

RESEARCH ARTICLE

# Overcoming the widespread flaws in the annotation of vertebrate selenoprotein genes in public databases

Max Ticó<sup>1,2</sup>, Emerson Sullivan<sup>1,3</sup>, Roderic Guigó<sup>2</sup>, Marco Mariotti<sup>1\*</sup>

**1** Department of Genetics, Microbiology and Statistics, Universitat de Barcelona, Barcelona, Catalonia, Spain, **2** Centre for Genomic Regulation (CRG), The Barcelona Institute for Science and Technology (BIST), Barcelona, Catalonia, Spain, **3** Johns Hopkins Whiting School of Engineering, Baltimore, Maryland, United States of America

\* [marco.mariotti@ub.edu](mailto:marco.mariotti@ub.edu)



## Abstract

Genome annotations provide the essential framework for genomic analyses, capturing our current knowledge of gene structure and function as inferred from computational predictions and experimental evidence. Even as automated annotation pipelines become more sophisticated, their accuracy in representing unconventional gene expression events remains largely untested. Here, we address this gap by examining the most common form of translational recoding: the insertion of selenocysteine (Sec), a non-canonical amino acid incorporated into selenoproteins, oxidoreductase enzymes carrying essential roles in redox homeostasis. Sec insertion occurs in response to UGA, normally interpreted as stop codon, but recoded in selenoprotein mRNAs. Owing to the dual function of UGA, the identification of selenoprotein genes poses a challenge. We show that the vertebrate selenoprotein genes are widely misannotated in major public databases. Only 11% and 5% of selenoprotein genes are well annotated in Ensembl and NCBI GenBank, respectively, due to the lack of dedicated selenoprotein annotation pipelines. In most cases (81% and 84%), overlapping flawed annotations are present which lack the Sec-encoding UGA. In contrast, NCBI RefSeq employs a dedicated selenoprotein pipeline, yet with some shortcomings: its selenoprotein annotations are correct in 77% of cases, and most errors affect families with a C-terminal Sec residue. We argue that selenoproteins must be correctly annotated in public databases and that must occur via automated pipelines, to keep the pace with genome sequencing. To facilitate this task, we present a new version of Selenoprofiles, an homology based tool for selenoprotein prediction that produces predictions with accuracy comparable to manual curation, and can be easily deployed and integrated in existing annotation pipelines.

## OPEN ACCESS

**Citation:** Ticó M, Sullivan E, Guigó R, Mariotti M (2026) Overcoming the widespread flaws in the annotation of vertebrate selenoprotein genes in public databases. *PLoS Comput Biol* 22(1): e1013885. <https://doi.org/10.1371/journal.pcbi.1013885>

**Editor:** Iddo Friedberg, Iowa State University College of Veterinary Medicine, UNITED STATES OF AMERICA

**Received:** August 28, 2025

**Accepted:** December 31, 2025

**Published:** January 12, 2026

**Copyright:** © 2026 Ticó et al. 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 use, distribution, and reproduction in any medium, provided the original author and source are credited.

**Data availability statement:** The latest Selenoprofiles code is available at <https://github.com/marco-mariotti/selenoprofiles4> and documented at <https://selenoprofiles4.readthedocs.io/>. Profile alignments, anchor alignments, and expectation tables of selenoproteins per

lineage are provided as companion package [https://github.com/marco-mariotti/seleno-protein\\_profiles](https://github.com/marco-mariotti/seleno-protein_profiles). Selenoprofiles predictions on Ensembl and NCBI genomes, and related annotation status labels, can be found at Zenodo Doi: 10.5281/zenodo.17828122, and also online at [https://mariottigenomicslab.bio.ub.edu/selenoprofiles\\_data/selenoprofiles\\_predictions/](https://mariottigenomicslab.bio.ub.edu/selenoprofiles_data/selenoprofiles_predictions/).

**Funding:** This work was supported by grants RYC2019-027746-I funded by MICIU/AEI (Ministry of Science, Innovation and Universities; <https://www.ciencia.gob.es/en/>; Agencia Estatal de Investigación; <https://www.aei.gob.es/en/>) /10.13039/501100011033 and “ESF Investing in your future” (European Social Fund; <https://european-social-fund-plus.ec.europa.eu/en/>); grant PID2020-115122GA-I00 funded by MICIU/AEI, and PID2023-147164NB-I00 funded by MICIU/AEI and ERDF/EU (European Regional Development Fund / European Union; [https://ec.europa.eu/regional\\_policy/funding/erdf\\_en/](https://ec.europa.eu/regional_policy/funding/erdf_en/)), to MM. Funders did not play any role in the study design, data collection and analysis, decision to publish, or preparation of the manuscript.

**Competing interests:** The authors have declared that no competing interests exist.

## Author summary

Genome annotations serve as a map of the genetic information encoded in our DNA. They are essential for numerous applications, working as a reference to interpret high throughput experiments. While the annotation of the human genome is manually curated, the pace of sequencing far surpassed our collective capacity for manual annotation, so that we mostly rely on our automated methods for the majority of available genomes. While these methods may work correctly for genes that follow the standard “rules of life”, they may overlook the many exceptions found in nature. Once such exception is translational recoding, wherein the canonical rules of protein synthesis (the “genetic code”) are bent for specific genes. Here, we analysed the most prevalent form of translational recoding: selenoproteins. These are proteins which contain the unusual amino acid selenocysteine, which is inserted into proteins using a genetic code signal that usually means “stop”. We show that the vast majority of selenoproteins in vertebrate genomes are misrepresented in public databases, meaning that essential genes are either missing or incomplete in widely used resources. To address this issue, we developed an improved version of *Selenoprofiles*, a software tool that automatically identifies and annotates selenoprotein genes with high accuracy.

## Introduction

Selenoproteins are a group of proteins that contain the amino acid selenocysteine (Sec or U), the twenty-first amino acid in the genetic code. Sec is co-translationally inserted via recoding of UGA, a codon which normally terminates protein synthesis [1,2]. Sec insertion primarily depends on a *cis*-acting RNA structure called SECIS element, located in 3' UTR regions of selenoprotein mRNAs of eukaryotes [3]. The human selenoproteome consist of 25 selenoprotein genes, which have important roles in various biological processes including redox homeostasis, hormone maturation, immune response, and others [2,4,5]. Thioredoxin reductases (TXNRD), glutathione peroxidases (GPX) and iodothyronine deiodinases (DIO) are represented by 3, 5, and 3 selenoprotein genes in human, respectively, and they are arguably the most well studied eukaryotic selenoprotein families. TXNRDs and GPXs are essential components of the two major antioxidant systems in mammals, the thioredoxin and glutathione system respectively, while DIOs control thyroid hormone maturation. The functions of approximately a third of human selenoproteins are yet uncharacterized or poorly defined. Most selenoprotein families include homologs that substitute Sec with cysteine (Cys), i.e., “Cys homologs”. These may be found as orthologs (e.g., *Drosophila* thioredoxin reductase *Trxr1*) and/or as paralogs (e.g., human *GPX5*, *GPX7*, *GPX8*).

Selenoprotein genes have been identified in diverse branches of the tree of life. In this work, we focus on eukaryotes only. Vertebrate have relatively well-characterized selenoproteomes encompassing 19 protein families, some represented by several

homologs per genome [5]. Tetrapods typically encode 24–25 selenoproteins, whereas fish show a broader variation, with counts ranging from 28 to 55 [6,7]. Sec usage was also observed in other eukaryote species, including most metazoa, algae, some protists, and a handful of fungi [8–10]. On the other hand, Sec is absent in higher plants, yeasts, many insects and few nematodes [8,10–12].

Owing to the rare dual function of the UGA codon, selenoproteins are usually neglected by standard gene annotation programs, so that specialized software is required [8,13]. Relevantly for this work, our group previously developed Selenoprofiles, an automated pipeline to predict selenoprotein genes by homology to a set of built-in profiles representing known selenoprotein families [14,15].

In the present work, we assessed selenoprotein gene annotation in two major public databases, Ensembl and NCBI. We focused on vertebrates, arguably the best annotated of all genomes. By using Selenoprofiles to predict selenoproteins in genome assemblies and assessing the corresponding gene annotations, we show evidence for widespread misannotation of selenoprotein genes. To remedy this important issue, we have upgraded Selenoprofiles with several key enhancements: a streamlined, one-command installer; an automatic orthology-assignment utility that labels each prediction with gene names; and a lineage-aware filter that retains only candidates consistent with the expected selenoproteome. Together, these features yield results that match the precision of manual curation, and facilitate the integration of Selenoprofiles in automated genome annotation pipelines.

## Materials and methods

### Genomic sequences and resources

We downloaded all assemblies and annotations of vertebrate organisms available in Ensembl version v.108 [16], comprising of 309 genomes, each corresponding to a different species or strain. For the data at the National Center for Biotechnology Information (NCBI), we used the tool “NCBI datasets” to obtain all genome assemblies corresponding to vertebrates, then considered only those with an available gene annotation [17,18]. Because some annotations were available but limited to the mitochondrial genome or scaffold structure, we only considered GTF files whose archive size was >3MB. This resulted in a set of 1,172 NCBI vertebrate genomes, each representing a distinct species. Selenoprofiles (see below) was successfully run on all genomes, except for 25 NCBI and 3 Ensembl assemblies. These assemblies consistently produced non-recoverable errors and were consequently excluded from our dataset. The list of genome assemblies is available in [S1 Table](#). For analysis, we differentiated between RefSeq (ids starting with “GCF\_”) from GenBank (“GCA\_”) entries. The phylogenetic tree of Ensembl species we downloaded from Ensembl Compara [19]. The analogous NCBI tree was obtained from NCBI Taxonomy [20].

### Selenoprotein gene prediction

Selenoprotein gene prediction was performed using Selenoprofiles [14,15] version v4.5.7, available at <https://github.com/marco-mariotti/selenoprofiles4>. Selenoprofiles is a computational pipeline to identify members of the known selenoprotein families and related proteins. It uses several homology-based gene prediction programs, and employs a set of manually curated multiple sequence alignments of selenoprotein families to scan genomes for homologues. The profiles we used are available at [https://github.com/marco-mariotti/selenoprotein\\_profiles](https://github.com/marco-mariotti/selenoprotein_profiles) (v1.4.1). We scanned genome assemblies for all known selenoprotein families in metazoa and focused our analyses on the predictions labelled as “selenocysteine”, i.e., the genes encoding for selenoproteins, with no apparent pseudogene features (e.g., frameshifts). Our results in vertebrate genomes originated from 21 profiles representing 19 protein families. Multimember families are represented by a single profile that can correctly predict all its paralogs (e.g., the GPX profile can predict Sec-encoding *GPX1*, *GPX2*, *GPX3*, *GPX4*, *GPX6* and Cys-paralogs *GPX5*, *GPX7*, *GPX8*). *SELENOW* and *SELENOV* proteins share high homology and are therefore considered as groups within the same protein family, represented by a single profile (SELENOW). On the

other hand, though we consider *SELENOF* and *SELENOE* proteins to be part of the same homologous family, they are dissimilar enough to warrant two different profiles, for better prediction performance. Lastly, *SELENOK* genes are highly divergent in insects, which led to the development of two distinct profiles (*SELENOK* and *SelKi*) whose results are joined in post-processing. Whenever a profile encompasses multiple groups or paralogs, the *orthology* routine (described below) allows to classify predicted genes in the relevant subfamilies.

### Annotation assessment

To assess the annotation of selenoproteins in public databases, we developed a new utility called *Selenoprofiles assess*, now integrated in *Selenoprofiles* v4.5.7. This tool takes as input two GTF/GFF files, one from *Selenoprofiles* and one from another source, and matches gene annotations based on their overlap in genomic coordinates. Overlapping annotations must occur on the same strand, and on the same frame (except for the “Out of frame” classification). After analysis of the Sec-encoding codon and the rest of coding sequence, each *Selenoprofiles* prediction is assigned an annotation status label, as shown in Fig 1 and described in Results. To develop this utility, we took advantage of *PyRanges*, a package to represent and manipulate genomic data in Python [21]. Plots were built using the R package *ggplot2* [22].

### Phylogenetic reconstruction and orthology

The protein sequences of all *Selenoprofiles* predictions across all genomes were collected using the *Selenoprofiles join* utility. These were used as input to phylogenetic reconstruction workflows, available in the ETE3 v3.1.3 build framework [23]. We ultimately analyzed trees obtained with the “*eggnoG41*” workflow [24], which internally uses *PhyML* v20160115. patched for phylogenetic reconstruction [25]. We visualized gene trees with the *show\_tree* program in the *biotree\_tools* package v0.0.6 ([https://github.com/marco-mariotti/biotree\\_tools.git](https://github.com/marco-mariotti/biotree_tools.git)), based on the ETE3 python module [23].



**Fig 1. Classification of selenoprotein annotation/misannotations.** Segments represent coding sequence (CDS) regions of selenoprotein genes. *Selenoprofiles* predicted CDSs are colored in dark grey, and annotated Ensembl/NCBI CDSs are in light grey. Dotted vertical lines represent start and end positions of the Sec encoding codon. Introns were shrunk to a fixed size for visualization purposes. The classes in blue, i.e., all other than “Well annotated” and “Absent”, are collectively referred to with the term “Misannotation” in this manuscript.

<https://doi.org/10.1371/journal.pcbi.1013885.g001>

For every selenoprotein family with multiple paralogs in vertebrates (GPX, TXNRD, DIO, SELENOW/V), we inspected phylogenetic trees and manually partitioned them into orthologous groups, here referred to as subfamilies, by applying the species overlap method and following maximum parsimony criteria [26]. To assign subfamily names, we performed similarity searches against Uniprot and consulted relevant literature [6,27,28].

Next, we built a new utility called *Selenoprofiles orthology*, now available in v4.5.7, constituting a fast alternative to tree building for subfamily assignment. To this aim, we first generated reference alignments (“anchors”) for each multimember selenoprotein family including representatives of each subfamily, which were assigned to orthologous groups beforehand via phylogenetic reconstruction as detailed above. Alignments were performed using Mafft v7.45331 [29]. Anchor alignments were reduced to a maximum of 4, 8, or 12 representative sequences for each subfamily for each lineage, automatically selected using TrimAl [30] to preserve sequence diversity. In the *Selenoprofiles orthology* procedure, candidate sequences are aligned to reference alignments using Mafft, and Pyaln v0.1.4 (<https://github.com/marco-mariotti/pyaln.git>) is used for fast calculation of similarity scores of candidate sequences against the representative sequences of each subfamily. These scores are used for subfamily assignment of each candidate. Similarity score parameters are displayed in S2 Table. Results of our benchmark with 3,952 Selenoprofiles predictions in Ensembl genomes are shown in S3 Table, which led to choosing anchor alignments with size of 12 sequences, and this Pyaln sequence comparison method: gap positions are excluded (gaps=n), and the average weighted sequence identity “AWSI” (metrics=‘w’) is computed with weights based on the maximum frequency of non-gap characters (weights=‘m’). The same similarity metrics is used for ranking gene predictions in the novel filtering procedure *Selenoprofiles lineage* (see Results). Sequence and tree manipulations throughout the project were carried out using the programs *alignment\_tools* and *tree\_tools* from the *biotree\_tools* package.

## Results

### Evaluation of selenoprotein annotation

We set to evaluate the status of selenoprotein annotations in Ensembl and NCBI. The identification of all selenoprotein genes in a genome is not trivial, and it is likely that many selenoprotein families remain undiscovered in less studied organisms, e.g., unicellular eukaryotes [8]. To elude this problem, we focused on vertebrates: their selenoproteome has been thoroughly studied [4,6,7,31–36], so that, we argue, it can be reliably identified by homology to known selenoprotein families. Our dataset consisted of paired genome assemblies and genome annotations, 306 from Ensembl and 1,147 from NCBI (Methods). We subdivided NCBI genomes into two categories: GenBank (413 entries), which includes genome and annotation data submitted directly by researchers without further curation, and RefSeq (734 entries), which comprises a non-redundant selection of genomes consistently annotated using the NCBI Eukaryotic Genome Annotation Pipeline (EGAP) [37]. As *bona fide* selenoprotein annotation, we ran Selenoprofiles, which yielded 8,827 and 31,572 “raw” selenoprotein predictions in Ensembl and NCBI genomes, respectively, corresponding to 19 protein families (see Methods, Data Availability). Our choice is justified by the high accuracy of Selenoprofiles, whose benchmark demonstrated sensitivity and specificity above 90% at the nucleotide level when compared to manually annotated gene structures [15]. Later on, we critically evaluated and refined this method by adding a novel *Selenoprofiles lineage* utility, detailed in later sections of this manuscript, which yielded a higher-quality set of “filtered” predictions (8,085 and 29,242 in Ensembl and NCBI, respectively).

We inferred the status of selenoprotein annotation by matching genomic coordinates from Selenoprofiles predictions with annotated coding sequence (CDS) regions and evaluating whether the Sec codon was included in the annotated CDS. If it was not the case, we further classified them into different misannotation cases depending on the position of Sec and annotated CDS (Fig 1). CDS annotations that ended precisely at the Sec UGA were classified as “Stop codon”. The misannotation classes “Upstream”, “Downstream”, and “Skipped” were applied when there were portions of annotated CDS at the 5’, 3’, or both 5’ and 3’ of the Sec UGA, respectively. Some selenoprotein annotations exhibited a distinct

frame compared to Selenoprofiles predictions, and were classified as “Out of frame” cases. Transcripts which contained a splicing site in the Sec codon were classified as “Spliced”. Instances where no annotation was overlapping to a Selenoprofiles selenoprotein CDS were classified as “Absent”.

The great majority of selenoprotein families contain a single Sec residue. Notable exceptions in vertebrates are SELENOP, which contains 10 Secs in human and up to 37 in fish [38], and SELENOL, which contains two nearby Sec forming a diselenide bond [36]. While the first Sec residue of SelenoP is embedded in a conserved thioredoxin-like motif, the rest are located in a Sec-rich C-terminal tail, which is particularly challenging to identify in genomes [38]. For simplicity, here we considered only the first Sec residue of each selenoprotein gene to define its annotation status.

### The status of selenoprotein annotation in Ensembl

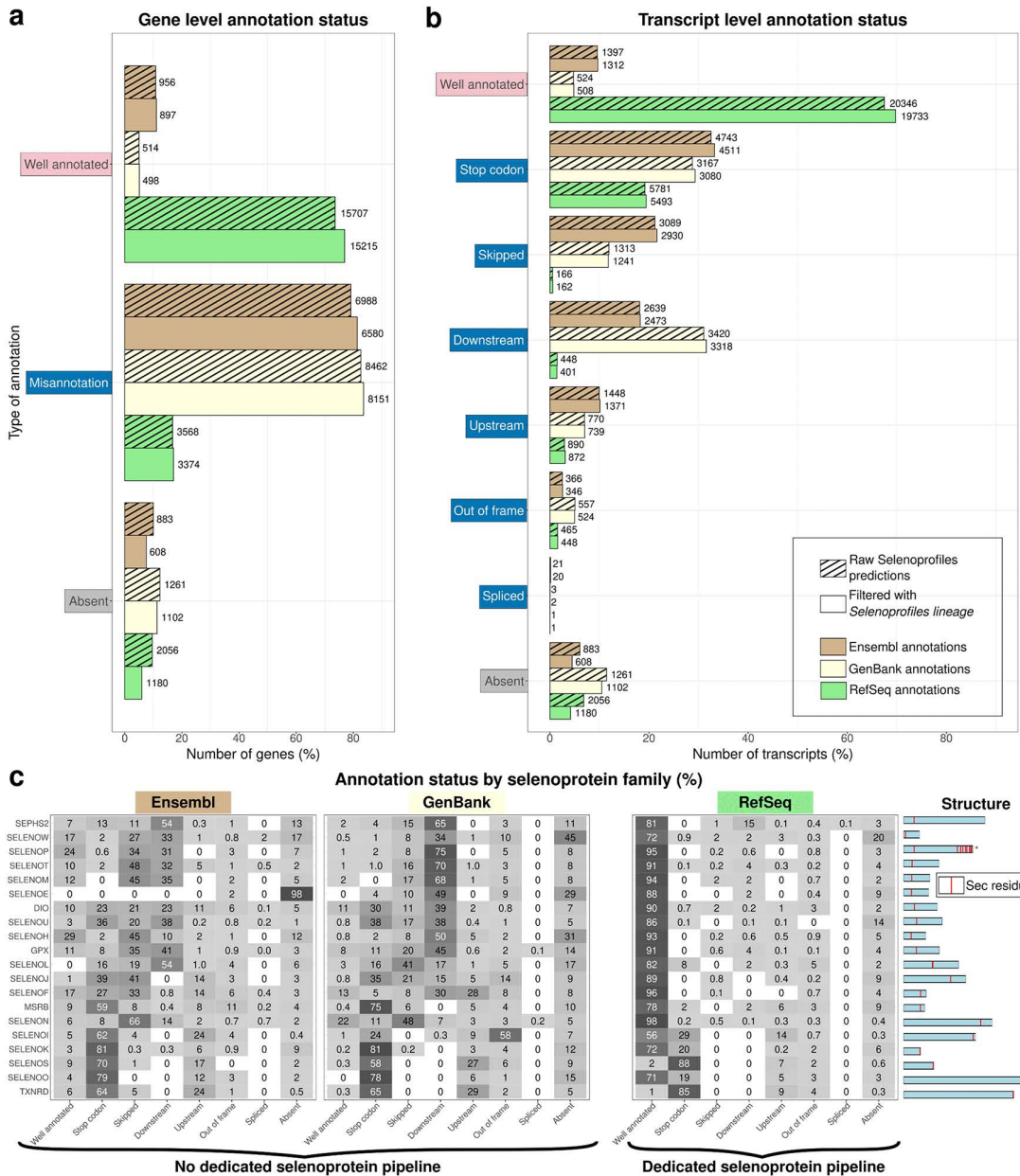
Our results show that most selenoprotein genes have defective annotations in Ensembl: only ~11% of predicted selenoproteins present at least one “Well annotated” entry (Fig 2A). Around 8% of selenoproteins have no annotation at all (“Absent”). The majority (~81%) are “Misannotations”, i.e., there are only flawed annotations which lack the Sec-encoding UGA. We next examined annotation at transcript level, wherein a single Selenoprofiles gene prediction may match multiple annotated transcripts, particularly in well studied model organisms. This analysis allowed us to distinguish among misannotation types (Fig 1). We observed that almost 22% of transcripts skip the Sec-UGA-containing exon, or at least its relevant portion (Fig 2B). There are also numerous cases where the annotated CDS ends at the Sec codon (33%), or further upstream (10%), or starts downstream of it (18%). Other cases were more rare.

Next, we examined the annotation status for each selenoprotein family separately. Surprisingly, the most well-known selenoprotein families do not exhibit the best annotation status in Ensembl (Fig 2C); rather, SELENOH and SELENOP (first Sec) exhibit the highest percentage of “Well annotated” cases (29% and 24%, respectively). Consistent with intuition, we observed a clear correlation between misannotation type and the position of Sec in the gene structure. In protein families that carry Sec in close proximity to the C-terminus, such as TXNRDs and selenoproteins S, O, K, and I, the Sec residue tends to be annotated as a “Stop codon” signal, or the CDS ends further “Upstream” (Fig 2C). On the other hand, “Downstream” cases are more abundant in families containing Sec in their N-terminal part, such as GPX, DIO, SEPHS2, and selenoproteins H, M, N, T, V, W. Unsurprisingly, the selenoproteins SELENOE [32], SELENOL [36], and SELENOJ [31], that are not present in human and mouse, exhibit the worst status in Ensembl, with <1% “Well annotated” genes.

We noticed that misannotations are not homogeneously distributed across species. Instead, “Well annotated” selenoprotein genes are concentrated in few model species (e.g., human, mouse, zebrafish, pig). In contrast, almost all selenoproteins are misannotated in the rest of species. The full tree of Ensembl species with the complete classification of their selenoprotein genes is shown in S1 Fig, while Fig 3A shows a reduced version with a representative set of species. Altogether, we attribute the presence of well annotated selenoprotein genes in model organisms to the work of manual curators, while the pervasive mistakes in remaining taxa stem from the lack of a dedicated automated pipeline at Ensembl. In some cases, this led to the paradoxical situation of closely related species exhibiting selenoproteome annotations of very different quality.

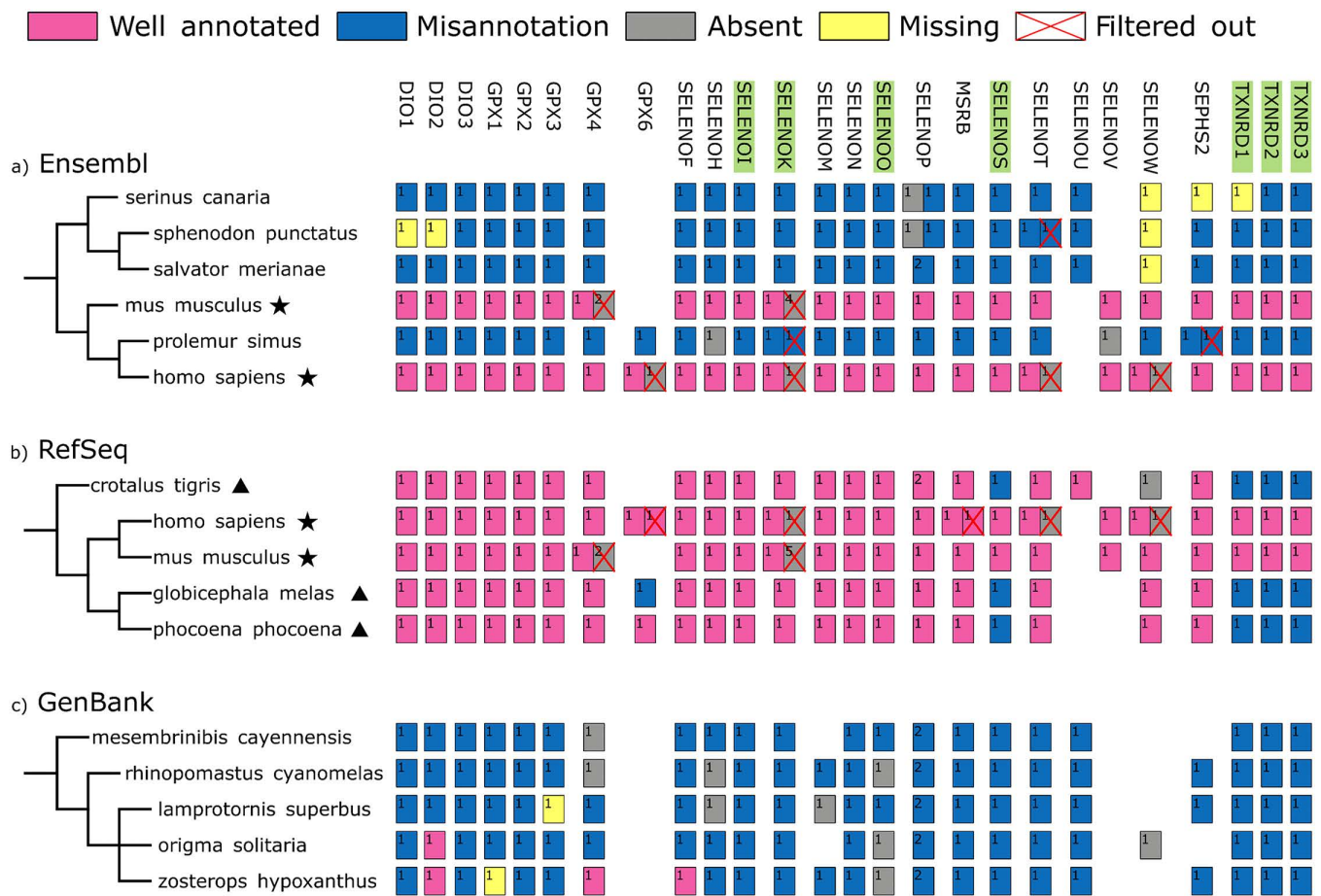
### The status of selenoprotein annotation in NCBI GenBank and RefSeq

Our analysis reveals that selenoprotein annotation is flawed in NCBI GenBank, too: only ~5% of selenoprotein genes are “Well annotated”, while 84% fall into the “Misannotation” category and 11% are completely “Absent”. At the transcript level, the most frequent error remains treating the Sec residue as a stop codon, but the other classes are also common. Well annotated cases are spread across diverse species with no apparent underlying pattern. Correct models appear sporadically across a wide taxonomic range with no discernible trend, reflecting the fact that GenBank records deposited by the research community are generated through diverse pipelines that do not incorporate dedicated selenoprotein routines.



**Fig 2. Status of selenoprotein annotation in Ensembl, RefSeq and GenBank.** **a.** Bar plot showing the annotation status of selenoproteins (Selenoprofiles predictions) at gene level in different databases. Three main annotation cases are shown, with “Misannotation” aggregating all classes (shown in Fig 1) other than “Well annotated” and “Absent”. Bar shading differentiates the raw (i.e., unfiltered) Selenoprofiles predictions and those filtered with *Selenoprofiles lineage* as explained later in the manuscript. Bar height represents percentages, and numbers on top of bars show actual counts. Bar colors correspond to each database. **b.** Bar plot showing the annotation status at transcript level in Ensembl, RefSeq and GenBank, with all possible annotation types. **c.** Heatmap representing the distribution of annotation/misannotation classes for each selenoprotein family, calculated at transcript level using filtered predictions. The right panel shows the gene structures of representatives of each family: the CDS is represented in blue, and Sec residues are displayed as red lines. Families are sorted by the position of Sec residue, with those carrying it at the C-terminus placed at the bottom. Note that the SELENOW family contains both SELENOW and SELENOV genes, since they are highly homologous and included in the same Selenoprofiles profile. (\*) SELENOV: Annotation status was determined using only the first Sec residue of each selenoprotein.

<https://doi.org/10.1371/journal.pcbi.1013885.g002>



**Fig 3. Distribution of misannotations across species.** The species tree of six representative organisms is shown with the annotation of selenoprotein genes in Ensembl (a), RefSeq (b) and GenBank (c), color-coded by annotation/misannotation labels grouped by subfamily. Gene-level annotation classes are used (Fig 1), wherein “Absent” (grey) means that a selenoprotein predicted by Selenoprofiles has no overlapping annotation in Ensembl/NCBI. Moreover, the yellow color is used to mark “Missing” genes, i.e., expected in this genome based on lineage but not detected by Selenoprofiles. Also, red crosses mark those genes “Filtered out” by *Selenoprofiles* lineage, i.e., unexpected genes. Subfamilies with C-terminal Sec are highlighted in green. In the species tree, stars denote model organisms with well annotated selenoproteomes, and a triangle denote NCBI genomes with well annotated selenoproteins except for the C-terminal Sec families.

<https://doi.org/10.1371/journal.pcbi.1013885.g003>

In contrast, vertebrate selenoprotein genes in RefSeq genome annotations have “Well annotated” Sec residues in 77% of cases. The rest of cases are either “Absent” (6%) or “Misannotations” (17%). By analyzing the distribution of misannotation types at transcript-level and by family (Fig 2), it appears evident that the main flaw of RefSeq annotation of selenoproteins concerns the protein families with C-terminal Sec residues. Indeed, the “Stop codon” class constitutes 75% of “Misannotation” cases. It also appears that, in RefSeq, the annotation status of fish-specific selenoprotein families is not worse than mammalian families. Analyzing of the distribution of misannotations across species (representative tree in Fig 3; full set in S2 Fig), we notice that we can divide RefSeq genome annotations in two groups: 1. A handful of model species have selenoproteomes that are completely well annotated; 2. Some species have well annotated selenoproteins, except for a few C-terminal Sec protein families, namely SELENOS, TXNRD, and, in a lesser extent, SELENOI, SELENOK and SELENOO. Taken together, our analysis shows that RefSeq relies on manual curated entries for a few model species, while it deploys an automated pipeline for the rest. The RefSeq pipeline evidently has a procedure in-place for selenoproteins that is overall effective, but falters when the Sec residue is followed by a very short C-terminal extension.

To better illustrate the status of selenoprotein annotation in vertebrates, we summarized our results in 18 organisms commonly used in research (Table 1).

### Improving selenoprofiles for fully automated reliable annotation of selenoproteins

We set to provide an fully automated computational tool to produce reliable selenoprotein annotations. Previously, we developed Selenoprofiles, a pipeline for selenoprotein gene prediction in genomes featuring remarkable gene prediction accuracy [15]. Yet, Selenoprofiles was not integrated in the annotation pipeline of public databases, to our knowledge. We ascribe this at least partially to the characteristics of the early versions of Selenoprofiles, which were hard to install and run. Here, we present a new version which remedied these shortcomings and introduce various other improvements. First, Selenoprofiles can now be installed via a single conda command, or by pulling a ready-made Docker image, as explained in the new comprehensive online documentation (<https://selenoprofiles4.readthedocs.io/>). Moreover, code consistency and reproducibility are ensured through a robust continuous integration system implemented on GitHub, which automatically runs validation tests at every code change.

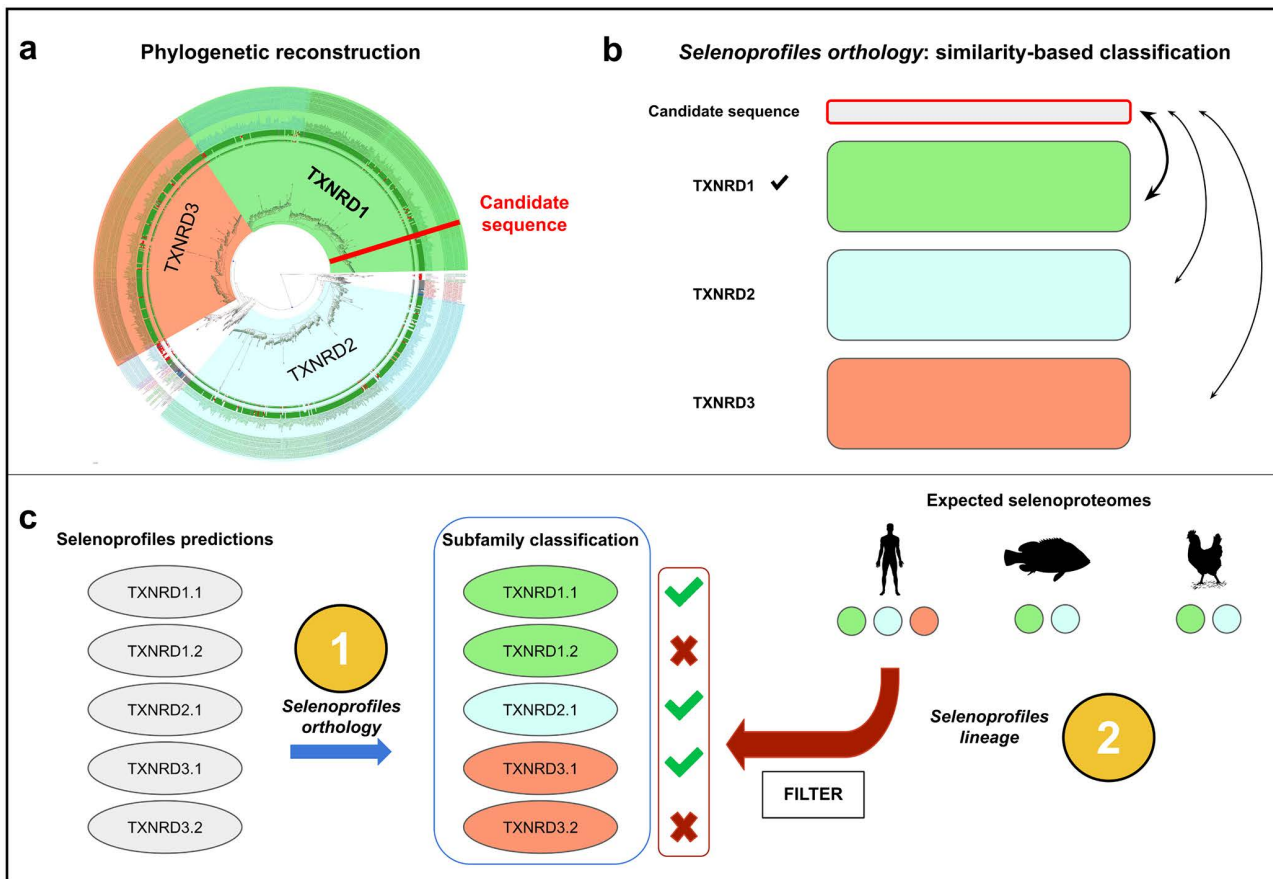
### Selenoprofiles orthology: Automatic assignment of selenoprotein subfamilies

Second, we created a new utility called *Selenoprofiles orthology* to assign gene predictions to specific subfamilies, i.e., orthologous groups/ paralogs. Indeed, Selenoprofiles predicts all paralogous genes of a given family (e.g., DIO1, DIO2, DIO3) using the same profile (DIO), so that, in previous versions, the user could not differentiate among paralogs without a dedicated posterior analysis. The *Selenoprofiles orthology* utility addresses this task by adopting a fast alternative to phylogenetic tree reconstruction, the state-of-the-art for orthology assignment (Fig 4). Our procedure is based on curated

**Table 1. Selenoprotein annotation status in common vertebrate models in Ensembl and NCBI RefSeq. The table reports selenoprotein gene counts for each species in each of these categories: W=well annotated; Ma=misannotated; A=absent from annotation; Mi=missing from assembly. Expected=the total number of selenoproteins expected in that genome.**

Species	Ensembl				NCBI Refseq				Expected
	W	Ma	A	Mi	W	Ma	A	Mi	
Bos taurus	12	8	5	0	24	0	0	1	25
Canis lupus familiaris	1	21	2	1	24	1	0	0	25
Capra hircus	0	24	1	0	17	6	1	1	25
Coturnix japonica	1	22	0	2	19	4	0	2	25
Danio rerio	23	10	11	1	38	0	0	7	45
Equus caballus	0	20	3	2	22	2	0	1	25
Felis catus	1	21	0	2	17	5	0	2	24
Gallus gallus	8	14	1	2	22	1	0	2	25
Homo sapiens	25	0	0	0	25	0	0	0	25
Macaca mulatta	0	19	6	0	25	0	0	0	25
Mus musculus	24	0	0	0	24	0	0	0	24
Oncorhynchus kisutch	11	37	4	3	19	8	26	2	55
Oryctolagus cuniculus	2	15	2	5	19	3	0	2	24
Ovis aries	0	20	2	3	19	5	0	1	25
Rattus norvegicus	14	6	3	1	24	0	0	0	24
Salmo salar	20	16	17	2	24	9	19	3	55
Sus scrofa	15	5	1	4	20	1	0	4	25
Xenopus tropicalis	3	14	6	1	22	0	0	2	24

<https://doi.org/10.1371/journal.pcbi.1013885.t001>



**Fig 4. New Selenoprofiles utilities for subfamily classification pipeline and prediction refinement.** Initially, we built a phylogenetic tree for every multimember family to produce a bona-fide classification of orthologous groups (a). Next, we devised an automated, fast similarity-based procedure that does not require tree building, now implemented as the *Selenoprofiles orthology* utility (b). Finally, we developed a filter based on the expected selenoproteome for each lineage, now implemented as the *Selenoprofiles lineage* utility (c).

<https://doi.org/10.1371/journal.pcbi.1013885.g004>

reference alignments referred to as “anchors”, containing selected representatives for all subfamilies for each multimember selenoprotein family. New selenoprotein predictions are aligned to anchors, and their protein sequence similarity is computed against each subfamily, ultimately assigning each prediction to its most similar subfamily.

To develop this functionality, we built gene trees and manually partitioned them into orthologous groups (Methods), which replicated well our previous phylogenetic analysis of vertebrate selenoproteins [6]. From these, we derived anchor alignments with reliable subfamily labels, which we then reduced for faster processing by automatic selection of sequence representatives. For subfamily assignment of candidates, we tested various metrics of sequence similarity, differing by how gaps are scored and how different columns are aggregated. We then selected the method that maximized the accuracy of orthology assignment in a set of 3,952 selenoprotein predictions from four multimember families (GPX, TXNRD, DIO and SELENOW/V) in Ensembl genomes (Methods; S2 and S3 Tables).

#### ***Selenoprofiles lineage*: Filtering based on expected selenoproteomes per lineage**

Lastly, we decided to improve Selenoprofiles so that its predictions roughly match the accuracy of experienced manual curation. While our pipeline already features high sensitivity, it can lead to some false positives, typically

in presence of abundant retrotransposed pseudogenes in mammalian genomes [15]. In Selenoprofiles, most pseudogenes are already labelled and discarded due to the presence of in-frame stop codons or frameshifts (e.g., 888 predictions across 306 Ensembl genomes, an average of 2.9 per genome). Yet, some happen to lack these features and are output. At manual inspection, these may be identified by three criteria: 1. they are typically intronless, 2. most lack evidence of transcription, and 3. they are not conserved. Because some functional genes are also intronless (e.g., *Seps2* [6,39]), the first criterion is not universally applicable and cannot be easily employed in an automated pipeline. The second criterion is also not suitable for Selenoprofiles, because it requires additional data (transcriptomics) not universally available. We thus decided to implement the conservation criteria in a novel utility called *Selenoprofiles lineage*, which consists in a filter based on the set of selenoproteins previously observed in other species related to the one under analysis. This functionality depends on a manually curated table defining the precise selenoprotein content of all major vertebrate lineages (Data Availability). We compiled this data based on our recent analysis of vertebrate selenoproteomes [7].

*Selenoprofiles lineage* (Fig 4C) is built on top of the *orthology* utility: gene expectations are defined at the subfamily level, for improved precision. First, the species under analysis is mapped into a vertebrate lineage present in the expectation table (e.g., placentals, birds), either manually or automatically via NCBI taxonomy. Next, the program checks for extra predicted genes: for example, two copies of *GPX4* are found in the mouse genome, but a single one is expected in placentals (Fig 3). Finally, genes are selected based on sequence similarity to the subfamily anchor alignment: only the N top-scoring candidates are kept, where N is the number of expected genes. Besides marking some genes as “Filtered out”, *Selenoprofiles lineage* may also report some subfamilies as “Missing”, if a given genome contains fewer predictions than expected.

Our new utility allows to refine predicted selenoproteomes in full automation. Results are shown in Fig 2, which shows the annotation status of Ensembl and NCBI before and after lineage-based filtering, and in Fig 3, S1 Fig and S2 Fig, wherein filtered and missing genes are shown for each species. Mostly, filtered out predictions corresponded to “Absent” and “Misannotation” rather than “Well annotated” cases, consistent with the notion that functional genes are carefully annotated more often than pseudogenes. Altogether, the *lineage* utility yields predictions with quality comparable to manual analysis. On the other hand, it may miss out on lineage-specific duplications that are not reported in the expectation table, although this was compiled via a thorough analysis of hundreds of vertebrates [7]. We also note that this scored-based approach based may occasionally favour recent retrotransposed pseudogenes over functional copies, when these are highly similar. In the future, we may incorporate further selection criteria, such as exon–intron structure expectations.

The new *orthology* and *lineage* functionalities are available in Selenoprofiles v4.5.7. In the future, we may expand their use to non-vertebrate genomes.

## Discussion

Despite the massive advances in sequencing, genome annotation remains challenging. Due to the recoding of the UGA codon, this problem is aggravated in the case of selenoprotein genes, leading to mispredictions and omissions in public databases. The current study confirmed clear support to this hypothesis: we report that the extent of flawed annotations for selenoprotein genes reach 88% for Ensembl, 23% for RefSeq and 95% for GenBank. In most cases, an overlapping gene annotation is present but does not include the Sec codon, with many possible gene structure arrangements observed (Figs 1 and 2). In summary, accurate annotation of selenoprotein genes in Ensembl is currently largely restricted to a few model organisms; the RefSeq EGAP pipeline performs well overall, owing to their dedicated efforts to improve selenoprotein annotation [13,37], though it consistently struggles with families bearing a C-terminal Sec residue; whereas GenBank displays widespread annotation errors across the vast majority of species and selenoprotein families.

Genome annotations are essential for understanding the information encoded in genomes. This includes elucidating evolutionary relationships, molecular functions, gene expression and regulation, among others, and non-model organisms in particular are subject of unprecedented scientific attention. Indeed, we entered an era of massive sequencing and biodiversity genomics, and in the near future the goal is to sequence all eukaryotic life on Earth, driven by initiatives under the umbrella of the Earth Biogenome Project [40–42]. Naturally, manual curation cannot keep the pace of modern genomics, and having accurate automated pipelines is compelling.

It is now more important than ever that genome annotations of all organisms are as complete as possible, to enable the effective application of methodologies that make use of annotations as a key source of information. This includes a wide array of omics technologies, phylogenetic studies, and also unbiased experimental and computational approaches such as CRISPR-based screening for genes involved in a pathway of interest, and genome scans for positive selection. In these contexts, omissions and misannotations may lead to biases, incomplete design, and misguided interpretation of results.

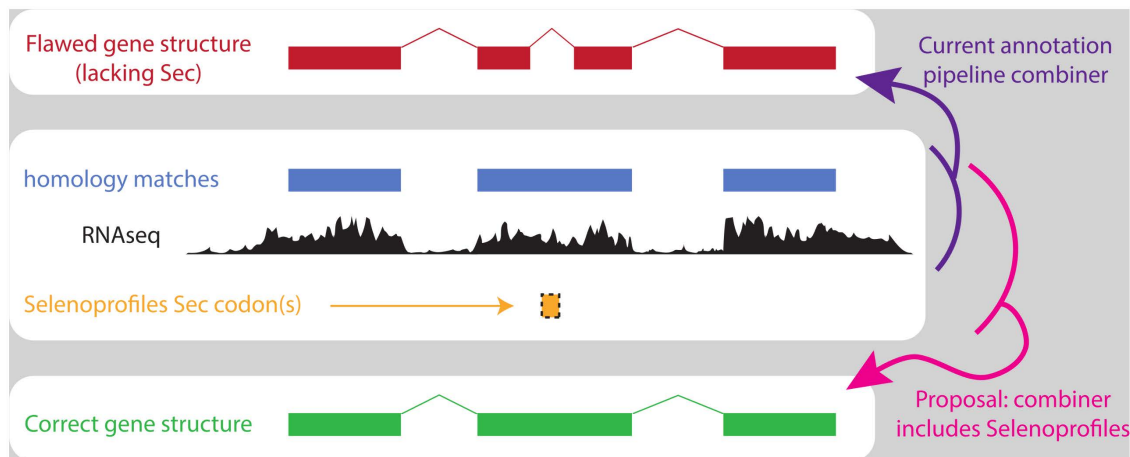
One may think that Sec residues represent only a tiny fraction of our genomes, and may unwisely dismiss their importance. Yet, selenoprotein genes are responsible for many essential functions: most notably, they are core components of the antioxidant and redox homeostasis systems of all vertebrates, which are fundamental for life. In these enzymes, Sec is located almost invariantly in the catalytic site, where it provides a selective advantage over its canonical analog cysteine, which may be related to increased resistance to oxidative inactivation, or enhanced catalysis [2]. Therefore, Sec is arguably the most important residue of many essential proteins, and its correct annotation is compelling.

In previous attempts to address the issue of selenoprotein prediction, a dedicated database was created, called SelenoDB [43,44], now hosted at <http://selenodb.crg.eu/>. Yet, we reckon that a major flaw of this approach is that, in practice, only experts of selenium biology or translational recoding typically use this database. This leaves the great majority of the scientific community, which overwhelmingly uses public annotation databases, unaware and ill-equipped towards this issue. We argue that the only possible effective strategy is to correct selenoproteins in the public annotations that researchers are already using, such as Ensembl and NCBI.

To facilitate this task, we developed a new version of Selenoprofiles, a pipeline that we already showed as effective for homology-based prediction of selenoprotein in genomes [15]. Selenoprofiles is now straightforward to install and run, it has a new extensive online documentation, and features new functionalities: assignment of subfamilies (i.e., precise gene naming) and lineage-based filtering, which improved specificity to a level comparable to manual curation. These are implemented for vertebrates, acting as straightforward test cases, and may be eventually extended to other eukaryotic lineages.

We advocate that the annotation pipelines used for public databases must account for the peculiar genetics of selenoproteins, and we argue that Selenoprofiles can be integrated in existing pipelines to remedy their shortcomings. We reckon that misannotations are mainly due to the in-frame UGA being penalized in the algorithms used, leading to alternative Sec-lacking gene models being selected (Fig 5). Therefore, we suggest that the output by Selenoprofiles should be fed to the combiner module that ultimately generate gene models in annotation pipelines, as an additional information layer besides the genomic data already used (e.g., for Ensembl, GeneWise predictions using Uniprot proteins and Exonerate mapping of available cDNAs) [45]. Selenoprofiles only interrogates the genome sequence, while annotation pipelines typically also use transcriptomic information, so that Selenoprofiles predictions may be less accurate in regions far from Sec residue. We thus recommend that an optimal strategy would be to integrate in public annotation pipelines only the “raw” Sec encoding codons output by Selenoprofiles, prioritizing them as a high-score layer to build gene models (Fig 5).

Having created scalable computational tools suitable for correct annotation of selenoproteins in genomes, we will advocate for their use in public databases. We also argue that, in the near future, other well characterized cases of translational recoding [46] should receive a similar attention.



**Fig 5. Proposed integration of Selenoprofiles in gene annotation pipelines.** While existing pipelines differ in their structure, all follow the same principle in that various information layers are ultimately combined into gene structures. We propose to integrate Sec codons predicted by Selenoprofiles as high-scored layer to correct selenoprotein annotation.

<https://doi.org/10.1371/journal.pcbi.1013885.g005>

## Supporting information

**S1 Fig. Selenoproteome across Ensembl genomes.** Each row represents a species and each column corresponds to a selenoprotein family. Rectangles represent individual “selenocysteine” Selenoprofiles predictions, color-labeled according to their gene annotation in Ensembl: pink for correctly annotated genes, blue for misannotated genes, grey for genes absent from the Ensembl annotation. Two graphical features are used to indicate the departure from the expected selenoproteome described in 2012 by Mariotti et al (coded in the Expectation\_table.csv file at [https://github.com/marco-mariotti/selenoprotein\\_profiles](https://github.com/marco-mariotti/selenoprotein_profiles) v.1.2.0): yellow represent genes expected in that lineage but undetected by Selenoprofiles; and red crosses mark extra genes compared to expectations, therefore filtered by the Selenoprofiles lineage tool.

(PDF)

**S2 Fig. Selenoproteome across RefSeq genomes.** Each row represents a RefSeq genome assembly and each column corresponds to a selenoprotein family. See caption of S1 Fig for further details.

(PDF)

**S3 Fig. Selenoproteome across GenBank genomes.** Each row represents a GenBank genome assembly and each column corresponds to a selenoprotein family. See caption of S1 Fig for further details.

(PDF)

**S1 Table. List of genome assemblies used in this study.** Separated sheets are used for Ensembl and NCBI genomes. For the latter, a column clarifies if the assembly is from GenBank or RefSeq. The column “Crashed” marks as “True” those assemblies that were excluded from analyses due to errors in our selenoprotein prediction pipeline.

(XLSX)

**S2 Table. Explanation of score similarity parameters optimized in Selenoprofiles orthology.** These are arguments available in the Pyaln function score\_similarity. For details, see [https://pyaln.readthedocs.io/en/latest/alignment.html#pyaln.Alignment.score\\_similarity](https://pyaln.readthedocs.io/en/latest/alignment.html#pyaln.Alignment.score_similarity)

(PDF)

**S3 Table. Results of benchmark for Selenoprofiles orthology.** See manuscript text for explanation. (PDF)

### Author contributions

**Conceptualization:** Roderic Guigó, Marco Mariotti.

**Data curation:** Max Ticó.

**Formal analysis:** Max Ticó, Emerson Sullivan.

**Funding acquisition:** Marco Mariotti.

**Methodology:** Max Ticó, Marco Mariotti.

**Supervision:** Marco Mariotti.

**Visualization:** Max Ticó.

**Writing – original draft:** Max Ticó.

**Writing – review & editing:** Roderic Guigó, Marco Mariotti.

### References

1. Labunskyy VM, Hatfield DL, Gladyshev VN. Selenoproteins: molecular pathways and physiological roles. *Physiol Rev.* 2014;94(3):739–77. <https://doi.org/10.1152/physrev.00039.2013> PMID: [24987004](https://pubmed.ncbi.nlm.nih.gov/24987004/)
2. Mariotti M, Gladyshev VN. Selenocysteine-containing proteins. *Redox chemistry and biology of thiols.* Elsevier; 2022. 405–21. <https://doi.org/10.1016/b978-0-323-90219-9.00012-1>
3. Howard MT, Copeland PR. New directions for understanding the codon redefinition required for selenocysteine incorporation. *Biol Trace Elem Res.* 2019;192(1):18–25. <https://doi.org/10.1007/s12011-019-01827-y> PMID: [31342342](https://pubmed.ncbi.nlm.nih.gov/31342342/)
4. Kryukov GV, Castellano S, Novoselov SV, Lobanov AV, Zehtab O, Guigó R, et al. Characterization of mammalian selenoproteomes. *Science.* 2003;300(5624):1439–43. <https://doi.org/10.1126/science.1083516> PMID: [12775843](https://pubmed.ncbi.nlm.nih.gov/12775843/)
5. Ticó M, Mariotti M. The metazoan selenoproteome. *Annu Rev Anim Biosci.* 2025. <https://doi.org/10.1146/annurev-animal-030424-072943> PMID: [41212935](https://pubmed.ncbi.nlm.nih.gov/41212935/)
6. Mariotti M, Ridge PG, Zhang Y, Lobanov AV, Pringle TH, Guigo R, et al. Composition and evolution of the vertebrate and mammalian selenoproteomes. *PLoS One.* 2012;7(3):e33066. <https://doi.org/10.1371/journal.pone.0033066> PMID: [22479358](https://pubmed.ncbi.nlm.nih.gov/22479358/)
7. Ticó M, Mariotti M. Tracing the vertebrate selenoproteome reveals expansions in ray-finned fishes and convergent depletions in tetrapods. Cold Spring Harbor Laboratory; 2025. <https://doi.org/10.1101/2025.05.29.656587>
8. Santesmasses D, Mariotti M, Gladyshev VN. Bioinformatics of selenoproteins. *Antioxid Redox Signal.* 2020;33(7):525–36. <https://doi.org/10.1089/ars.2020.8044> PMID: [32031018](https://pubmed.ncbi.nlm.nih.gov/32031018/)
9. Lobanov AV, Hatfield DL, Gladyshev VN. Eukaryotic selenoproteins and selenoproteomes. *Biochim Biophys Acta.* 2009;1790(11):1424–8. <https://doi.org/10.1016/j.bbagen.2009.05.014> PMID: [19477234](https://pubmed.ncbi.nlm.nih.gov/19477234/)
10. Mariotti M, Salinas G, Gabaldón T, Gladyshev VN. Utilization of selenocysteine in early-branching fungal phyla. *Nat Microbiol.* 2019;4(5):759–65. <https://doi.org/10.1038/s41564-018-0354-9> PMID: [30742068](https://pubmed.ncbi.nlm.nih.gov/30742068/)
11. Chapple CE, Guigó R. Relaxation of selective constraints causes independent selenoprotein extinction in insect genomes. *PLoS One.* 2008;3(8):e2968. <https://doi.org/10.1371/journal.pone.0002968> PMID: [18698431](https://pubmed.ncbi.nlm.nih.gov/18698431/)
12. Otero L, Romanelli-Cedrez L, Turanov AA, Gladyshev VN, Miranda-Vizuete A, Salinas G. Adjustments, extinction, and remains of selenocysteine incorporation machinery in the nematode lineage. *RNA.* 2014;20(7):1023–34. <https://doi.org/10.1261/rna.043877.113> PMID: [24817701](https://pubmed.ncbi.nlm.nih.gov/24817701/)
13. Rajput B, Pruitt KD, Murphy TD. RefSeq curation and annotation of stop codon recoding in vertebrates. *Nucleic Acids Res.* 2019;47(2):594–606. <https://doi.org/10.1093/nar/gky1234> PMID: [30535227](https://pubmed.ncbi.nlm.nih.gov/30535227/)
14. Santesmasses D, Mariotti M, Guigó R. Selenoprofiles: a computational pipeline for annotation of selenoproteins. *Methods Mol Biol.* 2018;1661:17–28. [https://doi.org/10.1007/978-1-4939-7258-6\\_2](https://doi.org/10.1007/978-1-4939-7258-6_2) PMID: [28917034](https://pubmed.ncbi.nlm.nih.gov/28917034/)
15. Mariotti M, Guigó R. Selenoprofiles: profile-based scanning of eukaryotic genome sequences for selenoprotein genes. *Bioinformatics.* 2010;26(21):2656–63. <https://doi.org/10.1093/bioinformatics/btq516> PMID: [20861026](https://pubmed.ncbi.nlm.nih.gov/20861026/)
16. Cunningham F, Allen JE, Allen J, Alvarez-Jarreta J, Amode MR, Armean IM, et al. Ensembl 2022. *Nucleic Acids Res.* 2022;50(D1):D988–95. <https://doi.org/10.1093/nar/gkab1049> PMID: [34791404](https://pubmed.ncbi.nlm.nih.gov/34791404/)

17. Sayers EW, Beck J, Bolton EE, Bourexis D, Brister JR, Canese K, et al. Database resources of the National center for biotechnology information. *Nucleic Acids Res.* 2021;49(D1):D10–7. <https://doi.org/10.1093/nar/gkaa892> PMID: 33095870
18. O'Leary NA, Cox E, Holmes JB, Anderson WR, Falk R, Hem V, et al. Exploring and retrieving sequence and metadata for species across the tree of life with NCBI Datasets. *Sci Data.* 2024;11(1):732. <https://doi.org/10.1038/s41597-024-03571-y> PMID: 38969627
19. Herrero J, Muffato M, Beal K, Fitzgerald S, Gordon L, Pignatelli M, et al. Ensembl comparative genomics resources. *Database (Oxford).* 2016;2016:bav096. <https://doi.org/10.1093/database/bav096> PMID: 26896847
20. Schoch CL, Ciufo S, Domrachev M, Hotton CL, Kannan S, Khovanskaya R, et al. NCBI taxonomy: a comprehensive update on curation, resources and tools. *Database (Oxford).* 2020;2020:baaa062. <https://doi.org/10.1093/database/baaa062> PMID: 32761142
21. Stovner EB, Sætrom P. PyRanges: efficient comparison of genomic intervals in Python. *Bioinformatics.* 2020;36(3):918–9. <https://doi.org/10.1093/bioinformatics/btz615> PMID: 31373614
22. Wickham H. *ggplot2: Elegant Graphics for Data Analysis.* New York: Springer-Verlag. 2009.
23. Huerta-Cepas J, Serra F, Bork P. ETE 3: reconstruction, analysis, and visualization of phylogenomic data. *Mol Biol Evol.* 2016;33(6):1635–8. <https://doi.org/10.1093/molbev/msw046> PMID: 26921390
24. Huerta-Cepas J, Szklarczyk D, Forslund K, Cook H, Heller D, Walter MC, et al. eggNOG 4.5: a hierarchical orthology framework with improved functional annotations for eukaryotic, prokaryotic and viral sequences. *Nucleic Acids Res.* 2016;44(D1):D286–93. <https://doi.org/10.1093/nar/gkv1248> PMID: 26582926
25. Guindon S, Dufayard J-F, Lefort V, Anisimova M, Hordijk W, Gascuel O. New algorithms and methods to estimate maximum-likelihood phylogenies: assessing the performance of PhyML 3.0. *Syst Biol.* 2010;59(3):307–21. <https://doi.org/10.1093/sysbio/syq010> PMID: 20525638
26. Gabaldón T. Large-scale assignment of orthology: back to phylogenetics?. *Genome Biol.* 2008;9(10):235. <https://doi.org/10.1186/gb-2008-9-10-235> PMID: 18983710
27. Hirt RP, Müller S, Embley TM, Coombs GH. The diversity and evolution of thioredoxin reductase: new perspectives. *Trends Parasitol.* 2002;18(7):302–8. [https://doi.org/10.1016/s1471-4922\(02\)02293-6](https://doi.org/10.1016/s1471-4922(02)02293-6) PMID: 12379950
28. Gladyshev VN, Arnér ES, Berry MJ, Brigelius-Flohé R, Bruford EA, Burk RF, et al. Selenoprotein gene nomenclature. *J Biol Chem.* 2016;291(46):24036–40. <https://doi.org/10.1074/jbc.M116.756155> PMID: 27645994
29. Katoh K, Standley DM. MAFFT multiple sequence alignment software version 7: improvements in performance and usability. *Mol Biol Evol.* 2013;30(4):772–80. <https://doi.org/10.1093/molbev/mst010> PMID: 23329690
30. Capella-Gutiérrez S, Silla-Martínez JM, Gabaldón T. trimAl: a tool for automated alignment trimming in large-scale phylogenetic analyses. *Bioinformatics.* 2009;25(15):1972–3. <https://doi.org/10.1093/bioinformatics/btp348> PMID: 19505945
31. Castellano S, Lobanov AV, Chapple C, Novoselov SV, Albrecht M, Hua D, et al. Diversity and functional plasticity of eukaryotic selenoproteins: identification and characterization of the SelJ family. *Proc Natl Acad Sci U S A.* 2005;102(45):16188–93. <https://doi.org/10.1073/pnas.0505146102> PMID: 16260744
32. Novoselov SV, Hua D, Lobanov AV, Gladyshev VN. Identification and characterization of Fep15, a new selenocysteine-containing member of the Sep15 protein family. *Biochem J.* 2006;394(Pt 3):575–9. <https://doi.org/10.1042/BJ20051569> PMID: 16236027
33. Davy T, Castellano S. The genomics of selenium: its past, present and future. *Biochim Biophys Acta Gen Subj.* 2018;1862(11):2427–32. <https://doi.org/10.1016/j.bbagen.2018.05.020> PMID: 29859258
34. Sarangi GK, Romagné F, Castellano S. Distinct patterns of selection in selenium-dependent genes between land and aquatic vertebrates. *Mol Biol Evol.* 2018;35(7):1744–56. <https://doi.org/10.1093/molbev/msy070> PMID: 29669130
35. Castellano S, Novoselov SV, Kryukov GV, Lescure A, Blanco E, Krol A, et al. Reconsidering the evolution of eukaryotic selenoproteins: a novel nonmammalian family with scattered phylogenetic distribution. *EMBO Rep.* 2004;5(1):71–7. <https://doi.org/10.1038/sj.embor.7400036> PMID: 14710190
36. Shchedrina VA, Novoselov SV, Malinouski MY, Gladyshev VN. Identification and characterization of a selenoprotein family containing a diselenide bond in a redox motif. *Proc Natl Acad Sci U S A.* 2007;104(35):13919–24. <https://doi.org/10.1073/pnas.0703448104> PMID: 17715293
37. Goldfarb T, Kodali VK, Pujar S, Brover V, Robbertse B, Farrell CM, et al. NCBI RefSeq: reference sequence standards through 25 years of curation and annotation. *Nucleic Acids Res.* 2024;53: D243. <https://doi.org/10.1093/NAR/GKAE1038>
38. Baclaocos J, Santesmasses D, Mariotti M, Bierla K, Vetick MB, Lynch S, et al. Processive recoding and metazoan evolution of selenoprotein P: up to 132 UGAs in molluscs. *J Mol Biol.* 2019;431: 4381–407. <https://doi.org/10.1016/j.jmb.2019.08.007>
39. Mariotti M, Santesmasses D, Capella-Gutierrez S, Mateo A, Arnan C, Johnson R, et al. Evolution of selenophosphate synthetases: emergence and relocation of function through independent duplications and recurrent subfunctionalization. *Genome Res.* 2015;25(9):1256–67. <https://doi.org/10.1101/gr.190538.115>
40. Lewin HA, Robinson GE, Kress WJ, Baker WJ, Coddington J, Crandall KA, et al. Earth BioGenome project: sequencing life for the future of life. *Proc Natl Acad Sci U S A.* 2018;115(17):4325–33. <https://doi.org/10.1073/pnas.1720115115> PMID: 29686065
41. Corominas M, Marquès-Bonet T, Arnedo MA, Bayés M, Belmonte J, Escrivà H, et al. The catalan initiative for the earth BioGenome project: contributing local data to global biodiversity genomics. *NAR Genom Bioinform.* 2024;6(3):lqae075. <https://doi.org/10.1093/nargab/lqae075> PMID: 39022326

42. Mc Cartney AM, Formenti G, Mouton A, De Panis D, Marins LS, Leitão HG, et al. The European reference genome atlas: piloting a decentralised approach to equitable biodiversity genomics. *NPJ Biodivers.* 2024;3(1):28. <https://doi.org/10.1038/s44185-024-00054-6> PMID: [39289538](https://pubmed.ncbi.nlm.nih.gov/39289538/)
43. Castellano S, Gladyshev VN, Guigó R, Berry MJ. SelenoDB 1.0 : a database of selenoprotein genes, proteins and SECIS elements. *Nucleic Acids Res.* 2008;36(Database issue):D332-8. <https://doi.org/10.1093/nar/gkm731> PMID: [18174224](https://pubmed.ncbi.nlm.nih.gov/18174224/)
44. Romagné F, Santesmasses D, White L, Sarangi GK, Mariotti M, Hübler R, et al. SelenoDB 2.0: annotation of selenoprotein genes in animals and their genetic diversity in humans. *Nucleic Acids Res.* 2014;42(Database issue):D437-43. <https://doi.org/10.1093/nar/gkt1045> PMID: [24194593](https://pubmed.ncbi.nlm.nih.gov/24194593/)
45. Aken BL, Ayling S, Barrell D, Clarke L, Curwen V, Fairley S, et al. The Ensembl gene annotation system. *Database (Oxford).* 2016;2016:baw093. <https://doi.org/10.1093/database/baw093> PMID: [27337980](https://pubmed.ncbi.nlm.nih.gov/27337980/)
46. Rodnina MV, Korniy N, Klimova M, Karki P, Peng B-Z, Senyushkina T, et al. Translational recoding: canonical translation mechanisms reinterpreted. *Nucleic Acids Res.* 2020;48(3):1056–67. <https://doi.org/10.1093/nar/gkz783> PMID: [31511883](https://pubmed.ncbi.nlm.nih.gov/31511883/)