

The Gene Ontology Knowledgebase in 2023

The Gene Ontology Consortium*

Corresponding author: Paul D Thomas, Department of Population and Public Health Sciences, University of Southern California, Los Angeles, CA USA, pdthomas@usc.edu

Authors

GO Central: Suzi A. Aleksander (ORCID:0000-0001-6787-2901), James Balhoff (ORCID:0000-0002-8688-6599), Seth Carbon (ORCID:0000-0001-8244-1536), J.Michael Cherry (ORCID:0000-0001-9163-5180), Harold J Drabkin (ORCID 0000-0003-2689-5511), Dustin Ebert (ORCID 0000-0002-6659-0416), Marc Feuermann (ORCID: 0000-0002-4187-2863), Pascale Gaudet (ORCID:0000-0003-1813-6857), Nomi L Harris (ORCID:0000-0001-6315-3707), David P Hill (ORCID 0000-0001-7476-6306), Raymond Lee (ORCID 0000-0002-8151-7479), Huaiyu Mi (ORCID:0000-0001-8721-202X), Sierra Moxon (ORCID: 0000-0002-8719-7760), Christopher J Mungall (ORCID 0000-0002-6601-2165), Anushya Muruganugan (ORCID 0000-0001-7169-5864), Tremayne Mushayahama (ORCID 0000-0002-2874-6934), Paul W. Sternberg (ORCID 0000-0002-7699-0173), Paul D Thomas (ORCID 0000-0002-9074-3507), Kimberly Van Auken (ORCID 0000-0002-8151-747).

CACAO, EcoliWiki: Jolene Ramsey (ORCID: 0000-0002-3774-5896), Deborah A. Siegele (ORCID: 0000-0001-8935-0696)

dictyBase: Rex L. Chisholm (ORCID 0000-0002-5638-3990), Petra Fey (ORCID 0000-0002-4532-2703)

DisProt, Dept. of Biomedical Sciences, University of Padova, Italy: Maria Cristina Aspromonte (ORCID: 0000-0002-4937-6952), Maria Victoria Nugnes (ORCID: 0000-0001-8399-7907), Federica Quaglia (ORCID: 0000-0002-0341-4888), Silvio Tosatto (ORCID: 0000-0003-4525-7793)

ECO: Michelle Giglio (ORCID: 0000-0001-7628-5565), Suvarna Nadendla (ORCID: 0000-0003-3643-281X)

FlyBase: Giulia Antonazzo (ORCID 0000-0003-0086-5621), Helen Attrill (ORCID 0000-0003-3212-6364), Gil dos Santos (ORCID 0000-0003-3507-8273), Steven Marygold (ORCID 0000-0003-2759-266X), Victor Strelets (ORCID 0000-0001-6556-9335), Christopher J. Tabone (ORCID 0000-0001-8746-0680), Jim Thurmond (ORCID 0000-0002-5142-2583), Pinglei Zhou

© The Author(s) 2023. Published by Oxford University Press on behalf of the Genetics Society of America. This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

1 (ORCID 0000-0002-3012-1044).

2 **Functional Gene Annotation, Institute of Cardiovascular Science, University College**
3 **London (London, UK):** Saadullah H. Ahmed, Praoparn Asanitthong (ORCID:0000-0002-6286-
4 7959), Diana Luna Buitrago (0000-0003-2010-7174), Meltem N. Erdol (ORCID: 0000-0002-
5 6574-5519), Matthew C. Gage (ORCID: 0000-0002-6668-6573), Mohamed Ali Kadhum, Kan
6 Yan Chloe Li (ORCID:0000-0001-7737-2118), Miao Long, Aleksandra Michalak, Angeline
7 Pesala (ORCID: 0000-0002-4639-2295), Armalya Pritazahra, Shirin C.C. Saverimuttu
8 (ORCID:0000-0003-1191-2681), Renzhi Su (0000-0002-6819-5356), Kate E. Thurlow (ORCID:
9 0000-0001-9985-3684), Ruth C. Lovering (ORCID: 0000-0002-9791-0064)

10 **GREEKC:** Colin Logie (ORCID: 0000-0002-8534-6582)

11 **JaponicusDB (London, UK):** Snezhana Oliferenko (ORCID 0000-0002-8138-6851)

12
13 **MGI:** Judith Blake (ORCID 0000-0001-8522-334X), Karen Christie (ORCID 0000-0001-5501-
14 853X), Lori Corbani (ORCID 0000-0002-2366-557X), Mary E Dolan (ORCID 0000-0001-7732-
15 3295), Harold J Drabkin (ORCID 0000-0003-2689-5511), David P Hill (ORCID 0000-0001-7476-
16 6306), Li Ni (ORCID 0000-0002-9796-7693), Dmitry Sitnikov (ORCID 0000-0003-3394-9805),
17 Cynthia Smith (ORCID 0000-0003-3691-0324)

18 **PHI-base:** Alayne Cuzick (ORCID 0000-0001-8941-3984), James Seager (ORCID 0000-0001-
19 7487-610X)

20
21 **Planteome:** Laurel Cooper (ORCID 0000-0002-6379-8932, Justin Elser (ORCID 0000-0003-
22 0921-1982), Pankaj Jaiswal (ORCID 0000-0002-1005-8383),

23 **Plant Reactome:** Parul Gupta (ORCID 0000-0002-0190-8753), Pankaj Jaiswal (ORCID 0000-
24 0002-1005-8383), Sushma Naithani (ORCID 0000-0001-7819-4552),

25 **PomBase:** Manuel Lera-Ramirez (ORCID 0000-0002-8666-9746), Kim Rutherford (ORCID
26 0000-0001-6277-726X), Valerie Wood (ORCID 0000-0001-6330-7526)

27 **RGD:** Jeffrey L. De Pons, Melinda R. Dwinell (ORCID [0000-0002-9528-3618](#)), G. Thomas
28 Hayman (ORCID 0000-0002-9553-7227), Mary L. Kaldunski (ORCID 0000-0003-3645-6803),
29 Anne E. Kwitek (ORCID 0000-0003-1024-4116), Stanley J. F. Laulederkind (ORCID 0000-0001-
30 5356-4174), Marek A. Tutaj (ORCID 0000-0002-1025-101X), Mahima Vedi (ORCID 0000-0001-
31 5361-6739), Shur-Jen Wang (ORCID 0000-0001-5256-8683)

32 **Reactome:** Peter D'Eustachio (ORCID 0000-0002-5494-626X)

33 **Rhea, Swiss-Prot group, SIB Swiss Institute of Bioinformatics (SIB) (Geneva,**
34 **Switzerland):** Lucila Aimo (ORCID: 0000-0003-0943-6401), Kristian Axelsen (ORCID: 0000-
35 0003-3889-2879), Alan Bridge (ORCID: 0000-0003-2148-9135, Nevila Hyka-Nouspikel (ORCID:
36 0000-0001-7855-209X), Anne Morgat (ORCID: 0000-0002-1216-2969)

1 **SGD, Department of Genetics, Stanford University (Stanford, CA):** Suzi A. Aleksander
2 (ORCID:0000-0001-6787-2901), J.Michael Cherry (ORCID:0000-0001-9163-5180), Stacia R.
3 Engel (<https://orcid.org/0000-0001-5472-917X>), Kalpana Karra (ORCID:0000-0002-6816-2458),
4 Stuart R. Miyasato (ORCID:0000-0001-5250-8920), Robert S. Nash (ORCID:0000-0002-3726-
5 7441), Marek S. Skrzypek (ORCID:0000-0001-6749-615X), Shuai Weng (ORCID:0000-0003-
6 4233-0772), Edith D. Wong (ORCID:0000-0001-9799-5523)

7 **The Arabidopsis Information Resource (TAIR), Phoenix Bioinformatics (Newark, CA,**
8 **USA):** Erika Bakker (ORCID: 0000-0001-7933-3817), Tanya Z. Berardini (ORCID: 0000-0002-
9 3837-8864), Leonore Reiser (ORCID: 0000-0003-0073-0858)

10 **UniProt, Swiss-Prot group, SIB Swiss Institute of Bioinformatics (SIB) (Geneva,**
11 **Switzerland):** Andrea Auchincloss (ORCID: 0000-0002-5297-5390), Kristian Axelsen (ORCID:
12 0000-0003-3889-2879), Ghislaine Argoud-Puy (ORCID: 0000-0002-2979-8613), Marie-Claude
13 Blatter (ORCID: 0000-0002-7474-1499), Emmanuel Boutet (ORCID: 0000-0002-3743-4203),
14 Lionel Breuza (ORCID: 0000-0002-8075-8625), Alan Bridge (ORCID: 0000-0003-2148-9135),
15 Cristina Casals-Casas (ORCID: 0000-0001-8769-177X), Elisabeth Coudert (ORCID: 0000-
16 0001-8314-404X), Anne Estreicher (ORCID: 0000-0001-6828-2508), Maria Livia Famiglietti
17 (ORCID: 0000-0002-5283-6593), Marc Feuermann (ORCID: 0000-0002-4187-2863), Arnaud
18 Gos (ORCID: 0000-0002-5018-1378), Nadine Gruaz-Gumowski (ORCID: 0000-0002-4699-
19 4907), Chantal Hulo (ORCID: 0000-0001-8176-7999), Nevila Hyka-Nouspikel (ORCID: 0000-
20 0001-7855-209X), Florence Jungo (ORCID: 0000-0002-7456-8390), Philippe Le Mercier
21 (ORCID: 0000-0001-8528-090X), Damien Lieberherr (ORCID: 0000-0002-9724-1710), Patrick
22 Masson (ORCID: 0000-0001-7646-0052), Anne Morgat (ORCID: 0000-0002-1216-2969), Ivo
23 Pedruzzi (ORCID: 0000-0001-8561-7170), Lucille Pourcel (ORCID: 0000-0003-1522-9900),
24 Sylvain Poux (ORCID: 0000-0001-7299-6685), Catherine Rivoire (ORCID: 0000-0002-5979-
25 8382), Shyamala Sundaram (ORCID: 0000-0003-4209-460X).

26 **UniProt, EMBL-EBI, (Hinxton, UK):** Alex Bateman (ORCID: 0000-0002-6982-4660), Emily
27 Bowler-Barnett (ORCID: 0000-0003-4785-7231), Hema Bye-A-Jee (ORCID: 0000-0003-2464-
28 7688), Paul Denny (ORCID: 0000-0003-4659-6893), Alexandr Ignatchenko (ORCID: 0000-
29 0002-6083-941X), Rizwan Ishtiaq (ORCID: 0000-0001-8041-7321), Antonia Lock (ORCID:
30 0000-0003-1179-5999), Yvonne Lussi (ORCID: 0000-0002-5753-0235), Michele Magrane
31 (ORCID: 0000-0003-3544-996X), Maria J. Martin (ORCID: 0000-0001-5454-2815), Sandra
32 Orchard (ORCID: 0000-0002-8878-3972), Pedro Raposo (ORCID: 0000-0001-6149-9456),
33 Elena Speretta (ORCID: 0000-0003-1506-7438), Nidhi Tyagi (ORCID: 0000-0002-2065-9051),
34 Kate Warner (ORCID: 0000-0001-8705-181X), Rossana Zaru (ORCID: 0000-0002-3358-4423)

35 **University at Buffalo, Department of Biomedical Informatics (Buffalo, NY, USA):** Alexander
36 D. Diehl (ORCID: 0000-0001-9990-8331)

37
38 **WormBase:** Raymond Lee (ORCID 0000-0002-8151-7479), Juancarlos Chan (ORCID:0000-
39 0002-7259-8107), Stavros Diamantakis (ORCID:0000-0002-0273-3406), Daniela Raciti
40 (ORCID:0000-0002-4945-5837), Magdalena Zarowiecki (ORCID:0000-0001-6102-7731)

41

1 **Xenbase:** Malcolm Fisher (ORCID:0000-0003-1074-8103), Christina James-Zorn
2 (ORCID:0000-0001-5495-4588), Virgilio Ponferrada (ORCID:0000-0002-8590-7183), Aaron
3 Zorn (ORCID:0000-0003-3217-3590)

4
5 **ZFIN:** Sridhar Ramachandran (ORCID:0000-0002-2246-3722), Leyla Ruzicka (ORCID:0000-
6 0002-1009-339X), Monte Westerfield (ORCID:0000-0003-1187-7839)

7

8

9 1. **Keywords:** gene ontology, gene annotation, gene function, knowledgebase, knowledge
10 graphs

11

1 Abstract

2 The Gene Ontology (GO) knowledgebase (<http://geneontology.org>) is a comprehensive
3 resource concerning the functions of genes and gene products (proteins and non-coding RNAs).
4 GO annotations cover genes from organisms across the tree of life as well as viruses, though
5 most gene function knowledge currently derives from experiments carried out in a relatively
6 small number of model organisms. Here, we provide an updated overview of the GO
7 knowledgebase, as well as the efforts of the broad, international consortium of scientists that
8 develops, maintains and updates the GO knowledgebase. The GO knowledgebase consists of
9 three components: 1) the Gene Ontology – a computational knowledge structure describing
10 functional characteristics of genes; 2) GO annotations – evidence-supported statements
11 asserting that a specific gene product has a particular functional characteristic; and 3) GO
12 Causal Activity Models (GO-CAMs) – mechanistic models of molecular “pathways” (GO
13 biological processes) created by linking multiple GO annotations using defined relations. Each
14 of these components is continually expanded, revised and updated in response to newly
15 published discoveries, and receives extensive QA checks, reviews and user feedback. For each
16 of these components, we provide a description of the current contents, recent developments to
17 keep the knowledgebase up to date with new discoveries, as well as guidance on how users
18 can best make use of the data we provide. We conclude with future directions for the project.

19 Introduction

20 Genes encode gene products, often proteins but also non-coding RNA molecules (ncRNAs),
21 that perform functions at the molecular, cellular, and organismal levels. The GO knowledgebase
22 provides a comprehensive, structured, computer-accessible representation of gene function, for
23 genes from any cellular organism or virus. The GO knowledgebase has become a critical
24 component of life science research, supporting analysis of large-scale experiments and
25 biological systems (Duck *et al.* 2016). It is designed to make expert knowledge of gene function
26 accessible for bench scientists as well as computational analyses. The basic model underlying
27 GO is the “molecular biology paradigm” (Ashburner *et al.* 2000; Thomas 2017), in which there
28 are three types (“aspects”) of functional characteristics used to describe gene function:

- 29 ● Molecular function (MF): the activities performed by a gene product at the molecular
30 level
- 31 ● Cellular component (CC): the locations, relative to cellular structures, where molecular
32 functions are performed
- 33 ● Biological process (BP): a “biological program” comprising molecular activities acting in
34 concert to achieve a particular outcome; this program can be at the cellular level or at
35 the organism level of multicellular organisms.

36
37 The GO knowledgebase consists of three components: the Gene Ontology, GO annotations and
38 GO Causal Activity models (GO-CAMs) (Figure 1). The **Gene Ontology** (Figure 1A) structures
39 our current knowledge of the types of functional characteristics a gene product may possess
40 into a connected graph-based representation. Each ontology term (called “class” in the field of

1 ontologies) represents a functional characteristic that can be attributed to a gene product.
2 Terms can have relationships between them, such as one term being more specific than
3 another term (also called "subclass"); e.g. *DNA-binding transcription factor activity* is a subclass
4 of *transcription regulator activity*. A **GO annotation** (Figure 1B) is an association between a
5 specific gene (or gene product) and a GO term, and should be interpreted as a statement that
6 the specified gene product possesses the specified functional characteristic represented by the
7 GO term. Each GO annotation includes the evidence upon which it is based. Because each GO
8 annotation covers only a single characteristic of gene function, multiple GO annotations are
9 generally required to completely describe the function of a gene product. **GO-CAMs** (Figure 1C)
10 link multiple GO annotations together to create models of biological processes by 1) connecting
11 the activities of more than one gene product together into causal networks, and 2) allowing the
12 specification of the biological context (e.g. cell type, tissue type) in which the processes occur.
13
14

15 The GO knowledgebase is large and dynamic

16 For applications that use the components of the GO knowledgebase, it is crucial that the
17 ontology and associated annotations represent the current state of knowledge, and are not just
18 an archive of all public data. Therefore, all aspects of the GO knowledgebase are dynamic
19 (ontology, annotations, GO-CAMs, links to external ontologies, etc.), and citable, versioned
20 updates are released on a monthly basis. Below, we describe each component of the
21 knowledgebase, focusing on recent changes made to improve the resource during the past two
22 years. Statistics and descriptions given here are based on the GO release 2022-11-03
23 (<http://release.geneontology.org/2022-11-03>, doi:10.5281/zenodo.7407024).

24 Ontology

25 The ontology component of the GO knowledgebase consists of the terms used to describe
26 functional characteristics of gene products, which are linked together by relations into a labeled
27 directed acyclic graph (like a hierarchy but with multiple parentages allowed). It also includes
28 term definitions, synonyms, and relations to terms from external ontologies. The GO is available
29 in different editions, including (i) the "basic" edition, which includes only core relationship types;
30 (ii) the core ontology, including additional relationship types; and (iii) the "go-plus" edition which
31 includes relationships to terms in other ontologies. These editions are explained on the GO
32 downloads page <http://geneontology.org/docs/download-ontology/>. The ontology contains
33 43,303 terms (Table 2), linked together by 88,099 relationships in the basic edition. When
34 relationships to external terms are included, there are 121698 relationships; release statistics
35 can be viewed at: <http://geneontology.org/stats>.
36

37 The GO is subject to constant review and revision to most accurately model current biological
38 knowledge. Revision of the ontology includes the addition or obsolescence of terms and re-
39 organization of the relationship structure. New GO terms are added to represent concepts

1 previously missing from the GO in response to published findings or when a branch of GO is
2 revised. Terms may be obsoleted when unused or inconsistently used in annotation, when they
3 are redundant with other terms, or during revision of specific branches of the ontology.

4
5 Most of the revisions in the structure of GO are in response to advances in biological
6 knowledge, as well as improvements in the precision of newer experimental approaches. In
7 addition, because many branches of the ontology have grown organically in a bottom-up fashion
8 by accumulating specific individual term requests, we also perform systematic review aimed at
9 improving consistency and clarity while reducing redundancy. Additional revisions are initiated
10 by internal review and consistency and quality assurance checks. Revisions are also made
11 following feedback from users. Whenever possible, changes are performed in collaboration with
12 expert biocurators or domain specialists; recent examples include blood-brain barrier-related
13 functions (Saverimuttu *et al.* 2021) and transcription factors (Gaudet *et al.* 2021).

14
15 At each release, we track all changes, and report on our website the number of added,
16 obsoleted, and merged terms in the ontology. Table 1 shows the number of GO terms added
17 and removed (merged or obsoleted) over the past two-year period, for each aspect of GO. In
18 the molecular function and cellular component aspects of the ontology, term creation versus
19 term obsolescence have approximately balanced each other, such that the number of terms in
20 these two branches has remained roughly constant. The most significant changes have been in
21 the biological process aspect of GO, with a net decrease of over 800 terms.

22
23
24 Many of these revisions result from global reviews of the ontology to address clear
25 inconsistencies in usage, and changes in annotation practices. Terms that have been removed
26 from the ontology over the last two years fall into several different categories, including:

- 27 - **terms that correspond to phenotypes**, and for which the understanding of the process
28 was previously too incomplete to annotate to a different term. Examples include:
29 *regulation of spindle density* (GO:0090225); *age-dependent general metabolic decline*
30 (GO:0007571).
- 31 - terms that are combinations of multiple GO terms, that **can now be represented more**
32 **precisely using GO-CAM models**. Examples: *chromatin remodeling in response to*
33 *cation stress* (GO:0043156) and *regulation of cyclin-dependent protein serine/threonine*
34 *kinase activity involved in G2/M transition of mitotic cell cycle* (GO:0031660).
- 35 - **revisions based on updated knowledge**, either by GO editors, by authoritative
36 databases, or in the literature. For example alpha-taxilin (UniProt:P40222) was originally
37 thought to be the high molecular weight interleukin-14 (Ambrus *et al.* 1993); an erratum
38 was later published (Ambrus *et al.* 1996) indicating that the open reading frame had
39 been incorrectly predicted. Hence all terms mentioning interleukin-14 have been
40 obsoleted (*interleukin-14 binding* (GO:0019974), *interleukin-14 production*
41 (GO:0032617), and 5 other terms). GO terms created from articles that have been
42 retracted and for which no other supporting evidence exists are obsoleted, for example
43 *CDP-acylglycerol O-arachidonoyltransferase activity* (GO:0047193) (Thompson and Zuk
44 1983). Some terms have been obsoleted from authoritative databases, for example

1 EC:1.3.1.59 was removed from the Enzyme Commission database, and the
 2 corresponding GO term *1,6-dihydroxy-5-methylcyclohexa-2,4-dienecarboxylate*
 3 *dehydrogenase activity* (GO:0018512) was obsoleted in GO.

- 4 - **single step reactions in the biological process aspect of the ontology:** There were
 5 many instances in the GO where a molecular function could be represented as both a
 6 molecular function and a biological process, for example 'histone kinase activity' and
 7 'histone phosphorylation'. This was useful when fewer activities were characterized at
 8 the molecular level, and the best level of resolution for many experiments was that the
 9 gene has some uncharacterized role that led to histone phosphorylation, for example.
 10 However, with increasingly detailed molecular data, the redundancy between MF and BP
 11 annotations became unnecessary and the value of having a similar term in both aspects
 12 of the ontology led to inconsistency. This is an ongoing project and many BP terms still
 13 need to be obsoleted for this reason.
- 14 - **terms that conflate more than one ontological aspect:** *ubiquinone biosynthetic*
 15 *process monooxygenase activity* (GO:0015997) included a biological process within a
 16 molecular function; *MAP kinase phosphatase activity involved in regulation of innate*
 17 *immune response* (GO:0038078) included a molecular function within a biological
 18 process; and *histone deacetylation at centromere* (GO:0031059) represented all three
 19 aspects: a molecular function in the biological process branch of the ontology (histone
 20 acetylation) that also included cellular component information (centromere).
- 21 - **misclassified terms**, for example *urea homeostasis* (GO:0097274) and *creatinine*
 22 *homeostasis* (GO:0097273): while these compounds are important medical biomarkers,
 23 the normal process they measure is proper renal function, therefore these terms have
 24 been obsoleted. Annotations have been re-housed under *renal tubular secretion*
 25 (GO:0097254) (or one of its children) or removed if the paper supporting the annotation
 26 did not allow one to infer the process that affected the circulating levels of urea or
 27 creatinine.
- 28 - **reaction mechanisms:** *primary charge separation* (GO:0009766) and *enzyme active*
 29 *site formation* (GO:0018307) were obsoleted because they represent substeps of
 30 reactions which are beyond the scope of GO.
- 31 - **protein modifying activity terms that mention specific substrates**, for example
 32 *[cytochrome c]-arginine N-methyltransferase activity* (GO:0016275), which is captured
 33 by the more general *arginine N-methyltransferase activity* (GO:0016274). Substrates can
 34 be captured with the 'has input' relationship in GO-CAM models and in annotation
 35 extensions. The exception to this is the histone code: for GO to represent this important
 36 mechanism of gene expression and chromatin structure mechanism, specific activities
 37 are created for known histone modifications, for example: *histone H2AR3*
 38 *methyltransferase activity* (GO:0070612) and *histone H3T3 kinase activity*
 39 (GO:0072354).
- 40 - **experimental assays and non-physiological substrates:** some experiments are
 41 easier to perform using analogs of physiological substrates. GO used to capture this
 42 information, but is now moving away from this and removing any term that represents an
 43 experiment rather than its biological conclusion. An example is *rubidium ion transport*
 44 (GO:0035826): rubidium is used as a tracer for potassium ions (Gill *et al.* 2004), but has

1 no physiological role in itself. Another example is *regulation of nucleosome density*
2 (GO:0060303), which measures the degree of compaction of chromatin, and is a readout
3 for heterochromatin assembly or disassembly.
4

5 Concomitantly with these term obsoletions, many new terms have been added to the ontology in
6 the past two years. An example is *molecular condensate scaffold activity* (GO:0140693) for
7 proteins that nucleate condensates that mediate liquid phase transition. This latter term
8 represents a recent advance in the understanding of the organization of cellular biochemistry
9 (Banani *et al.* 2017).
10

11 We have also clarified the level of specificity at which molecular function terms should be
12 represented in GO. For example, we now strive to create GO terms that represent the range of
13 *in vivo* substrate specificity of an enzyme or transporter. This is in contrast to earlier guidelines,
14 in which a GO term was created for each separate molecular substrate tested in a single,
15 isolated experimental assay or result, which could include non-physiological substrates. With
16 recent improvements in experimental technologies and practices, it is now often possible to
17 annotate with a concept that more closely matches the biological substrate specificity range of a
18 protein. Therefore, while GO makes cross-references to EC (Enzyme Commission) (McDonald
19 and Tipton 2014), Rhea (Bansal *et al.* 2022), KEGG (Kanehisa *et al.* 2022), and MetaCyc
20 (Altman *et al.* 2013), GO does not necessarily create a different term for each of the reactions
21 represented in these resources for each substrate on which a molecular function acts. For
22 example, the GO term *3-oxoacyl-[acyl-carrier-protein] reductase (NADPH) activity*
23 (GO:0004316) represents the fact that the same gene product has a broad specificity toward 3-
24 oxo-acyl groups, and therefore we have obsoleted the more specific GO terms that refer to only
25 one specific substrate, such as *3-oxo-cis-Delta9-hexadecenoyl-[acp] reductase activity*
26 (GO:0102072), *3-oxo-glutaryl-[acp] methyl ester reductase activity* (GO:0102131), and *3-oxo-*
27 *pimeloyl-[acp] methyl ester reductase activity* (GO:0102132). For broad-specificity enzymes and
28 transporters, the activity on a specific substrate in a specific pathway can be captured by
29 biocurators in a GO-CAM (Thomas *et al.* 2019) or an annotation extension (Gene Ontology
30 Consortium 2010) rather than in a GO term.
31

32 Annotations

33 A GO annotation is a statement asserting that a particular gene or gene product has a particular
34 functional characteristic (GO term), examples are shown in Figure 1B. New annotations are
35 continually added to the knowledgebase. In the past two years, experimentally-supported gene
36 function annotations have been added from over 10,000 scientific papers. As of November
37 2022, the GO knowledgebase contains experimental knowledge from almost 173,000 papers.
38 GO annotations derived from experimental data are added primarily by the annotation groups in
39 the GO Consortium, which typically curate biological knowledge by organism (Table 2).
40

41 GO annotations are also regularly reviewed, and may be edited or removed from the
42 knowledgebase for various reasons, particularly when ontology terms are revised (see

1 "Ontology" section above) or when annotations are invalidated by later experimental data.
2 Annotations to terms that will be obsoleted are manually reviewed and annotations are made to
3 a different term whenever possible. For example, when we edited the ontology for histone
4 modifications, over 2,000 annotations to the obsoleted terms were manually reviewed, and
5 histone modifying enzymes were re-annotated to the appropriate molecular function term, while
6 annotations from indirect effects were either removed or re-annotated to different, appropriate
7 GO terms. More minor annotation reviews occur regularly.

8
9 The Phylogenetic Annotation with GO project (see below) involves an integrated biocurator
10 review of annotations that has provided additional quality control. The GO user community also
11 plays an important role in identifying incorrect annotations. Because each annotation can be
12 traced to the published paper containing the underlying evidence or describing a method used
13 to infer the annotation, users can quickly verify the accuracy of a given annotation. Potential
14 errors can be reported by clicking on the "Help" link at the top of the GO homepage
15 (<http://geneontology.org>). In addition, authors of a paper used to create GO annotations can
16 easily retrieve and review all annotations from a given paper and suggest changes; this can be
17 done from the PubMed abstract page (e.g. the PubMed page for PMID:20516198 (Lydeard *et*
18 *al.* 2010)) by clicking on LinkOut and then the "Gene Ontology" link.

19
20 **Phylogenetic Annotations as a source of highly reviewed annotations.** The Phylogenetic
21 Annotation using GO (PAN-GO) project creates a set of biocurator-reviewed, selected GO
22 annotations. This creates sets of improved, augmented GO annotations for genes in the
23 reviewed families. The PAN-GO process is described in detail in (Gaudet *et al.* 2011). Briefly,
24 using the PAINT software tool, a biocurator reviews all experimentally-supported GO
25 annotations collected for all members of a protein family, in the context of a phylogenetic tree
26 from the PANTHER resource (Thomas *et al.* 2022). They then select the most informative and
27 non-redundant GO terms that represent the gene's functional characteristics. Biocurators then
28 model the evolution of these characteristics in the tree by specifying branches along which the
29 GO terms were gained or lost, taking into account events such as duplications, mutations,
30 horizontal gene transfers, as well as taxonomic specificity. This allows for different members of
31 the same family to be annotated with different GO terms when justified by the experimental
32 data. All PAN-GO annotations can be traced to experimental evidence in one or more related
33 genes. To date, a total of 8,196 protein families (out of 11,719 families with experimental data)
34 have been curated. The PAN-GO curation effort has prioritized human gene-containing families,
35 though many other families have also been curated. As a result, annotation coverage of a
36 genome generally depends on how closely related it is to humans. PAN-GO annotations are
37 available for 82% of human genes (compared to 68% with experimental evidence alone). Other
38 vertebrate genomes have similarly high coverage, with genomes from other taxa covered at
39 lower but still substantial levels (Tables 3 and 4). PAN-GO annotations are updated at each GO
40 release, and are included in the standard, downloadable GO annotation files. These annotations
41 can be identified by the "IBA" (inferred from biological ancestor) evidence code, and are
42 available for the 142 organisms included in PANTHER gene families
43 (<http://pantherdb.org/panther/summaryStats.jsp>).

44

1
2 **Protein binding and protein-containing complex annotations.** We suggest that users should
3 be particularly cautious when using GO annotations directly to the term *protein binding*
4 (GO:0005515; see Table 2). These are highly specific annotations that include the protein
5 binding partner in another field of the annotation (not in the GO term itself), and should not be
6 used in applications such as gene set enrichment analysis. Instead, they are recommended for
7 applications such as protein-protein interaction network construction for human proteins (which
8 represent the vast majority of direct protein binding annotations in the knowledgebase). Since all
9 protein functions encompass some type of binding (to a substrate, or to another protein), GO
10 strives to describe the molecular activity of proteins using at least one term that is not only
11 under the *binding* branch of GO; see also "non-catalytic molecular functions" section above.
12 Therefore, *binding* (GO:0005488) in isolation can be considered a limited functional description
13 and is represented as a distinct branch of GO molecular function.

14
15 **Annotation evidence.** All annotations are supported by evidence, comprising two fields in the
16 annotation file (Figure 1B): an *evidence code* that describes the type of evidence and a
17 *reference* that lists a persistent identifier for tracing the source (provenance) of the original data.
18 It has often been asserted that the most reliable annotations are those made using an
19 experimental evidence code. However, we suggest that users take into account **the type of**
20 **experimental evidence** and the **level of review of the annotation** (Table 5). Some types of
21 experimental evidence, such as inference from a gene expression pattern (IEP), mutant
22 phenotype (IMP) or genetic interaction (IGI) can often be suggestive of function but not definitive
23 when considered in isolation; other annotations for the same gene are often useful to help
24 interpret these annotations. "High-throughput" evidence codes should be treated with particular
25 care. These codes (beginning with the letter H) denote experiments in which many genes are
26 analyzed at the same time, and these annotations are not individually reviewed by either the
27 paper's authors or GO Consortium biocurator (Attrill *et al.* 2019). Conversely, many non-
28 experimental evidence types are carefully reviewed by experts. Phylogenetic annotations (IBA
29 evidence code) are based on integration and expert assessment of experimental annotations,
30 and thus are *individually reviewed* twice: once in making the annotation from published
31 experimental results, and once in the context of all annotations for related genes (Gaudet *et al.*
32 2011). While annotations using the Inferred from electronic annotation (IEA) evidence code are
33 considered automated, most implement expert review of a subset of annotations to minimize
34 false positives (for example, UniRule (MacDougall *et al.* 2021) and InterPro2GO (Paysan-
35 Lafosse *et al.* 2022)). The GOC considers these annotations to be accurate though they are
36 often less specific than other annotations

37
38 GO evidence codes correspond to a subset of the terms found in the Evidence and Conclusion
39 Ontology (ECO) (Nadendla *et al.* 2022). Combinations of particular GO internal references
40 (GO_REFs) and evidence codes are also mapped to specific ECO terms
41 (<https://github.com/evidenceontology/evidenceontology/blob/master/gaf-eco-mapping.txt>). Users
42 needing to map granular ECO terms to GO evidence code abbreviations can use the mapping
43 file provided by ECO ([https://github.com/evidenceontology/evidenceontology/blob/master/gaf-
44 eco-mapping-derived.txt](https://github.com/evidenceontology/evidenceontology/blob/master/gaf-eco-mapping-derived.txt)).

1 GO Causal Activity Models (GO-CAMs)

2
3 GO-CAMs are models of causal influences between gene products (Thomas *et al.* 2019), or
4 pathways. More precisely, a GO-CAM links the **activities** (GO molecular functions) of gene
5 products together by **causal relations** that specify the effect of one activity on the other. Each
6 element of a GO-CAM is an instance of an ontology class or other standard database identifier,
7 so GO-CAMs are highly structured and amenable to computational analysis. The basic unit of a
8 GO-CAM is a “gene product activity unit”, which combines a GO molecular function annotation
9 (molecular activity), together with GO cellular component (location) and GO biological process
10 (larger functional module) annotations that provide the biological context of the activity. The
11 context can be further specified with other ontologies to capture cell type (using the Cell Type
12 Ontology (Diehl *et al.* 2016)), tissue/anatomical location (using several different ontologies
13 depending on the species, e.g. Uberon (Mungall *et al.* 2012) for most vertebrates, other
14 metazoan ontologies such as the *Drosophila* anatomy ontology (Costa *et al.* 2013), *C. elegans*
15 anatomy ontology (Lee and Sternberg 2003), or non-animal ontologies as the Plant Ontology
16 (Cooper *et al.* 2018)), or a temporal period (e.g. GO biological phase). Activity units are linked
17 together by causal relationships from the Relations Ontology (Smith *et al.* 2005) to capture how
18 they interact to impact larger pathways, modules or processes.
19

20
21 As of November 2022, GO Consortium annotation groups have created over 300 GO-CAM
22 models that describe molecular pathways (defined as containing at least three distinct gene
23 product activities linked into a causal chain). These models reflect curation priorities of the
24 contributing groups. Most of the available GO-CAMs are for processes in human or mouse, with
25 a limited number in zebrafish, *D. melanogaster* and *C. elegans*. Many of the human GO-CAMs
26 describe chromatin-mediated regulation of gene expression and immune response pathways,
27 while the mouse GO-CAMs focus on metabolic and signaling pathways. GO-CAMs are
28 accessible from the GO website homepage, by clicking on the “Browse GO-CAMs” link. GO-
29 CAMs can be viewed as pathway diagrams (Figure 2) and are currently available on GitHub at
30 <https://github.com/geneontology/noctua-models>.

31 Community Collaborations

32 The GO Consortium collaborates with experts in specific areas of molecular and cellular biology
33 to systematically update and improve their representation in the ontology and the corresponding
34 GO annotations and GO-CAMs. We recently revised the representation for transcription factors
35 and transcriptional regulation in collaboration with the GREEKC Consortium (Kuiper *et al.* 2022).
36 Additional collaborative projects include working with the DisProt project (Quaglia *et al.* 2022) on
37 improving the ontology and annotations for intrinsically disordered proteins; revising processes
38 that involve molecular pathways between interacting species, such as viral infection processes;
39 and integrating the gene ontology and annotations with external biochemical databases.
40

1 In 2021 the Gene Ontology started a collaboration with DisProt (<https://disprot.org/>) – the gold
2 standard database of manually curated annotations from the literature for Intrinsically
3 Disordered Proteins (IDPs). IDPs lack a stable three-dimensional structure and are
4 characterized by highly flexible and unstructured segments, i.e. intrinsically disordered regions
5 (IDRs). DisProt has developed a custom ontology, the Intrinsically Disordered Proteins Ontology
6 (IDPO), and used it to annotate the structural states of IDPs. The GO Consortium and DisProt
7 have collaborated to refactor IDPO and map the IDPO terms to GO terms whenever possible
8 (those related to functions and interactions of IDPs). The collaboration between the GO
9 Consortium and the DisProt database included the creation and addition of new GO terms to
10 align with already existing IDPO terms that were not yet available in GO. These newly created
11 terms also include the *molecular function activator* (GO:0140677) and *molecular function*
12 *inhibitor* (GO:0140678) terms, used to annotate molecular function regulators that
13 activate/inhibit or increase/decrease the activities of their targets via non-covalent binding that
14 does not result in covalent modification to the target. This collaboration resulted in more
15 accurate and detailed annotation of the modes of action of IDPs, e.g. *localization* (GO:0051179,
16 IDPO:00010) and *DNA binding* (GO:0003677, IDPO:00065), as well as providing GO
17 annotations. Currently more than 1,000 expert-curated annotations from DisProt are available in
18 the GO knowledgebase, comprising more than 860 molecular functions, 200 biological
19 processes and 10 cellular component annotations. The only terms in IDPO that could not be
20 mapped to GO were those describing self-regulatory (e.g. *self-activation* and *self-inhibition*) and
21 intrinsic disorder-specific functions (i.e. *entropic chains*), so these annotations are available only
22 in DisProt.

23 Multiorganism interactions

24 A group that includes experts from within and outside the GO Consortium have been working
25 together to improve and simplify the representation of interactions between organisms, including
26 medically and agriculturally important host-pathogen interactions. Examples of these
27 interactions include how a symbiont such as a virus enters its host, how the host's immune
28 response recognizes and defends the body against a potentially harmful organism, and also
29 beneficial interactions such as how plants form a symbiosis with nitrogen-fixing bacteria. The
30 goal of this project is to revise the host-symbiont branch of GO *biological process* to reflect the
31 current scientific knowledge in the field, and to ensure that genes are properly annotated to the
32 new ontology terms and structure, building on previous work undertaken as part of the PAMGO
33 consortium (Tyler *et al.* 2009). Symbionts in GO are broadly defined to include pathogens that
34 infect a host organism. We expect that this revision will improve GO-based analyses of
35 molecular studies of pathogens, the mechanisms by which they infect host cells, and host
36 response processes. A major change is that the branch of GO under *biological process involved*
37 *in interspecies interaction between organisms* (GO:0044419) has been reorganized. It now
38 reflects important concepts such as the types of biological programs used by symbionts to
39 enable infection, and by hosts to prevent or manage infection, such as *disruption of cellular*
40 *component of another organism* (GO:0140975), *formation of structure involved in a symbiotic*
41 *process* (GO:0044111), *killing of cells of another organism* (GO:0031640), and *modulation of*

1 *process of another organism* (GO:0035821). Each of these terms has multiple, more specific
2 subclass terms.

3
4 One challenge in this area was that some previous GO annotations for pathogen genes used
5 terms that apply to normal host processes, such as regulation of defense response processes.
6 Thus, it was not clear whether the pathogen gene was regulating its own defense process or
7 that of a host. With the new ontology terms and structure, these distinctions are clear for both
8 GO biocurators and users of GO. In general, it was important to clearly represent that certain
9 symbiont-initiated processes hijack various host cellular processes. This includes mechanisms
10 to enter and exit the cell, either by binding to host membrane proteins or using the intracellular
11 transport machinery, and using the host cellular machinery for genome replication, as well as
12 transcription and translation. We have obsoleted terms that do not clearly distinguish hijacking
13 with the functions that a host gene performs for the host organism, such as *dissemination or*
14 *transmission of symbiont from host by vector* (GO:0044008) and *positive regulation of viral*
15 *release from host cell* (GO:1902188). Conversely, a pathogenic symbiont triggers innate
16 responses in the host that are not the evolved role of these symbiont proteins, such as *induction*
17 *by symbiont of host cytokine production* (GO:0036523) and *pathogen-associated molecular*
18 *pattern dependent induction by symbiont of host innate immune response* (GO:0052033) - these
19 are not functions that a symbiont protein performs to enable its own survival and reproduction.

20 Integration with biochemical knowledgebases

21 GO also works closely with specialized databases and knowledgebases to ensure knowledge is
22 both complete and consistent. For accurate representation of biochemical aspects of gene
23 function, we work closely with the Rhea database of reactions (Bansal *et al.* 2022) and the
24 ChEBI ontology of chemical entities (Hastings *et al.* 2016). Rhea provides precise
25 representations of *in vivo* biochemical reactions, including precise chemical entity participants
26 and their stoichiometry. Rhea uses ChEBI terms to represent chemical entities in a
27 standardized, consistent manner. The Rhea database overlaps in content with the catalytic
28 activity branch of the Gene Ontology, but provides additional detailed reaction information, and
29 in some cases provides additional specificity. We have improved GO mappings to Rhea, which
30 now covers 4399 GO catalytic activities (in the MF branch of GO). These mappings allow for
31 non-exact matches when the chemical specificity differs between GO and Rhea. For example,
32 Rhea has two reactions, each referring to a different type of beta glucoside (RHEA:69647 and
33 RHEA:69655, narrow match), whereas GO:0008422, beta-glucosidase activity, covers both
34 substrates, as no known enzyme is specific for just one of them. We have recently used the
35 Rhea-GO mappings to include additional linkages between GO molecular function terms and
36 ChEBI terms in the go-plus release (see below). Previously ChEBI terms were linked only to
37 general terms in the GO biological process branch (e.g. between folate transport and folate), but
38 the additional Rhea linkages have added a total of 4334 distinct chemical entities linked via
39 20,307 relationships. The extensive linkage to chemical entities opens opportunities for using
40 GO in other applications, e.g. metabolomics analyses.

1 Accessing and downloading GO data

2 Browsing GO and its annotations

3 GO and associated annotations can be searched directly from the Gene Ontology home page
4 (<http://geneontology.org/>), queried using the AmiGO browser
5 (<http://amigo.geneontology.org/amigo>) or the QuickGO tool (<https://www.ebi.ac.uk/QuickGO/>)
6 (Munoz-Torres and Carbon 2017). Gene set enrichment analysis is also directly accessible from
7 the Gene Ontology home page, which runs the PANTHER gene analysis tool at
8 <http://pantherdb.org/webservices/go/overrep.jsp> (Mi *et al.* 2019).

9 Ontology downloads

10 GO provides three editions of the ontology on its download site
11 (<http://geneontology.org/docs/download-ontology/>) to accommodate various applications: go-
12 basic, go, and go-plus (Table 6). All GO terms, including obsolete terms and term metadata
13 such as definitions, cross references, synonyms, are available in all three editions. These
14 editions differ in the set of relations they contain:

- 15 • **go-basic** contains the types of information that has been available for GO from the
16 beginning of the project, hence only contains is a, part of, regulates, negatively regulates
17 and positively regulates relationships and excludes relationships that cross different
18 aspects (BP, MF or CC) of the ontology. This edition of the ontology is guaranteed to be
19 acyclic and can safely be used to selectively propagate annotations across any relation.
20 It is recommended for most GO-based annotation tools.
- 21 • **go** additionally includes has part and occurs in relationships that link terms across
22 different aspects of the ontology (for example, a biological process can have a *has part*
23 relation to a molecular function term, or an *occurs in* relation to a cellular component).
24 This edition is not acyclic and annotations should not be propagated across all the
25 relationship types it contains. This edition should not be used in most software tools that
26 rely on the Gene Ontology.
- 27 • **go-plus** is the fully axiomatized edition of the ontology, and includes cross-ontology
28 relationships to external ontologies including ChEBI, Cell Ontology and Uberon.

29 Ontology subsets (GO slims)

30 GO subsets are condensed versions of the GO containing a portion of the terms, which are
31 specified by tags within the ontology that indicate if a given term is a member of a particular
32 subset. GO subsets are particularly useful for providing a global overview of the functions of all
33 the genes in a genome, and even for all the functions of a single gene. range of functions and
34 processes found in a given clade or organism's genome. We have recently revised the "GO
35 Generic subset", a subset maintained by the GO consortium that aims to be general and
36 applicable to any species. We have tested that the subset covers as many gene products as
37 possible in various organisms (human, *D. melanogaster*, fission yeast, *A. thaliana*, *E. coli*) with

1 as little redundancy as possible. This new GO Generic Subset contains 75 biological process
2 terms, 40 molecular function terms, and 29 cellular component terms. The GO generic subset
3 can be accessed at http://current.geneontology.org/ontology/subsets/goslim_generic.obo.
4 Versions in .owl, .json and .tsv are also available from
5 <http://current.geneontology.org/ontology/subsets/index.html>.

6
7 As part of the Alliance of Genome Resources, we have developed a widget that provides a
8 graphical visualization of a gene's function in a 'ribbon'-like display (Figure 3). The widget can
9 be customized to use any GO subset, and uses the goslim_agr subset by default. This widget is
10 implemented in the Alliance gene pages and in the UniProt entry pages. It accesses GO
11 annotations using the GO API (application programming interface) and can be easily added to
12 any webpage.
13

14 GO Annotations

15 The two major sites for downloading GO annotations are geneontology.org and UniProt-GOA.
16 **geneontology.org** is the website developed by the GO Consortium. This downloads site
17 (<http://current.geneontology.org/products/pages/downloads.html>) provides a total of 7.5 million
18 human and model organism annotations contributed by multiple groups. It contains all manually
19 reviewed GO annotations, and electronic (computationally predicted) annotations for the most
20 commonly used organisms. For model organisms, all annotations use gene identifiers from the
21 authoritative database (for example, FlyBase (FBgn), WormBase (WBGene), and SGD (S)).
22 Human and other organisms without an authoritative dedicated database are represented by
23 UniProtKB accession numbers. For these organisms, the GO website provides annotations to
24 UniProt reference proteomes (https://www.uniprot.org/help/reference_proteome), which are
25 generally one entry per gene, thus limiting redundancy in annotations. **UniProt-GOA**
26 (https://www.ebi.ac.uk/GOA/uniprot_release) contains 1 billion annotations for all entries in
27 UniProt (1,264,340 taxa), covering both reviewed (Swiss-Prot) entries of UniProt, and
28 unreviewed (TrEMBL) entries. All annotations for model organism genes are converted to
29 UniProt protein identifiers. For most organisms, all annotations are electronic annotations
30 generated via various pipelines (see above for evidence codes and references for different
31 methods). In addition to these resources, GO annotations are also viewable in a number of
32 biological databases, including model organism databases, UniProt (UniProt: the universal
33 protein knowledgebase 2017), NCBI (Sayers *et al.* 2020), and The Alliance of Genome
34 Resources (Alliance of Genome Resources Consortium 2022). These sites show GO
35 annotations in the broader context of a gene product's expression pattern, phenotypes,
36 metabolic and signaling pathways, etc.
37

1 Conclusions and future directions

2 The extensive and wide-ranging use of the GO knowledgebase, evidenced by its recent, peer-
3 reviewed designation as a Global Core Biodata Resource ([https://globalbiodata.org/scientific-](https://globalbiodata.org/scientific-activities/global-core-biodata-resources/)
4 [activities/global-core-biodata-resources/](https://globalbiodata.org/scientific-activities/global-core-biodata-resources/)), demands its continued development and expansion.
5 We are focusing on several high priority areas of development for the near future. For pathways,
6 we will continue to accumulate GO-CAM models. The UniProt/Swiss-Prot curation team has
7 ramped up their production of GO-CAM models and we expect to add models at a rapid rate. In
8 parallel, we have started converting Reactome pathways into GO-CAMs (Good *et al.* 2021) and
9 expect to release GO-CAM representations of most Reactome metabolic pathways in the near
10 future. This will provide a complementary, causal flow representation of the chemical reaction-
11 centered representation in Reactome. Conversion of Reactome signaling pathways is more
12 challenging, and will be released somewhat later. We are also working on converting the
13 YeastPathways resource (<https://pathways.yeastgenome.org>) into GO-CAMs, making a large
14 number of yeast metabolic pathways available. The increasing number of GO-CAM models will
15 allow us to expand on the utility of these highly structured pathway and process representations.
16 Some potential areas are automated pathway visualization; using the causal links and more
17 granular gene sets to enhance enrichment analysis; and better generation of automated
18 descriptions of gene function (e.g., (Kishore *et al.* 2020)).
19

20 With respect to ontology development, in addition to continuing to revise the ontology in
21 response to recent discoveries, we see an immediate need for clearly delineating the level of
22 biological organization at which a function is described. This includes distinguishing molecular
23 functions from biological processes, and distinguishing biological processes that occur at the
24 level of individual cells, versus those that occur at the level of multicellular organisms. For
25 example, the term homeostasis—the maintenance of a roughly steady level of a molecule or ion—
26 is used very broadly in the literature to refer to both processes that maintain a steady-state level
27 within a cell, and processes that maintain a steady state in blood or other fluid that is
28 transported within a multicellular organism. Even in some publications, it is difficult to know
29 which type of homeostasis is being tested.
30

31 We will continue to make the GO knowledgebase easier to use, and more community-driven.
32 One near-term priority is to make annotations available for download by species, with a single
33 identifier for each distinct gene. We are also planning to create quick-start guides for common
34 GO use cases, in both written and video form. The immense user base of the GO and the need
35 for much improvement and extension drives us to consider how to expand the number of people
36 that contribute to the GO. From its inception the GO has been a large, open, community project.
37 However, we are planning additional routes through which the broader GO user community can
38 contribute their expert feedback and knowledge to GO, improving the resource for all users. For
39 now, users are encouraged to contact the GO Helpdesk (<http://help.geneontology.org/>) with any
40 questions, or to report any GO ontology terms or annotations that may be inaccurate or difficult
41 to interpret.
42

1 Data Availability

2 All Gene Ontology code and resources are freely available for download and reuse. Software
3 (<https://github.com/geneontology/>) is under the BSD 3-Clause open-source license. Downloads
4 are available under the CC BY 4.0 license from <http://geneontology.org/docs/downloads/>

5 Funding

6 The core funding for the GOC is from the National Human Genome Research Institute
7 (U41HG002273, U24HG012212). Curation activities supported by National Human Genome
8 Research Institute grants U24HG002659 (ZFIN), U24HG002223 (WormBase), U41HG000739
9 (FlyBase), U24HG001315 (SGD), U24HG000330 (MGD), U24HG012198 (Reactome curation),
10 U24HG011851 (Reactome - GO harmonization) and grant R01HL064541 from the National Heart,
11 Lung and Blood Institute (RGD). Additional funding for GO curation at FlyBase is provided by UK
12 Medical Research Council Award MR/W024233/1. PomBase is supported by Wellcome Trust
13 218236/Z/19/Z. Xenbase is supported by grant P41 HD064556 from the Eunice Kennedy Shriver
14 National Institute of Child Health and Human Development. Functional Gene Annotation,
15 University College London is supported by National Institute for Health Research University
16 College London Hospitals Biomedical Research Centre. Planteome and Plant Reactome are
17 supported by National Science Foundation awards #1340112, #1127112, and USDA-ARS. Some
18 software development was funded by U24HG010859 (Alliance of Genome Resources Central).
19 The TAIR project is funded by academic, institutional, corporate, and individual subscriptions;
20 TAIR is administered by the 501(c)(3) non-profit Phoenix Bioinformatics. Chris Mungall, Seth
21 Carbon, and Sierra Moxon were supported in part by Director, Office of Science, Office of Basic
22 Energy Sciences of the U.S. Department of Energy Contract No. DE-AC02-05CH11231
23

24 UniProt is funded by National Human Genome Research Institute (NHGRI), Office of Director
25 [OD/DPCPSI/ODSS]; National Institute of Allergy and Infectious Diseases (NIAID), National
26 Institute on Aging (NIA), National Institute of General Medical Sciences (NIGMS), National
27 Institute of Diabetes and Digestive and Kidney Diseases (NIDDK), National Eye Institute (NEI),
28 National Cancer Institute (NCI), National Heart, Lung, and Blood Institute (NHLBI) of the National
29 Institutes of Health [U24HG007822]; Biotechnology and Biological Sciences Research Council
30 [BB/T010541/1, BB/S01781X/1]; Open Targets; Swiss Federal Government through the State
31 Secretariat for Education, Research and Innovation SERI; European Molecular Biology
32 Laboratory core funds.

33

34 Summary

35 The Gene Ontology (GO) knowledgebase has broad applications in genetic and genomic
36 research, and has been continually updated and improved for more than 20 years. We describe

1 the latest improvements to the GO resource, as well as giving an overview of the current contents
2 of the knowledgebase. We also include guidance on the use of GO, annotations and the growing
3 number of pathways represented as GO-Causal Activity Models (GO-CAMs).

4 References

- 5 Alliance of Genome Resources Consortium. Harmonizing model organism data in the Alliance
6 of Genome Resources. *Genetics* 2022;**220**, DOI: 10.1093/genetics/iyac022.
- 7 Altman T, Travers M, Kothari A *et al.* A systematic comparison of the MetaCyc and KEGG
8 pathway databases. *BMC Bioinformatics* 2013;**14**:112.
- 9 Ambrus JL Jr, Pippin J, Joseph A *et al.* Identification of a cDNA for a human high-molecular-
10 weight B-cell growth factor. *Proc Natl Acad Sci U S A* 1993;**90**:6330–4.
- 11 Ambrus JL Jr, Pippin J, Joseph A *et al.* Identification of a cDNA for a human high molecular-
12 weight B-cell growth factor. *Proc Natl Acad Sci U S A* 1996;**93**:8154.
- 13 Ashburner M, Ball CA, Blake JA *et al.* Gene ontology: tool for the unification of biology. The
14 Gene Ontology Consortium. *Nat Genet* 2000;**25**:25–9.
- 15 Attrill H, Gaudet P, Huntley RP *et al.* Annotation of gene product function from high-throughput
16 studies using the Gene Ontology. *Database* 2019;**2019**, DOI: 10.1093/database/baz007.
- 17 Banani SF, Lee HO, Hyman AA *et al.* Biomolecular condensates: organizers of cellular
18 biochemistry. *Nat Rev Mol Cell Biol* 2017;**18**:285–98.
- 19 Bansal P, Morgat A, Axelsen KB *et al.* Rhea, the reaction knowledgebase in 2022. *Nucleic
20 Acids Res* 2022;**50**:D693–700.
- 21 Basu S, Fey P, Jimenez-Morales D *et al.* dictyBase 2015: Expanding data and annotations in a
22 new software environment. *Genesis* 2015;**53**:523–34.
- 23 Bult CJ, Blake JA, Smith CL *et al.* Mouse Genome Database (MGD) 2019. *Nucleic Acids Res*
24 2019;**47**:D801–6.
- 25 Cerqueira GC, Arnaud MB, Inglis DO *et al.* The Aspergillus Genome Database: multispecies
26 curation and incorporation of RNA-Seq data to improve structural gene annotations. *Nucleic
27 Acids Res* 2014;**42**:D705–10.
- 28 Davis P, Zarowiecki M, Arnaboldi V *et al.* WormBase in 2022-data, processes, and tools for
29 analyzing *Caenorhabditis elegans*. *Genetics* 2022;**220**, DOI: 10.1093/genetics/iyac003.
- 30 Del Toro N, Shrivastava A, Ragueneau E *et al.* The IntAct database: efficient access to fine-
31 grained molecular interaction data. *Nucleic Acids Res* 2022;**50**:D648–53.
- 32 Diehl AD, Meehan TF, Bradford YM *et al.* The Cell Ontology 2016: enhanced content,
33 modularization, and ontology interoperability. *J Biomed Semantics* 2016;**7**:44.
- 34 Duck G, Nenadic G, Filannino M *et al.* A Survey of Bioinformatics Database and Software
35 Usage through Mining the Literature. *PLoS One* 2016;**11**:e0157989.

- 1 Fabregat A, Jupe S, Matthews L *et al.* The Reactome Pathway Knowledgebase. *Nucleic Acids*
2 *Res* 2018;**46**:D649–55.
- 3 Fortriede JD, Pells TJ, Chu S *et al.* Xenbase: deep integration of GEO & SRA RNA-seq and
4 CHIP-seq data in a model organism database. *Nucleic Acids Res* 2020;**48**:D776–82.
- 5 *Fungal-Anatomy-Ontology: A Structured Controlled Vocabulary for the Anatomy of Fungi.*
6 Github
- 7 Gaudet P, Livstone MS, Lewis SE *et al.* Phylogenetic-based propagation of functional
8 annotations within the Gene Ontology consortium. *Brief Bioinform* 2011;**12**:449–62.
- 9 Gaudet P, Logie C, Lovering RC *et al.* Gene Ontology representation for transcription factor
10 functions. *Biochim Biophys Acta Gene Regul Mech* 2021:194752.
- 11 Gene Ontology Consortium. The Gene Ontology in 2010: extensions and refinements. *Nucleic*
12 *Acids Res* 2010;**38**:D331–5.
- 13 Gill S, Gill R, Wicks D *et al.* Development of an HTS assay for Na⁺, K⁺-ATPase using
14 nonradioactive rubidium ion uptake. *Assay Drug Dev Technol* 2004;**2**:535–42.
- 15 Gkoutos GV, Schofield PN, Hoehndorf R. The anatomy of phenotype ontologies: principles,
16 properties and applications. *Brief Bioinform* 2018;**19**:1008–21.
- 17 Good BM, Van Auken K, Hill DP *et al.* Reactome and the Gene Ontology: Digital convergence
18 of data resources. *Bioinformatics* 2021, DOI: 10.1093/bioinformatics/btab325.
- 19 Haendel MA, Balhoff JP, Bastian FB *et al.* Unification of multi-species vertebrate anatomy
20 ontologies for comparative biology in Uberon. *J Biomed Semantics* 2014;**5**:21.
- 21 Haendel MA, Neuhaus F, Osumi-Sutherland D *et al.* CARO – The Common Anatomy Reference
22 Ontology. In: Burger A, Davidson D, Baldock R (eds.). *Anatomy Ontologies for Bioinformatics:*
23 *Principles and Practice*. London: Springer London, 2008, 327–49.
- 24 Harris MA, Rutherford KM, Hayles J *et al.* Fission stories: Using PomBase to understand
25 *Schizosaccharomyces pombe* biology. , DOI: 10.1101/2021.09.07.459264.
- 26 Hastings J, Owen G, Dekker A *et al.* ChEBI in 2016: Improved services and an expanding
27 collection of metabolites. *Nucleic Acids Res* 2016;**44**:D1214–9.
- 28 Howe DG, Ramachandran S, Bradford YM *et al.* The Zebrafish Information Network: major gene
29 page and home page updates. *Nucleic Acids Res* 2021;**49**:D1058–64.
- 30 Kanehisa M, Furumichi M, Sato Y *et al.* KEGG for taxonomy-based analysis of pathways and
31 genomes. *Nucleic Acids Res* 2022, DOI: 10.1093/nar/gkac963.
- 32 Keseler IM, Mackie A, Santos-Zavaleta A *et al.* The EcoCyc database: reflecting new
33 knowledge about *Escherichia coli* K-12. *Nucleic Acids Res* 2017;**45**:D543–50.
- 34 Kishore R, Arnaboldi V, Van Slyke CE *et al.* Automated generation of gene summaries at the
35 Alliance of Genome Resources. *Database* 2020;**2020**, DOI: 10.1093/database/baaa037.
- 36 Koopmans F, van Nierop P, Andres-Alonso M *et al.* SynGO: An Evidence-Based, Expert-

- 1 Curated Knowledge Base for the Synapse. *Neuron* 2019;**103**:217–34.e4.
- 2 Kuiper M, Bonello J, Fernández-Breis JT *et al.* The gene regulation knowledge commons: the
3 action area of GREEKC. *Biochim Biophys Acta Gene Regul Mech* 2022;**1865**:194768.
- 4 Lamesch P, Berardini TZ, Li D *et al.* The Arabidopsis Information Resource (TAIR): improved
5 gene annotation and new tools. *Nucleic Acids Res* 2012;**40**:D1202–10.
- 6 Lang OW, Nash RS, Hellerstedt ST *et al.* An Introduction to the Saccharomyces Genome
7 Database (SGD). *Methods Mol Biol* 2018;**1757**:21–30.
- 8 Lydeard JR, Lipkin-Moore Z, Sheu Y-J *et al.* Break-induced replication requires all essential
9 DNA replication factors except those specific for pre-RC assembly. *Genes Dev* 2010;**24**:1133–
10 44.
- 11 MacDougall A, Volynkin V, Saidi R *et al.* UniRule: a unified rule resource for automatic
12 annotation in the UniProt Knowledgebase. *Bioinformatics* 2021;**36**:5562.
- 13 McDonald AG, Tipton KF. Fifty-five years of enzyme classification: advances and difficulties.
14 *FEBS J* 2014;**281**:583–92.
- 15 McIntosh BK, Renfro DP, Knapp GS *et al.* EcoliWiki: a wiki-based community resource for
16 Escherichia coli. *Nucleic Acids Res* 2012;**40**:D1270–7.
- 17 Meldal BHM, Bye-A-Jee H, Gajdoš L *et al.* Complex Portal 2018: extended content and
18 enhanced visualization tools for macromolecular complexes. *Nucleic Acids Res* 2019;**47**:D550–
19 8.
- 20 Mi H, Muruganujan A, Huang X *et al.* Protocol Update for large-scale genome and gene function
21 analysis with the PANTHER classification system (v.14.0). *Nat Protoc* 2019;**14**:703–21.
- 22 Mungall CJ, Batchelor C, Eilbeck K. Evolution of the Sequence Ontology terms and
23 relationships. *J Biomed Inform* 2011;**44**:87–93.
- 24 Munoz-Torres M, Carbon S. Get GO! Retrieving GO Data Using AmiGO, QuickGO, API, Files,
25 and Tools. *Methods Mol Biol* 2017;**1446**:149–60.
- 26 Nadendla S, Jackson R, Munro J *et al.* ECO: the Evidence and Conclusion Ontology, an update
27 for 2022. *Nucleic Acids Res* 2022;**50**:D1515–21.
- 28 Natale DA, Arighi CN, Blake JA *et al.* Protein Ontology (PRO): enhancing and scaling up the
29 representation of protein entities. *Nucleic Acids Research* 2017;**45**:D339–46.
- 30 Paysan-Lafosse T, Blum M, Chuguransky S *et al.* InterPro in 2022. *Nucleic Acids Res* 2022,
31 DOI: 10.1093/nar/gkac993.
- 32 Quaglia F, Mészáros B, Salladini E *et al.* DisProt in 2022: improved quality and accessibility of
33 protein intrinsic disorder annotation. *Nucleic Acids Res* 2022;**50**:D480–7.
- 34 Radivojac P, Clark WT, Oron TR *et al.* A large-scale evaluation of computational protein
35 function prediction. *Nat Methods* 2013;**10**:221–7.
- 36 Ramsey J, McIntosh B, Renfro D *et al.* Crowdsourcing biocuration: The Community Assessment

- 1 of Community Annotation with Ontologies (CACAO). *PLoS Comput Biol* 2021;**17**:e1009463.
- 2 Saverimuttu SCC, Kramarz B, Rodríguez-López M *et al.* Gene Ontology curation of the blood-
3 brain barrier to improve the analysis of Alzheimer's and other neurological diseases. *Database*
4 2021;**2021**, DOI: 10.1093/database/baab067.
- 5 Sayers EW, Beck J, Brister JR *et al.* Database resources of the National Center for
6 Biotechnology Information. *Nucleic Acids Res* 2020;**48**:D9–16.
- 7 Sian L, Agapite J, Attrill H *et al.* FlyBase: a guided tour of highlighted features. *Genetics*
8 2022;**220**:iyac035.
- 9 Skrzypek MS, Binkley J, Binkley G *et al.* The Candida Genome Database (CGD): incorporation
10 of Assembly 22, systematic identifiers and visualization of high throughput sequencing data.
11 *Nucleic Acids Res* 2017;**45**:D592–6.
- 12 Smith B, Ceusters W, Klagges B *et al.* Relations in biomedical ontologies. *Genome Biol*
13 2005;**6**:R46.
- 14 Smith JR, Hayman GT, Wang S-J *et al.* The Year of the Rat: The Rat Genome Database at 20:
15 a multi-species knowledgebase and analysis platform. *Nucleic Acids Res* 2020;**48**:D731–42.
- 16 Thomas PD. The Gene Ontology and the Meaning of Biological Function. *Methods Mol Biol*
17 2017;**1446**:15–24.
- 18 Thomas PD, Ebert D, Muruganujan A *et al.* PANTHER: Making genome-scale phylogenetics
19 accessible to all. *Protein Sci* 2022;**31**:8–22.
- 20 Thomas PD, Hill DP, Mi H *et al.* Gene Ontology Causal Activity Modeling (GO-CAM) moves
21 beyond GO annotations to structured descriptions of biological functions and systems. *Nat*
22 *Genet* 2019;**51**:1429–33.
- 23 Thompson W, Zuk RT. Acylation of CDP-monoacylglycerol cannot be confirmed. *J Biol Chem*
24 1983;**258**:9623.
- 25 Tyler BM, Collmer A, Collmer CW *et al.* The PAMGO consortium: unifying themes in microbe-
26 host associations identified through the gene ontology. *BMC Microbiol* 2009;**9**.
- 27 UniProt: the universal protein knowledgebase. *Nucleic Acids Res* 2017;**45**:D158–69.
- 28 Walls RL, Cooper L, Elser J *et al.* The Plant Ontology Facilitates Comparisons of Plant
29 Development Stages Across Species. *Front Plant Sci* 2019;**10**:631.
- 30 Winsor GL, Lo R, Ho Sui SJ *et al.* Pseudomonas aeruginosa Genome Database and
31 PseudoCAP: facilitating community-based, continually updated, genome annotation. *Nucleic*
32 *Acids Res* 2005;**33**:D338–43.

33

1 Tables

2 Table 1. Changes to GO terms in the past two year period. The ontology has undergone
3 substantial revision and improvement, with nearly 2,000 terms added or removed.

GO aspect	Total number of terms	Added terms	Obsoleted terms	Merged terms ¹
Molecular Function	11,271	315	65	143
Cellular Component	4,039	34	19	162
Biological Process	27,993	217	782	254

4 ¹Also includes obsoleted terms that have been replaced by another term.
5
6

7 Table 2. Group contributing literature-based annotations. Includes all annotations traceable to
8 the literature (EXP, including HTP, TAS, NAS, IC, see [http://geneontology.org/docs/guide-go-](http://geneontology.org/docs/guide-go-evidence-codes)
9 [evidence-codes](http://geneontology.org/docs/guide-go-evidence-codes); see below for information). Direct annotations to the term "protein binding" are
10 listed separately, since without information about interacting partner(s), protein binding
11 represents an activity that most proteins possess and therefore the GO class itself provides little
12 information (see text for further description). The statistics for groups that have contributed more
13 than 700 manual annotations. Other contributing groups include: HGNC, JaponicusDB, PHI-
14 base, PAMGO, JCVI, MENGO, and GDB. Current GO Consortium members are labeled with an
15 asterisk. See <http://geneontology.org/docs/annotation-contributors/> for more details.
16

Group	Organism or area of focus	Number of literature-based annotations, excluding direct protein binding	Number of literature-based annotations directly to protein binding
UniProt* (UniProt: the universal protein knowledgebase 2017)	human, and also a wide variety of organisms not covered by other GOC members	185121	30927
MGI* (Bult et al. 2019)	mouse	106435	8051
Reactome* (Fabregat et al. 2018)	human pathways	92178	6
TAIR* (Lamesch et al. 2012)	A. thaliana (model plant)	64633	4695
FlyBase* (Sian et al. 2022)	D. melanogaster (fruit fly)	55203	892
UCL*	human	54595	2935
RGD* (Smith et al. 2020)	rat	47694	1894

SGD* (Lang et al. 2018)	S. cerevisiae (Baker's yeast)	48811	165
ZFIN* (Howe et al. 2021)	zebrafish	28261	488
PomBase* (Harris et al.)	S. pombe (fission yeast)	26128	2201
GeneDB	microbial pathogens	23884	756
ComplexPortal* (Meldal et al. 2019)	protein complexes	18343	0
WormBase* (Davis et al. 2022)	C. elegans (nematode)	17171	560
CGD* (Skrzypek et al. 2017)	C. albicans (yeast pathogen)	17113	0
EcoCyc* (Keseler et al. 2017)	E. coli (bacterium)	13372	829
AgBase	Agricultural animals, primarily chicken	11198	1110
dictyBase* (Basu et al. 2015)	D. discoideum (slime mold)	9615	844
HPA	human protein subcellular localization	9963	0
SynGO (Koopmans et al. 2019)	neuron-neuron synapses	9552	0
PINC	human and mouse	6746	0
MTBBASE	M. tuberculosis (bacterial pathogen)	6160	463
IntAct* (Del Toro et al. 2022)	protein-protein interactions	4849	216488
CAFA (Radivojac et al. 2013)	various	4818	371
CACAO* (Ramsey et al. 2021)	various	4382	0
AspGD (Cerqueira et al. 2014)	A. niger (fungal pathogen)	4099	0
PseudoCAP (Winsor et al. 2005)	P. aeruginosa (bacterium)	2323	0
EcoliWiki* (McIntosh et al. 2012)	E. coli (bacterium)	2123	55
TIGR	bacteria	2150	0
GO_Central*	various	3643	160

CollectTF	bacterial transcription factors	1850	0
NTNU_SB	human, mouse, rat transcription factors	1733	0
GR	rice	1260	0
SGN	tomato	1255	0
DisProt	disordered proteins	933	156
Xenbase* (Fortriede et al. 2020)	Xenopus (frog)	731	0

1
2
3
4
5

Table 3. Genome coverage of PAN-GO annotations. Percentage of protein-coding genes with at least one PAN-GO reviewed annotation, for different taxonomic groups

Taxonomic group	Number of individually reviewed, annotated genomes	Gene coverage of annotations
vertebrates	19	66% - 83%
invertebrates	15	40% - 68%
fungi	14	32% - 76%
plants	40	28% - 51%
protists, alveolates, amoebae	11	19% - 46%
archaea	8	23% - 34%
bacteria	35	20% - 57%

6
7
8

Table 4. Number of PAN-GO annotations for selected genomes.

Genome	Total IBA annotations	MF annotations	BP annotations	CC annotations
Danio rerio	82855	23049	33888	25918
Mus musculus	72554	19832	29947	22775
Rattus norvegicus	71276	19792	28738	22746
Homo sapiens	68695	18537	28177	21981
Gallus gallus	59293	15847	24010	19436
Xenopus tropicalis	43232	12687	16735	13810
Arabidopsis thaliana	37509	12106	12922	12481
Caenorhabditis elegans	31093	9021	11728	10344
Drosophila melanogaster	30331	8628	11202	10501
Dictyostelium discoideum	18966	5556	6815	6595
Saccharomyces cerevisiae	15675	4286	5763	5626

Schizosaccharomyces pombe	13723	3719	4953	5051
Escherichia coli	5681	2171	1916	1594

1
2
3
4
5
6
7

Table 5. GO evidence codes+reference combinations. Users should consider both the type of evidence and the review level. GO internal references (starting with GO_REF:) describe specific annotation methods, and are available at <https://github.com/geneontology/go-site/tree/master/metadata/gorefs/README.md>.

Evidence code	Reference	Evidence type description	Review level description
IDA Inferred from direct assay	Scientific publication providing experimental data	Experimental, most direct evidence for function	Individually expert-reviewed
IMP Inferred from mutant phenotype	Scientific publication providing experimental data	Experimental, from a perturbation in the normal function	Individually expert-reviewed
IGI Inferred from genetic interaction	Scientific publication providing experimental data	Experimental, from perturbations in normal functions of more than one gene	Individually expert-reviewed
IEP Inferred from expression pattern	Scientific publication providing experimental data	Experimental, used only for biological process annotations, from comparison with genes of known function	Individually expert-reviewed
IPI Inferred from protein interaction	Scientific publication providing experimental data	Experimental, used only for annotations to protein binding terms	Individually expert-reviewed
HDA Inferred from high throughput direct assay	Scientific publication providing experimental data	Experimental (high throughput), direct	Not individually reviewed; expert review methodology to exclude high false positive rate observations
HMP Inferred from high throughput mutant phenotype	Scientific publication providing experimental data	Experimental (high throughput), mutant phenotype	Not individually reviewed; expert review to exclude high false positive rate observations
HGI Inferred from high throughput genetic interaction	Scientific publication providing experimental data	Experimental (high throughput), genetic interaction	Not individually reviewed; expert review to exclude high false positive rate observations

HEP Inferred from high throughput expression pattern	Scientific publication providing experimental data	Experimental (high throughput), expression pattern	Not individually reviewed; expert review to exclude high false positive rate observations
IBA Inferred from biological ancestor	(Gaudet et al. 2011)	Homology, from experimental evidence propagated through a phylogenetic tree and/or from direct experimental evidence	Individually expert-reviewed in the context of all experimental annotations for related genes
ISS Inferred from sequence similarity	Scientific publication providing sequence similarity evidence	Homology, from experimental evidence propagated from one gene to one related gene, asserted in the publication	Individually expert-reviewed in the context of all experimental annotations for related genes
ISS Inferred from sequence similarity	GO_REF:0000024	Homology, from experimental evidence propagated from one gene to one related gene, asserted by a biocurator	Individually expert-reviewed
ISO Inferred from sequence orthology	Scientific publication providing orthology evidence	Homology, from experimental evidence propagated from one gene to one orthologous gene	Individually expert-reviewed
ISO Inferred from sequence orthology	GO_REF:0000008	Homology, from experimental evidence propagated from one mammalian gene to one orthologous mouse gene	Individually expert-reviewed
ISO Inferred from sequence orthology	GO_REF:0000024	Homology, from experimental evidence propagated from one gene to one orthologous gene	Individually expert-reviewed
ISO Inferred from sequence orthology	GO_REF:0000096	Homology, from experimental evidence propagated from one	Not individually reviewed; orthology manually reviewed

		gene to one orthologous gene among human, mouse, rat orthologs	
ISO Inferred from sequence orthology	GO_REF:0000101	Homology, from experimental evidence propagated from one gene to one orthologous gene	Not individually reviewed; orthology computed using OrthoMCL
IEA Inferred from electronic annotation	GO_REF:0000107	Homology, from experimental evidence propagated from one gene to one orthologous gene	Not individually reviewed; 1-to-1 orthology computed using Ensembl Compara phylogenetic trees
IEA Inferred from electronic annotation	GO_REF:0000002	Homology, from a hit to an InterPro signature	Not individually reviewed; expert review of annotations of signatures to ensure low or no false positives
IEA Inferred from electronic annotation	GO_REF:0000003	Imported from another resource, from mapping an EC number assigned in UniProt	Not individually reviewed; expert review of mappings; EC assignments are manually reviewed for Swiss-Prot and computationally inferred for TrEMBL
IEA Inferred from electronic annotation	GO_REF:0000004	Imported from another resource, from mapping a manually assigned Swiss-Prot keyword	Not individually reviewed; expert review of keywords and mappings
IEA Inferred from electronic annotation	GO_REF:0000104	Homology, from manually curated UniRule	Not individually reviewed; expert curation of UniRules to ensure low or no false positives
IEA Inferred from electronic annotation	GO_REF:0000108	Logical assertion using the ontology, from asserted relation between different aspects of GO	Not individually reviewed; expert curation of ontology links
IEA Inferred from electronic annotation	GO_REF:0000117	Computational, from machine learning	Not individually reviewed; assigned by machine learning from curated training sets

TAS Traceable author statement	Scientific publication citing original data	From a published statement referencing experimental evidence in a different paper	Individually reviewed
NAS Nontraceable author statement	Scientific publication with general biological knowledge statement	From an unreferenced published statement	Individually reviewed

1
2
3
4
5
6
7
8
9

Table 6. GO ontology editions. Editions are distinguished by the relations and metadata they include. All editions are updated at each GO release. External ontologies used in GO include: ChEBI, Uberon (Haendel *et al.* 2014), Relation Ontology (Smith *et al.* 2005), Cell Ontology (Diehl *et al.* 2016), Sequence Ontology (Mungall, Batchelor and Eilbeck 2011), Dicty Anatomy, CARO (Haendel *et al.* 2008), Fungal Anatomy Ontology (*Fungal-Anatomy-Ontology: A Structured Controlled Vocabulary for the Anatomy of Fungi*), Plant Ontology (Walls *et al.* 2019), PATO (Gkoutos, Schofield and Hoehndorf 2018), Protein Ontology (Natale *et al.* 2017).

GO edition	Format(s)	Relations included	Links to other ontologies
go-basic	OBO	is a, part of, regulates, negatively regulates and positively regulates	Not available
go	OBO and OWL-RDF/XML	Same as go-basic, plus has part and occurs in	Not available
go-plus	OWL-RDF/XML		ChEBI, Uberon, Cell Ontology, Sequence Ontology, Dicty Anatomy, CARO, Fungal Anatomy Ontology, Plant Ontology, PATO, Protein Ontology

10

11 Figure legends

12
13
14
15
16
17
18
19

Figure 1. Examples of the three components of the GO knowledgebase. A) The GO ontology consists of terms, e.g. *DNA binding transcription factor activity*, and relationships between the terms (arrows; black=*is a*, blue=*part of*, orange=*regulates*). B) GO annotations associate a specific gene product (here, human ZNF410) with GO terms asserting its functional aspects ("GO Class" column, e.g. sequence-specific double stranded DNA binding) and the evidence for each assertion with its traceable source ("Evidence" and "Reference" columns). C) The GO-CAM model combines individual GO annotations into a model, in this case a very

1 simple model describing how human ZNF410 acts as a transcription factor to positively regulate
2 (denoted by the green arrow) transcription of the *CHD4* gene, which in turn acts as a co-
3 repressor to repress (denoted by red dashed lines) transcription of fetal hemoglobin genes
4 (*HBG1* and *HBG2*) in erythroid lineage cells. In this view, each box in the GO-CAM is labeled
5 with the gene product and species abbreviation for simplicity.
6
7

8 **Figure 2.** GO-CAM model of the SARS-CoV2 - host interactions as displayed using the GO-
9 CAM Pathway Widget (code available at <https://github.com/geneontology/wc-gocam-viz>) on the
10 Alliance of Genome Resources gene pages
11 (<https://www.alliancegenome.org/gene/HGNC:20144#pathways>). The model includes proteins
12 from both humans (Hsap) and the SARS-CoV-2 virus (Scov2). A simplified representation of the
13 causal model is shown on the main figure, which is simplified by labeling with the gene and
14 organism. The model includes many additional details, which are displayed as “cards”; the
15 information for MAVS activity (inset) which normally acts as a signaling adaptor located in the
16 mitochondrial membrane. MAVS activity is suppressed directly by the SARS-CoV-2 M protein,
17 and indirectly by other SARS-CoV-2 proteins. Each of the “E” symbols on the right-hand side
18 can be clicked to see the evidence for each assertion in the model.
19
20

21 **Figure 3.** Alliance ribbon view for the yeast RPB7 gene. High-level GO categories annotated
22 are shown in blue squares (<https://www.alliancegenome.org/gene/SGD:S000002812>).
23
24

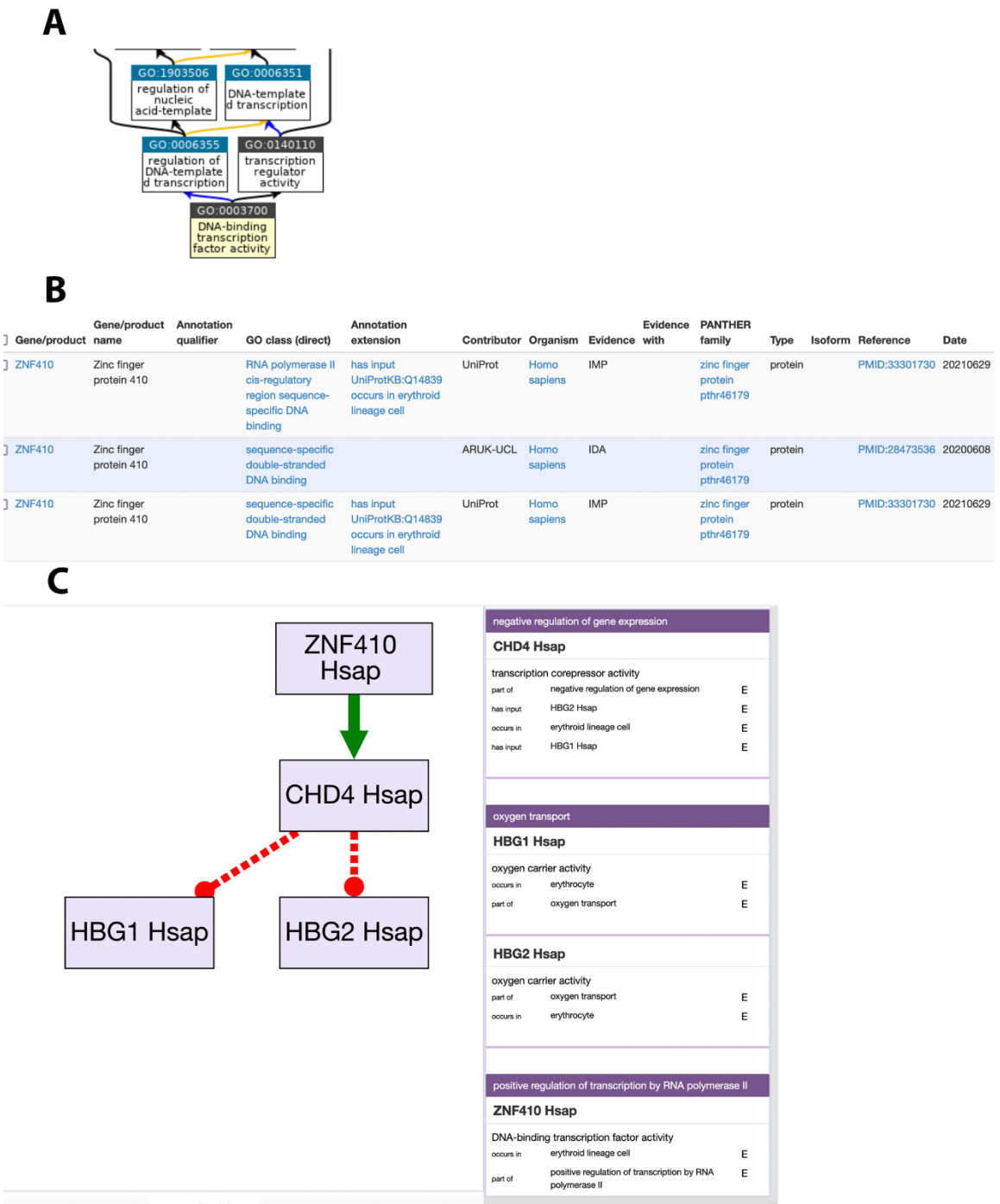


Figure 1
152x188 mm (x DPI)

1
2
3
4

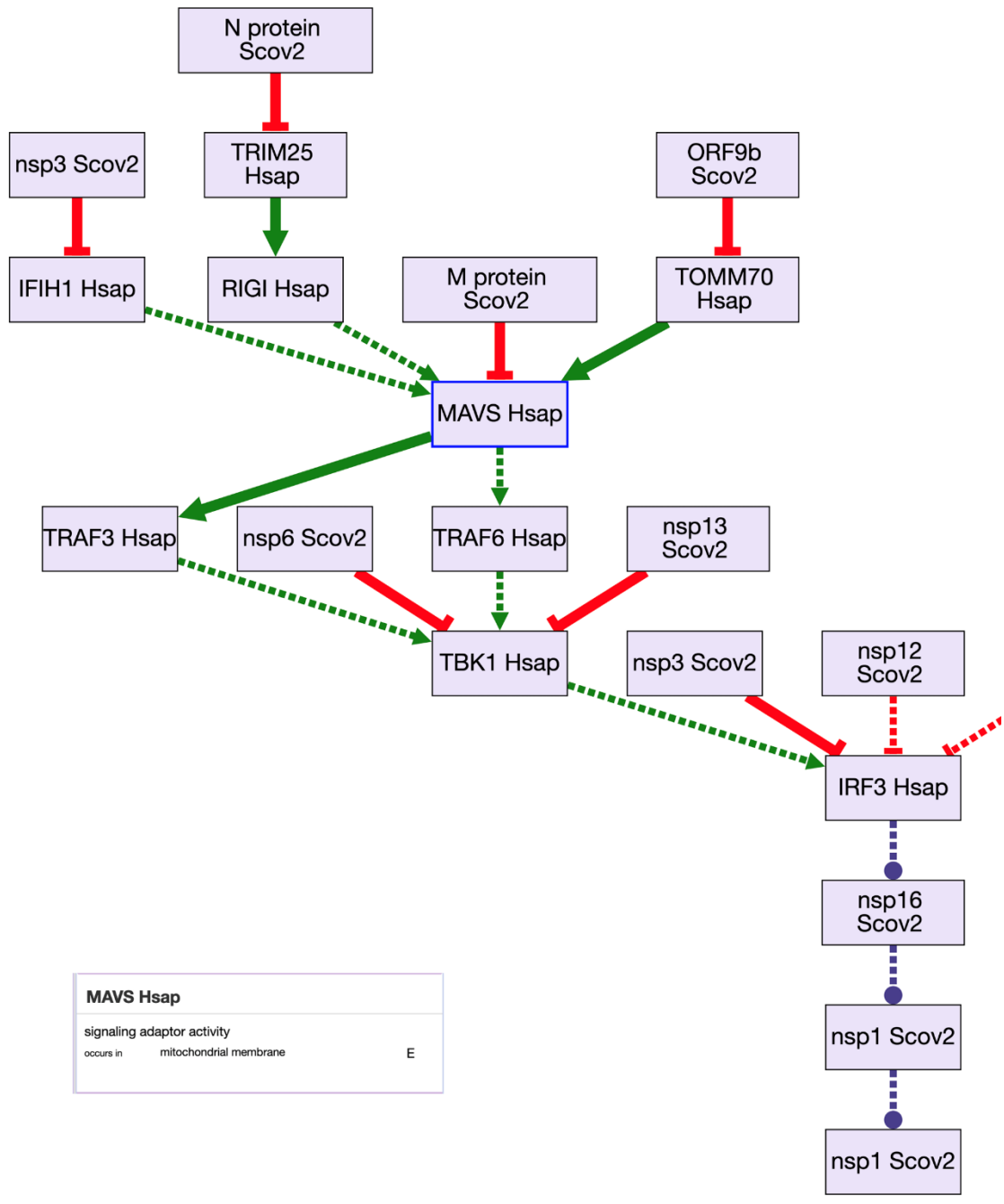
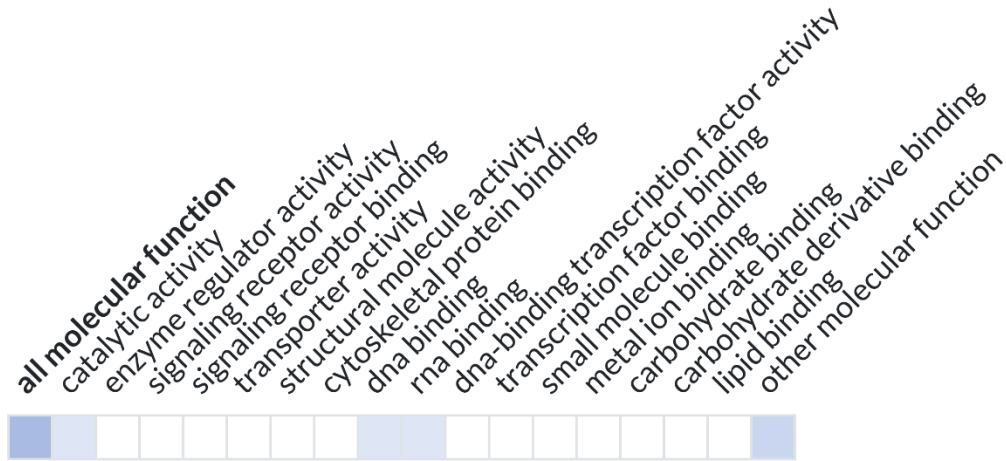


Figure 2
152x178 mm (x DPI)

1
2
3
4



1
2
3

Figure 3
138x67 mm (x DPI)