phylogenetic profile search interface. a. Definition of the phylogenetic profile. User selects: (1)
40 the database, (2) the query species in the drop-down menu and (3) the presence/absence constraints
41 using the phylogenetic tree. A summary of constraints is shown below the tree. Here, human proteins
42 absent in Prokaryotes except the archaeal Asgard group are selected. b. Output of the profile search.
43 Constraints are included on the top with the number of proteins found. Proteins are displayed in panels,
44 showing their distributions and functional information. Gene Ontology enrichment can be performed on the
45 protein list.
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

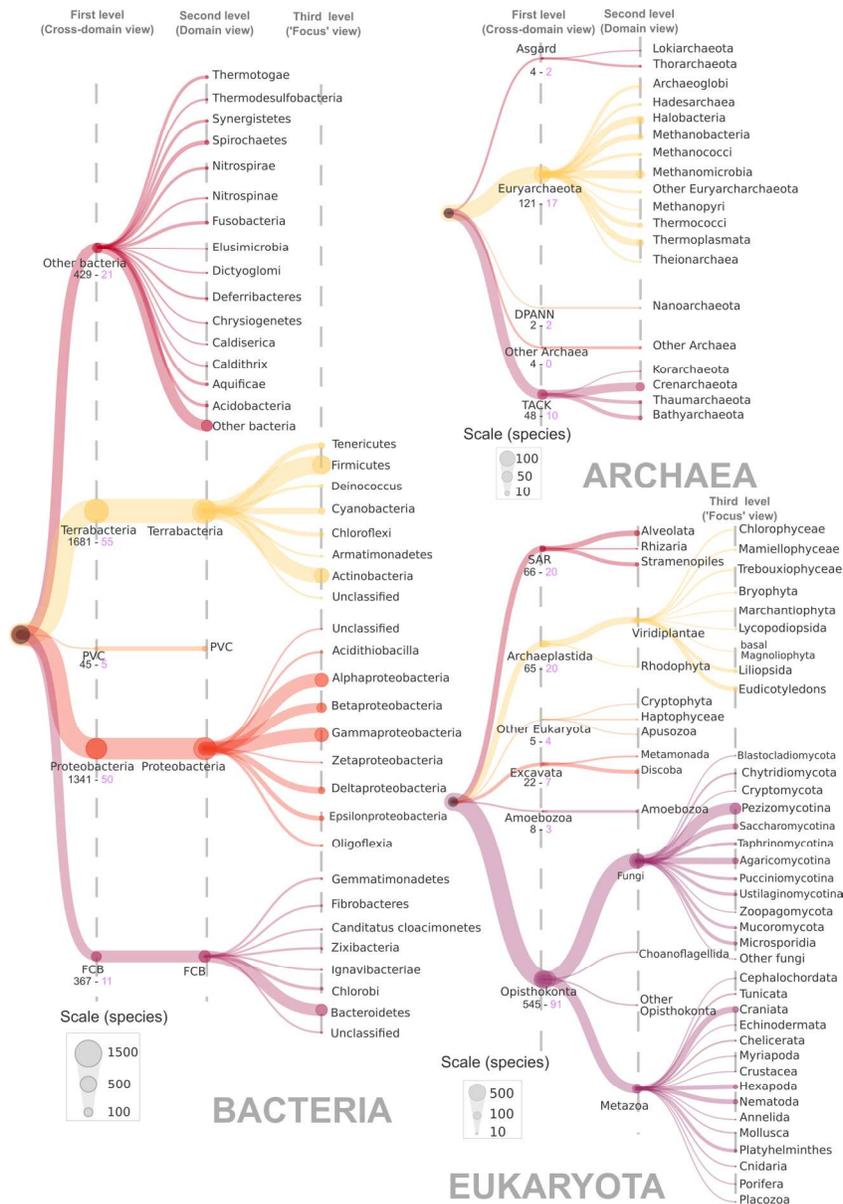


Figure 1. Taxonomic distribution of species represented in OrthoInspector. The domain trees are distributed on three 'levels'. The first level corresponds to the cross-domain taxonomic distribution heatmap shown when browsing the cross-domain database, the second level is shown on the heatmap for domain specific databases and the third level is the 'focus view' available for certain clades (See figure 2). The size of a node is proportional to the number of species in the corresponding clades according to indicated scales. The number of species and model species in first-level clades are displayed in black and pink respectively.



Figure 2. Taxonomic distribution heatmaps. Each labelled tile corresponds to a clade and is colored according to the proportion of species in the clade with at least one ortholog. Colors range from red (no species) to green (all species). a. Heatmap corresponding to the cross-domain database. The domain of life of the clades is shown by an additional label and a color code. Inparalogs distribution is indicated by a tick under each clade. b. Heatmap corresponding to the eukaryotic database. The box framed by a thin blue outline can be expanded to 'focus view'. c. Heatmap corresponding to the 'focus view' of Metazoa.

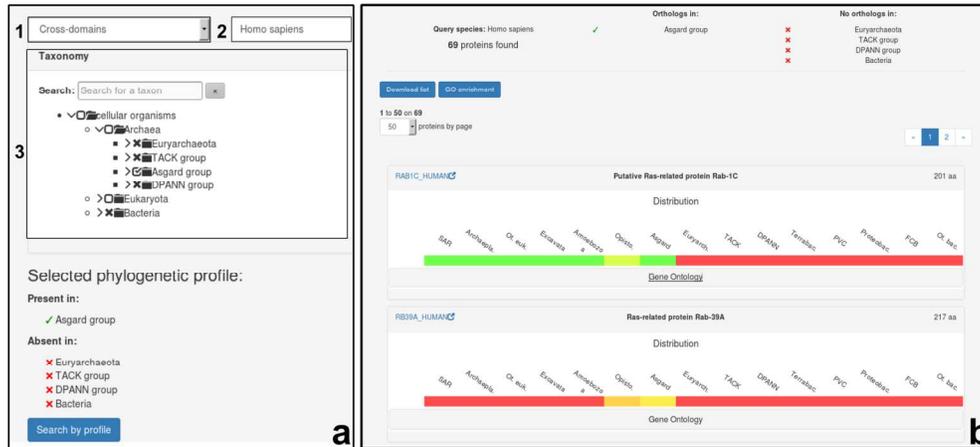


Figure 3. Phylogenetic profile search interface. a. Definition of the phylogenetic profile. User selects: (1) the database, (2) the query species in the drop-down menu and (3) the presence/absence constraints using the phylogenetic tree. A summary of constraints is shown below the tree. Here, human proteins absent in Prokaryotes except the archaeal Asgard group are selected. b. Output of the profile search. Constraints are included on the top with the number of proteins found. Proteins are displayed in panels, showing their distributions and functional information. Gene Ontology enrichment can be performed on the protein list.

Supplementary Materials

MATERIAL AND METHODS

Proteome selection

The goal in constructing the new OrthoInspector databases was to cover species representatives of most of the known tree of life, while avoiding incomplete proteomes and fragmentary or misannotated proteins. We generated four main databases: one dedicated to each of the three life domains (Eukaryotes, Bacteria, Archaea) and a cross-domain database. To construct them, we downloaded all Uniprot Reference Proteomes (1) that are specifically designed to adequately cover the tree of life (downloaded in November 2016). However, incomplete or fragmentary annotation of genomes (for example, due to lack of coverage) can be a source of orthology misprediction. Thus, all proteomes were submitted to a preliminary quality control based on three criteria: protein number, distribution of protein lengths and the nature of the first amino acid of the protein sequences.

First, we excluded proteomes with outlier numbers of proteins that are not explained by biological peculiarities (e.g. symbiosis or parasitism for small proteomes) and that did not correspond to the literature. Second, we examined the distribution of protein lengths in each proteome and the proportion of small proteins (< 100 aa) as well as the proportion of protein sequences starting with an amino acid other than methionine (hereafter referred to as 'false-start'). The last two parameters were used as a proxy to estimate the number of fragmentary proteins in proteomes and used to filter low quality ones. For Archaea and Bacteria, we excluded proteomes with more than 20% of small proteins and/or 10% of false-start proteins and/or more than 10% proteins annotated as fragments. For Eukaryotes, we kept the same threshold for small proteins and excluded proteomes with more than 55% of false start proteins. This more lenient threshold takes into account the higher proportion of false start proteins generally observed in Eukaryotes due to the complex nature of eukaryotic genomes, which are classically more challenging to assemble and annotate.

Orthology inference

For each precomputed database, we constructed a BLAST database with all the corresponding selected reference proteomes and performed all-vs-all BLASTp searches using BLAST+ (2, 3) with the following parameters: wordsize=3, output size=5000, expect threshold= 10^{-5} and compositional adjustment. BLAST outputs were processed to construct the precomputed databases using the latest OrthoInspector Suite release (available at <http://www.lbgi.fr/orthoinspectorv3/download>) with expect threshold set to 10^{-5} .

Phylogenetic profiling and distance computing

Phylogenetic (presence-absence) profiles were generated using an in-house python script, directly connected to the databases. For each protein of a query species, all relations (one-to-one, one-to-many, many-to-one, many-to-many) were retrieved from the databases. A presence was assigned for a given target species if an orthologous relationship exists between the query protein and one or more proteins of the target species. The output is a binary matrix saved in a tab separated file with row corresponding to proteins and columns corresponding to species. We ran this script for each species in OrthoInspector. Profiles were generated for all proteins in all species for the small Archaea database and in 'model species' (see below for a definition) for Bacteria and Eukaryotes.

The Jaccard distance, a distance measure suitable for binary vectors, was calculated for each presence/absence matrix using the `dist.binary` function in `ade4` packages (4) and saved in another tab separated file. Finally, we saved each pair of proteins with a distance less than a threshold of 0.4 in a dedicated PostgreSQL database. These distances can be accessed on the website from the corresponding protein pages.

Protein domains and Gene Ontology

1
2
3 Information about protein domains is retrieved from InterPro (5) via the EBI webservice. Functional
4 enrichment analyses available on the website are performed on the three Gene Ontology (GO) categories
5 (parameters: 'function', 'process' 'cellular_localisation') using the Panther (6) enrichment webservice. GO
6 annotations available on the websites are extracted from our local instance of the official GO database (7).
7

8 REFERENCES

- 9 1. UniProt: the universal protein knowledgebase (2017) *Nucleic Acids Res*, **45**, D158–D169.
 - 10 2. Altschul,S.F., Madden,T.L., Schäffer,A.A., Zhang,J., Zhang,Z., Miller,W. and Lipman,D.J. (1997)
11 Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids*
12 *Res.*, **25**, 3389–3402.
 - 13 3. Camacho,C., Coulouris,G., Avagyan,V., Ma,N., Papadopoulos,J., Bealer,K. and Madden,T.L. (2009)
14 BLAST+: architecture and applications. *BMC Bioinformatics*, **10**, 421.
 - 15 4. Dray,S. and Dufour,A.-B. (2007) The ade4 Package: Implementing the Duality Diagram for Ecologists.
16 *Journal of Statistical Software, Articles*, **22**, 1–20.
 - 17 5. Finn,R.D., Attwood,T.K., Babbitt,P.C., Bateman,A., Bork,P., Bridge,A.J., Chang,H.-Y., Dosztányi,Z., El-
18 Gebali,S., Fraser,M., *et al.* (2017) InterPro in 2017—beyond protein family and domain annotations.
19 *Nucleic Acids Res*, **45**, D190–D199.
 - 20 6. Mi,H., Poudel,S., Muruganujan,A., Casagrande,J.T. and Thomas,P.D. (2016) PANTHER version 10:
21 expanded protein families and functions, and analysis tools. *Nucleic Acids Res.*, **44**, D336-342.
 - 22 7. Gene Ontology Consortium (2015) Gene Ontology Consortium: going forward. *Nucleic Acids Res.*, **43**,
23 D1049-1056.
- 24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60