

Manuscript Number:	GIGA-D-17-00074	
Full Title:	GeneSeqToFamily: the Ensembl Compara GeneTrees pipeline as a Galaxy workflow	
Article Type:	Technical Note	
Funding Information:	Biotechnology and Biological Sciences Research Council (BBSRC strategic funds)	Mr. Anil S Thanki Dr Wilfried Haerty Dr Robert P Davey
	Biotechnology and Biological Sciences Research Council (BBSRC Biomathematics and Bioinformatics Training fund (2014).)	Dr Nicola Soranzo
Abstract:	<p>Background Gene duplication is a major factor contributing to evolutionary novelty, and the contraction or expansion of gene families has often been associated with morphological, physiological and environmental adaptations. The study of homologous genes helps us to understand the evolution of gene families. It plays a vital role in finding ancestral gene duplication events as well as identifying genes that have diverged from a common ancestor under positive selection. There are various tools available, such as MSOAR, OrthoMCL and HomoloGene, to identify gene families and visualise syntenic information between species, providing an overview of syntenic regions evolution at the family level. Unfortunately, none of them provide information about structural changes within genes, such as the conservation of ancestral exon boundaries amongst multiple genomes. The Ensembl GeneTrees computational pipeline generates gene trees based on coding sequences and provides details about exon conservation, and is used in the Ensembl Compara project to discover gene families.</p> <p>Findings A certain amount of expertise is required to configure and run the Ensembl Compara GeneTrees pipeline via command line. Therefore, we have converted the command line Ensembl Compara GeneTrees pipeline into a Galaxy workflow, called GeneSeqToFamily, and provided additional functionality. This workflow uses existing tools from the Galaxy ToolShed, as well as providing additional wrappers and tools that are required to run the workflow.</p> <p>Conclusions GeneSeqToFamily represents the Ensembl Compara pipeline as a set of interconnected Galaxy tools, so they can be run interactively within the Galaxy's user-friendly workflow environment while still providing the flexibility to tailor the analysis by changing configurations and tools if necessary. Additional tools allow users to subsequently visualise the gene families produced by the workflow, using the Aequatus.js interactive tool, which has been developed as part of the Aequatus software project.</p>	
Corresponding Author:	Anil S Thanki, MSc Earlham Institute Norwich, Norfolk UNITED KINGDOM	
Corresponding Author Secondary Information:		
Corresponding Author's Institution:	Earlham Institute	
Corresponding Author's Secondary Institution:		
First Author:	Anil S Thanki, MSc	
First Author Secondary Information:		

Order of Authors:	Anil S Thanki, MSc
	Nicola Soranzo, PhD
	Wilfried Haerty, PhD
	Robert P Davey, PhD
Order of Authors Secondary Information:	
Opposed Reviewers:	
Additional Information:	
Question	Response
Are you submitting this manuscript to a special series or article collection?	No
<p>Experimental design and statistics</p> <p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist. Information essential to interpreting the data presented should be made available in the figure legends.</p> <p>Have you included all the information requested in your manuscript?</p>	Yes
<p>Resources</p> <p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist?</p>	Yes
<p>Availability of data and materials</p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.</p>	Yes

Have you have met the above requirement as detailed in our [Minimum Standards Reporting Checklist](#)?

GeneSeqToFamily: the Ensembl Compara GeneTrees pipeline as a Galaxy workflow

Anil S. Thanki¹, Nicola Soranzo¹, Wilfried Haerty¹, Robert P. Davey¹

1. Earlham Institute (EI), Norwich Research Park, Norwich NR4 7UZ, UK

Abstract

Background

Gene duplication is a major factor contributing to evolutionary novelty, and the contraction or expansion of gene families has often been associated with morphological, physiological and environmental adaptations. The study of homologous genes helps us to understand the evolution of gene families. It plays a vital role in finding ancestral gene duplication events as well as identifying genes that have diverged from a common ancestor under positive selection. There are various tools available, such as MSOAR, OrthoMCL and HomoloGene, to identify gene families and visualise syntenic information between species, providing an overview of syntenic regions evolution at the family level. Unfortunately, none of them provide information about structural changes within genes, such as the conservation of ancestral exon boundaries amongst multiple genomes. The Ensembl GeneTrees computational pipeline generates gene trees based on coding sequences and provides details about exon conservation, and is used in the Ensembl Compara project to discover gene families.

Findings

A certain amount of expertise is required to configure and run the Ensembl Compara GeneTrees pipeline via command line. Therefore, we have converted the command line Ensembl Compara GeneTrees pipeline into a Galaxy workflow, called GeneSeqToFamily, and provided additional functionality. This workflow uses existing tools from the Galaxy ToolShed, as well as providing additional wrappers and tools that are required to run the workflow.

Conclusions

GeneSeqToFamily represents the Ensembl Compara pipeline as a set of interconnected Galaxy tools, so they can be run interactively within the Galaxy's user-friendly workflow environment while still providing the flexibility to tailor the analysis by changing configurations and tools if necessary. Additional tools allow users to subsequently visualise the gene families produced by the workflow, using the Aequatus.js interactive tool, which has been developed as part of the Aequatus software project.

Keywords

Galaxy, Pipeline, Workflow, Genomics, Comparative Genomics, Homology, Orthology, Paralogy, Phylogeny, Gene Family, Alignment, Compara, Ensembl

Introduction

The phylogenetic information inferred from the study of homologous genes helps us to understand the evolution of gene families, which plays a vital role in finding ancestral gene duplication events as well as identifying regions under positive selection within species [1]. In order to investigate these low-level comparisons between gene families, the Ensembl Compara GeneTrees gene orthology and paralogy prediction software suite [2] was developed as a pipeline that uses TreeBest [3] [4] (part of TreeFam [5]) to find internal structural-level synteny for homologous genes. TreeBeST implements multiple independent phylogenetic methods and can merge their results in a consensus tree whilst trying to minimise duplications and deletions relative to a known species tree. This allows TreeBeST to take advantage of the fact that DNA-based methods are often more accurate for closely related parts of trees, while protein-based trees are better at longer distances.

The Ensembl GeneTrees pipeline comprises seven basic steps, starting from a set of protein sequences and performing similarity searching and multiple large-scale alignments to infer homology among them, using various tools: BLAST [6], hcluster_sg [7], T-Coffee [8], and phylogenetic tree construction tools, including TreeBeST. Whilst all these tools are freely available, most are specific to certain computing environments, are only usable via the command line, and require many dependencies to be fulfilled. Therefore, users are not always sufficiently expert in system administration in order to install, run, and debug the various tools at each stage in a chain of processes. To help ease the complexity of running the GeneTrees pipeline, we have employed the Galaxy bioinformatics analysis platform to relieve the burden of managing these system-level challenges.

Galaxy is an open-source framework for running a broad collection of bioinformatics tools via a user-friendly web interface [9]. No client software is required other than a recent web browser, and users are able to run tools singly or aggregated into interconnected pipelines, called *workflows*. Galaxy enables users to not only create, but also share workflows with the community. In this way, it helps users who have little or no bioinformatics expertise to run potentially complex pipelines in order to analyse their own data and interrogate results within a single online platform. Furthermore, pipelines can be published in a scientific paper or in a repository such as myExperiment [10] to encourage transparency and reproducibility.

In addition to analytical tools, Galaxy also contains plugins [11] for data visualisation. Galaxy visualisation plugins may be interactive and can be configured to visualise various data types, for example, bar plots, scatter plots, and phylogenetic trees. It is also possible to develop custom visualisation plugins and easily integrate them into Galaxy. As the output of the GeneSeqToFamily workflow is not conducive to human readability, we also provide a data-to-visualisation plugin based on the Aequatus software [12]. Aequatus.js [13] is a new JavaScript library for the visualisation of homologous genes, which is extracted from the standalone Aequatus tool. It provides a detailed view of gene structure across gene families, including shared exon information within gene families alongside gene tree representations. It also shows

1
2
3
4 details about the type of interrelation event that gave rise to the family, such as speciation,
5 duplication, and gene splits.
6

7 8 9 Methods

10 The GeneSeqToFamily workflow has been developed to run the Ensembl Compara software
11 suite within the Galaxy environment, combining various tools alongside preconfigured
12 parameters obtained from the Ensembl Compara pipeline to produce gene trees. Among the
13 tools used in GeneSeqToFamily (listed in Table 1), some were existing tools in the Galaxy
14 ToolShed [14], such as NCBI BLAST, TranSeq, Tralign and various format converters.
15 Additional tools that are part of the pipeline were developed at the Earlham Institute (EI) and
16 submitted to the ToolShed, i.e. *BLAST parser*, *hcluster_sg*, *hcluster_sg parser*, *T-Coffee*,
17 *TreeBeST best* and *Gene Alignment and Family Aggregator*. Finally, we developed helper tools
18 that are not part of the workflow itself, but aid the generation of input data for the workflow and
19 these are also in the ToolShed, i.e. *Get features by Ensembl ID*, *Get sequences by Ensembl ID*,
20 *Select longest CDS per gene*, *ETE species tree generator* and *GeneSeqToFamily preparation*.
21
22
23
24
25

26 The workflow comprises 7 main steps, starting with translation from input coding sequences
27 (CDS) to protein sequences, finding subsequent pairwise alignments of those protein
28 sequences using BLASTP, and then the generation of clusters from the alignments using
29 *hcluster_sg*. The workflow then splits into two simultaneous paths, whereby in one path it
30 performs the multiple sequence alignment (MSA) for each cluster using T-Coffee, and in the
31
32
33

34 Figure 1: Overview of the GeneSeqToFamily workflow

35
36
37 other, generates a gene tree with TreeBeST taking the cluster alignment and a species tree as
38 input. Finally, these paths merge to aggregate the MSA, the gene tree and the gene feature
39 information (transcripts, exons, and so on) into an SQLite [15] database for visualisation and
40 downstream reuse. Each step of the workflow along with the data preparation steps is explained
41 in detail below.
42
43
44

45 Figure 2: Screenshot from the Galaxy Workflow Editor, showing the GeneSeqToFamily
46 workflow
47
48
49

50 Table 1: Galaxy tools used in the workflows
51
52

Tool Name	Tool ID	Version	Developed at EI		ToolShed Reference
			Tool	Wrapper	
Get sequences by Ensembl ID	get_sequences	0.1.2	Yes	Yes	[16]

53
54
55
56
57
58
59
60
61
62
63
64
65

Get features by Ensembl ID	get_feature_info	0.1.2	Yes	Yes	[17]
Select longest CDS per gene	ensembl_longest_cds_per_gene	0.0.2	Yes	Yes	[18]
ETE species tree generator	ete_species_tree_generator	3.0.0b35	Yes	Yes	[19]
GeneSeqToFamily preparation	gstf_preparation	0.3.0	Yes	Yes	[20]
Transeq	EMBOSS: transeq101	5.0.0	No	No	[21]
NCBI BLAST+ makeblastdb	ncki_makeblastdb	0.1.07	No	No	[22]
NCBI BLAST+ blastp	ncki_blastp_wrapper	0.1.07	No	No	[22]
BLAST parser	blast_parser	0.1.1	Yes	Yes	[23]
hcluster_sg	hcluster_sg	0.5.1	No	Yes	[24]
hcluster_sg parser	hcluster_sg_parser	0.2.0	Yes	Yes	[25]
Filter by FASTA IDs	filter_by_fasta_ids	1.0	No	No	[26]
T-Coffee	t_coffee	11.0.8	No	Yes	[27]
Tranalign	EMBOSS: tranalign100	5.0.0	No	No	[21]
TreeBeST best	treebest_best	1.9.2	No	Yes	[28]
Gene Alignment and Family Aggregator	gafa	0.3.0	Yes	Yes	[29]
UniProt ID mapping and retrieval	uniprot_rest_interface	0.1	No	No	[30]

Data generation and preparation

We have developed a number of tools that assist in preparing the datasets needed by the workflows.

Ensembl REST API tools

Galaxy tools were developed which utilise the Ensembl REST API [31] to retrieve sequence information (*Get sequences by Ensembl ID*) and feature information (*Get features by Ensembl ID*) by Ensembl ID from the Ensembl service. REST (REpresentational State Transfer) is an architecture style for designing networked applications [32] which encourages the use of standardised HTTP technology to send and receive data between computers. As such, these tools are designed to help users to retrieve existing data from Ensembl rather than requiring

1
2
3
4 them to manually download datasets to their own computers and then subsequently uploading
5 them into the workflow.
6

7
8 We have also developed:

- 9 ● *Select longest CDS per gene*, which filters a CDS FASTA file from Ensembl retaining
10 only the longest CDS sequence for each gene
- 11 ● *ETE species tree generator*, which uses the ETE toolkit [33] to generate a species tree
12 from a list of species names or taxon IDs through NCBI Taxonomy.
13
14

15 16 17 GeneSeqToFamily workflow

18 19 0. GeneSeqToFamily preparation

20 Before GeneSeqToFamily can be run, a data preparation step must be carried out. We have
21 developed a tool called *GeneSeqToFamily preparation* to preprocess the input datasets (gene
22 feature information and CDS) for the GeneSeqToFamily workflow. It converts a set of gene
23 feature information files in GFF3 [34] and/or JavaScript Object Notation (JSON) [35] format to
24 an SQLite database. It also modifies all CDS FASTA header lines by appending the species
25 name to the transcript identifier, as required by *TreeBeST best*. We decided to use an SQLite
26 database to store the gene feature information because:
27

- 28 ● the GFF3 format has a relatively inconvenient and unstructured additional information
29 field (9th column)
- 30 ● searching is much faster and more memory efficient in a database than in a text file like
31 JSON or GFF3, especially when dealing with feature information for multiple large
32 genomes
33
34
35
36
37

38 1. CDS translation

39 40 Transeq

41 Transeq, part of the European Molecular Biology Open Software Suite (EMBOSS) [36], is a tool
42 to generate six-frame translation of nucleic acid sequences to their corresponding peptide
43 sequences. Here we use Transeq to convert a CDS to protein sequences in order to run
44 BLASTP and find protein clusters. However, since downstream tools in the pipeline such as
45 TreeBeST require nucleotide sequences to generate a gene tree, the protein sequences cannot
46 be directly used as workflow input and are instead generated with Transeq.
47
48
49

50 2. Pre-clustering alignment

51 52 BLAST

53 This workflow uses the BLAST wrappers [37] developed to run BLAST+ tools within Galaxy.
54 BLASTP is run over the set of sequences against the database of the same input, as is the case
55 with BLAST-all, in order to form clusters of related sequences.
56
57
58
59
60
61
62
63
64
65

BLAST parser

BLAST parser is a small Galaxy tool to convert the BLAST output into the input format required by *hcluster_sg*. It takes the BLAST 12-column output [38] as input and generates a 3-column tabular file, comprising the BLAST query, the hit result, and the edge weight. The weight value is simply calculated as minus \log_{10} of the BLAST e-value, replacing this with 100 if this value is greater than 100. It also removes the self-matching BLAST results.

3. Cluster generation

hcluster_sg

hcluster_sg performs hierarchical clustering under mean distance for sparse graphs. It reads an input file that describes the similarity between two sequences, and iterates through the process of grouping two nearest nodes at each iteration. *hcluster_sg* outputs a single list of gene clusters, each comprising a set of sequence IDs present in that cluster. This list needs to be reformatted using the *hcluster_sg parser* tool in order to be suitable for input into T-Coffee and TreeBeST (see below).

hcluster_sg parser

hcluster_sg parser converts the *hcluster_sg* output into a collection of lists of IDs, one list for each cluster. Each of these clusters will then be used to generate a gene tree via TreeBeST. The tool can also filter out clusters with a number of elements outside a specified range. The IDs contained in all discarded clusters are collected in separate output dataset. Since TreeBeST requires at least 3 genes to generate a gene tree, we configured the tool to filter out clusters with less than 3 genes.

Filter by FASTA IDs, which is available from the Galaxy ToolShed, is used to create separate FASTA files using the sequence IDs listed in each gene cluster.

4. Cluster alignment

T-Coffee

T-Coffee is a MSA package, but can also be used to combine the output of other alignment methods (Clustal, MAFFT, Probcons, MUSCLE) into a single alignment. T-Coffee can align both nucleotide and protein sequences [8], and we use it to align the protein sequences in each cluster generated by *hcluster_sg*.

We modified the Galaxy wrapper for T-Coffee to take a single FASTA (as normal) and an optional list of FASTA IDs to filter. If a list of IDs is provided, the wrapper will pass only those sequences to T-Coffee, which will perform the MSA for that set of sequences, thus removing the need to create thousands of intermediate Galaxy datasets.

5. Gene tree construction

Tranalign

Tranalign [36] is a tool that reads a set of nucleotide sequences and a corresponding aligned set of protein sequences and returns a set of aligned nucleotide sequences. Here we use it to generate CDS alignments of gene sequences using the protein alignments produced by T-Coffee.

TreeBeST 'best'

TreeBeST (Tree Building guided by Species Tree) is a tool to generate, manipulate, and display phylogenetic trees and can be used to build gene trees based on a known species tree.

The 'best' command of TreeBeST builds 5 different gene trees from a FASTA alignment file using different phylogenetic algorithms, then merges them into a single consensus tree using a species tree as a reference. In GeneSeqToFamily, *TreeBeST best* uses the nucleotide MSAs generated by Tranalign (at least 3 sequences are required) and a user-supplied species tree in Newick format [39] (either produced by a third-party software or through the *ETE species tree generator* data preparation tool, described above) to produce a GeneTree for each family represented also in Newick format. The resulting GeneTree also includes useful annotations specifying phylogenetic information of events responsible for the presence/absence of genes, for example, 'S' means speciation event, 'D' means duplication, and 'DCS' denotes the duplication score.

6. Gene Alignment and Family Aggregation

Gene Alignment and Family Aggregator (GAFA)

GAFA is a Galaxy tool which generates a single SQLite database containing the gene trees and MSAs, along with gene features, in order to provide a reusable, persistent data store for visualisation of synteny information with Aequatus. GAFA requires:

- gene trees in Newick format,
- the protein MSAs in fasta_aln format from *T-Coffee* and
- gene feature information generated with the *GeneSeqToFamily preparation* tool.

Internally, GAFA converts each MSA from fasta_aln format to a simple CIGAR string [40]. An example of CIGAR strings for aligned sequences is shown in Figure 3, in which each CIGAR string subset changes according to other sequences.

The simple schema [41] for the generated SQLite database is shown in Figure 4.

Figure 3: Showing how CIGAR for multiple sequence alignment is generated

Figure 4: Schema of the GAFA SQLite database

7. Visualisation

Aequatus visualisation plugin

The SQLite database generated by the GAFA tool can be rendered using a new visualisation plugin, Aequatus.js. The Aequatus.js library, developed at EI as part of the Aequatus project, has been configured to be used within Galaxy to visualise homologous gene structure and gene family relationships. This allows users to interrogate not only the evolutionary history of the gene family but also the structural variation (exon gain/loss) within genes across the phylogeny. Aequatus.js is available to download from GitHub [41], as visualisation plugins cannot yet be submitted to the Galaxy ToolShed.

Finding homology information for orphan genes

Although the GeneSeqToFamily workflow will assign most of the genes to orthogroups, many genes within a species might appear to be unique without homologous relationship to any other genes from other species. This observation could be the consequence of the parameters selected, choice of species, incomplete annotations. This could also reflect real absence of homology such as for rapidly evolving gene families. In addition to the GeneSeqToFamily workflow, we also developed two associated workflows to further annotate these genes by:

- 1) Retrieving a list of orphan genes from the GeneSeqToFamily workflow (see Figure 5) as follows:
 - a) Find the IDs of the sequences present in the input CDS of the GeneSeqToFamily workflow, but not in the result of *BLAST parser* from the same workflow
 - b) Add to this list the IDs of the sequences discarded by *hcluster_sg parser*
 - c) From the input CDS dataset, retrieve the respective sequence for each CDS ID (from the step above) using *Filter by FASTA IDs*These unique CDS can be fed into the SwissProt workflow below to find homologous genes in other species.
- 2) Finding homologous genes for some genes of interest using SwissProt (see Figure 6) as follows:
 - a) Translate CDS into protein sequences using *Transeq*
 - b) Run BLASTP for the protein sequences against the SwissProt database (from NCBI)
 - c) Extract UniProt IDs from these BLASTP results, using the preinstalled Galaxy tool *Cut columns from a table* (tool id *Cut1*)
 - d) Retrieve Ensembl IDs (representing genes and/or transcripts) for each UniProt ID using *UniProt ID mapping and retrieval*
 - e) Get genomic information for each gene ID and CDS for each transcript ID from the core Ensembl database using *Get features by Ensembl ID* and *Get sequences by Ensembl ID* respectively.

The results from this second workflow can be subsequently used as input to GeneSeqToFamily for familial analysis.

1
2
3
4
5 Figure 5: Screenshot from the Galaxy Workflow Editor, showing the orphan gene finding
6 workflow

7
8 Figure 6: Screenshot from the Galaxy Workflow Editor, showing the SwissProt workflow
9

10 11 Example use case

12 Since BLASTP plays a crucial role in determining gene families, we tested this workflow on a
13 large dataset of CDS from three vertebrate species in order to set a benchmark and find
14 optimum parameters to run the workflow. We downloaded the CDS sequences for *Sarcophilus*
15 *harrisii* (Tasmanian devil), *Mus musculus* (Mouse), and *Ornithorhynchus anatinus* (Platypus)
16 from Ensembl (release 87) and filtered them to retain only the longest transcript per gene (as in
17 the Ensembl Compara pipeline), obtaining a total of 62,597 CDS. We then ran the
18 GeneSeqToFamily workflow on them using various BLASTP parameters (as shown in Table 3),
19 in order to identify the optimal values for the workflow to generate a gene tree in which
20 members are possibly evolved from a single ancestor gene, and usually with identical
21 biochemical functions such as proteins.
22
23
24
25
26
27

28 Our results show that the number of gene families can vary quite distinctly with different
29 BLASTP parameters. Stringent parameters (Analysis 6) result in a large number of smaller
30 families, while relaxed parameters (Analysis 1) generate a smaller number of large families,
31 which may include distantly related genes. By testing different parameters and comparing the
32 analyses with third party tools such as PantherDB to validate the results against known families,
33 we chose those parameters listed as Analysis 5. These values seem to consistently generate
34 legitimate sets of gene families with closely related family members based on the datasets we
35 tested.
36
37
38
39

40 There are caveats, however. BLASTP parameters used in Analysis 5 restrict the maximum
41 number of target sequences per query (max_target_seqs) to 3 (the first hit when using all-
42 versus-all BLAST will always be the query sequence itself). The minimum query coverage per
43 High-scoring Segment Pair (HSP) (qcovhsp) is set to 90% and e-value cut-off to 1e-10, in order
44 to find the HSP closest to the query thus allowing partial matches which could be seen in the
45 event of gene split. If the input CDS contain multiple alternative transcripts per gene, we
46 recommend setting the max_target_seqs parameter to 4 rather than 3 to get a wider range of
47 results from BLAST, thereby helping to generate gene families with matching genes together
48 with alternative transcripts. In contrast setting a value of 3 for the max_target_seqs parameter
49 will restrict the search to only 3 matches per query, and the presence of alternative transcripts
50 will decrease the likelihood of finding matches from other genes, thus increasing the likelihood
51 of splitting of a gene tree into multiple trees and adversely inflating the number of families.
52
53
54
55
56

57 Table 2: Results of the GeneSeqToFamily workflow run on 62,597 CDS from 3 species using 6
58 different BLAST parameter configurations, the complete list of which are shown in Table 3.
59
60
61
62
63
64
65

Summary						
Analysis	1	2	3	4	5	6
No of families	8,398	13,547	13,675	12,977	13,496	17,090
Filtered out (>200)	12	1	1	1	1	0
Filtered out (<3)	2,655	2,600	3,142	1,030	2,965	5,851
Filtered families to consider	5,731	10,946	10,532	10,946	10,530	11,099
Average Family size	7.15	5.11	4.50	5.11	4.50	3.65
Median Family Size	4	4	3	4	3	3
Largest Family Size	7,885	827	599	827	599	64
Smallest Family Size	1	1	1	1	1	1

Table 3: Complete list of BLAST parameter configurations. Analysis 5 is highlighted to denote those parameters that were chosen to be used as workflow defaults.

Analysis ID	e-value	Max targets per query	Min coverage
1	1e-03 (Default)	0 (Default)	0 (Default)
2	1e-03	3	0
3	1e-03	3	90%
4	1e-10	3	0
5	1e-10	3	90%
6	1e-10	2	90%

To validate the biological relevance of results from the GeneSeqToFamily workflow, we analysed a smaller set of 23 homologous genes (39 transcripts) from *Pan troglodytes* (chimpanzee), *Homo sapiens* (human), *Rattus norvegicus* (rat), *Mus musculus* (mouse), *Sus scrofa* (pig) and *Canis familiaris* (domesticated dog). These genes are a combination of those found in four gene families, i.e. monoamine oxidases (MAO) A and B, insulin receptor (INSR), BRCA1-associated ATM activator 1 (BRAT1), and were chosen because they are present in all 6 species yet distinct from each other. Though MAO gene variants (A and B) are 70% similar, a single gene tree for all MAO genes could be generated if appropriate parameters are not selected. As such, these genes represent a reliable dataset to test whether the GenSeqToFamily workflow can reproduce already known gene families.

1
2
3
4 Before running the workflow, feature information and CDS for the selected genes were retrieved
5 from the core Ensembl database using the helper tools described above (*Get features by*
6 *Ensembl ID* and *Get sequences by Ensembl ID* respectively). A species tree was generated
7 using *ETE species tree generator* and CDS were prepared with *GeneSeqToFamily preparation*.
8 We ran the GeneSeqToFamily workflow on these data using the parameters shown in Analysis
9 5 of Table 3, but we set `max_target_seqs` as 4 (as described in the previous use case) to get a
10 wider range of results from BLAST because our dataset includes alternative transcripts. This
11 workflow generated 4 different gene trees, one for each gene family. Figure 7, 8, 9 and 10 show
12 the resulting gene trees for MAOA, MAOB, BRAT1 and INSR gene families. Different colours of
13 the nodes in each gene tree on the left-hand-side highlight potential evolutionary events, such
14 as speciation, duplication, and gene splits. Homologous genes showing shared exons use the
15 same colour in each representation, including insertions (black blocks) and deletions (red lines).
16 The GeneTrees for these genes are already available in Ensembl and we used them to validate
17 our findings [42] [43] [44] [45]. Our gene trees exactly matched the Ensembl GeneTrees,
18 showing that the workflow generates biologically valid results. We have provided the underlying
19 data for this example along with the submitted workflow in figshare [46].
20
21
22
23
24
25

26
27 Figure 7: Homologous genes of MAOA of *Canis familiaris* from *Mus musculus*, *Pan*
28 *troglodytes*, *Homo sapiens*, *Rattus norvegicus*, *Sus scrofa* and *Canis familiaris*.
29

30
31 Figure 8: Homologous genes of MAOB of *Canis familiaris* from *Mus musculus*, *Pan*
32 *troglodytes*, *Homo sapiens*, *Rattus norvegicus*, *Sus scrofa* and *Canis familiaris*.
33

34
35 Figure 9: Homologous genes of BRAT1 of *Canis familiaris* from *Mus musculus*, *Pan*
36 *troglodytes*, *Homo sapiens*, *Rattus norvegicus*, *Sus scrofa* and *Canis familiaris*.
37

38
39 Figure 10: Homologous genes of INSR of *Canis familiaris* from *Mus musculus*, *Pan*
40 *troglodytes*, *Homo sapiens*, *Rattus norvegicus*, *Sus scrofa* and *Canis familiaris*.
41

42 Conclusion

43
44
45 The ultimate goal of the GeneSeqToFamily is to provide a user-friendly workflow to analyse and
46 discover homologous genes using the Ensembl Compara GeneTrees pipeline within the Galaxy
47 framework, where users can interrogate genes of interest without using the command-line whilst
48 still providing the flexibility to tailor analysis by changing configurations and tools if necessary.
49 We have shown it to be an accurate, robust, and reusable method to elucidate and analyse
50 potentially large numbers of gene families in a range of model and non-model organisms. The
51 workflow stores the resulting gene families into a SQLite database, which can be visualised
52 using the Aequatus.js interactive tool, as well as shared as a complete reproducible container
53 for potentially large gene family datasets.
54
55
56

57
58 Gradually, we hope that the Galaxy community will undertake their own analyses and feedback
59 improvements to various tools, and publish successful combinations of parameters used in the
60
61
62
63
64
65

1
2
3
4 GeneSeqToFamily workflow. We encourage this process by allowing users to share their own
5 version of GeneSeqToFamily workflow for appraisal by the community.
6
7

8 **Future directions**

9

10 In terms of core workflow functionality, we would like to incorporate pairwise alignment between
11 pairs of genes for closely related species in addition of the MSA for the gene family, which will
12 help users to compare orthologs and paralogs in greater detail.
13
14

15 We also plan to include explicit integration of the PantherDB resources [47]. Protein ANalysis
16 Through Evolutionary Relationships (PANTHER) is a classification system to characterise
17 known proteins and genes in order to certify genomic annotation. Association of PantherDB with
18 GeneSeqToFamily will enable the automation of gene family validation and add supplementary
19 information about those gene families, which could then be used in turn to further validate novel
20 genomics annotation.
21
22

23
24 We also plan to add the ability to query the *GAF*A SQLite database using keywords, to make it
25 easy for users to find gene trees which include their genes of interest without needing to delve
26 into the database itself.
27
28

30 **Availability and requirements**

31
32

33 **Project name:** GeneSeqToFamily

34 **Project home page:** [https://github.com/TGAC/earlham-
35 galaxytools/tree/master/workflows/GeneSeqToFamily](https://github.com/TGAC/earlham-galaxytools/tree/master/workflows/GeneSeqToFamily)

36 **Archived version:** 0.1.0

37 **Operating system(s):** Platform independent

38 **Programming language:** JavaScript, Perl, Python, XML, SQL

39 **Other Requirements:** Web Browser; for development: Galaxy

40 **Any restrictions to use by non-academics:** None

41 **License:** The MIT License
42
43
44
45
46

47 **Availability of supporting data**

48
49

50 The example files and additional data sets supporting the results of this article are available in
51 figshare [46].
52
53
54

55 **Acknowledgements**

56
57
58
59
60
61
62
63
64
65

1
2
3
4 AT, WH and RPD are supported by BBSRC institute strategic programme grant funds awarded
5 to EI. NS is funded under the BBSRC Biomathematics and Bioinformatics Training fund (2014).
6 This research was supported in part by the NBI Computing infrastructure for Science (CiS)
7 group who provide technical support and maintenance to EI's High Performance Computing
8 cluster and storage systems, enabling us to develop this workflow.
9

10
11
12 We would like to thank Matthieu Muffato from the European Bioinformatics Institute (EBI) for his
13 advices during the initial stage of the project.
14
15

16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65

1. Jensen JD, Wong A, Aquadro CF. Approaches for identifying targets of positive selection. *Trends Genet.* 2007;23:568–77.

2. Vilella AJ, Severin J, Ureta-Vidal A, Heng L, Durbin R, Birney E. EnsemblCompara GeneTrees: Complete, duplication-aware phylogenetic trees in vertebrates. *Genome Res.* 2008;19:327–35.

3. Ensembl. Ensembl/treebest. GitHub. <https://github.com/Ensembl/treebest>. Accessed 26 Jan 2016.

4. Heng L. Constructing the TreeFam database. The Institute of Theoretical Physics, Chinese Academic of Science; 2006. <http://pfigshare-u-files.s3.amazonaws.com/1421613/PhDthesisliheng2006English.pdf>.

5. Ruan J, Li H, Chen Z, Coghlan A, Coin LJM, Guo Y, et al. TreeFam: 2008 Update. *Nucleic Acids Res.* 2008;36 Database issue:D735–40.

6. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment search tool. *J Mol Biol.* 1990;215:403–10.

7. Li H et al. hcluster_sg: hierarchical clustering software for sparse graphs. <https://github.com/douglasgscfield/hcluster>. Accessed 26 Jan 2016.

8. Notredame C, Higgins DG, Heringa J. T-Coffee: A novel method for fast and accurate multiple sequence alignment. *J Mol Biol.* 2000;302:205–17.

9. Afgan E, Baker D, van den Beek M, Blankenberg D, Bouvier D, Čech M, et al. The Galaxy platform for accessible, reproducible and collaborative biomedical analyses: 2016 update. *Nucleic Acids Res.* 2016;44:W3–10.

10. Goble CA, Bhagat J, Aleksejevs S, Cruickshank D, Michaelides D, Newman D, et al. myExperiment: a repository and social network for the sharing of bioinformatics workflows. *Nucleic Acids Res.* 2010;38 Web Server issue:W677–82.

11. Goecks J, Eberhard C, Too T, Galaxy Team, Nekrutenko A, Taylor J. Web-based visual

1
2
3
4 analysis for high-throughput genomics. BMC Genomics. 2013;14:397.
5

6 12. Thanki AS, Ayling S, Herrero J, Davey RP. Aequatus: An open-source homology browser.
7 bioRxiv. 2016;:055632. doi:10.1101/055632.
8

9 13. TGAC. TGAC/aequatus.js. GitHub. <https://github.com/TGAC/aequatus.js>. Accessed 26 Jan
10 2016.
11

12 14. Blankenberg D, Von Kuster G, Bouvier E, Baker D, Afgan E, Stoler N, et al. Dissemination
13 of scientific software with Galaxy ToolShed. Genome Biol. 2014;15:403.
14

15 15. SQLite Home Page. <https://www.sqlite.org/>. Accessed 18 Nov 2016.
16

17 16. Get sequences by Ensembl ID : Galaxy Tool Shed.
18 https://toolshed.g2.bx.psu.edu/view/earlhaminst/ensembl_get_sequences/. Accessed 20 Dec
19 2016.
20

21 17. Get features by Ensembl ID : Galaxy Tool Shed.
22 https://toolshed.g2.bx.psu.edu/view/earlhaminst/ensembl_get_feature_info/. Accessed 20 Dec
23 2016.
24

25 18. Select longest CDS per gene : Galaxy Tool Shed.
26 https://toolshed.g2.bx.psu.edu/view/earlhaminst/ensembl_longest_cds_per_gene. Accessed 8
27 Mar 2017.
28

29 19. ETE species tree generator : Galaxy Tool Shed.
30 <https://toolshed.g2.bx.psu.edu/view/earlhaminst/ete/>. Accessed 20 Dec 2016.
31

32 20. GeneSeqToFamily preparation : Galaxy Tool Shed.
33 https://toolshed.g2.bx.psu.edu/view/earlhaminst/gstf_preparation/. Accessed 17 Mar 2017.
34

35 21. EMBOSS : Galaxy Tool Shed. https://toolshed.g2.bx.psu.edu/view/devteam/emboss_5/.
36 Accessed 21 Dec 2016.
37

38 22. NCBI BLAST plus : Galaxy Tool Shed.
39 https://toolshed.g2.bx.psu.edu/view/devteam/ncbi_blast_plus. Accessed 21 Dec 2016.
40

41 23. BLAST parser : Galaxy Tool Shed.
42 https://toolshed.g2.bx.psu.edu/view/earlhaminst/blast_parser/. Accessed 20 Dec 2016.
43

44 24. hcluster_sg : Galaxy Tool Shed.
45 https://toolshed.g2.bx.psu.edu/view/earlhaminst/hcluster_sg/. Accessed 20 Dec 2016.
46

47 25. hcluster_sg parser : Galaxy Tool Shed.
48 https://toolshed.g2.bx.psu.edu/view/earlhaminst/hcluster_sg_parser/. Accessed 20 Dec 2016.
49

50 26. Filter by FASTA IDs : Galaxy Tool Shed.
51 https://toolshed.g2.bx.psu.edu/view/galaxyp/filter_by_fasta_ids/. Accessed 21 Dec 2016.
52

53 27. T-Coffee : Galaxy Tool Shed. https://toolshed.g2.bx.psu.edu/view/earlhaminst/t_coffee/.
54 Accessed 20 Dec 2016.
55

56 28. TreeBeST best : Galaxy Tool Shed.
57
58
59
60
61
62
63
64
65

- 1
2
3
4 https://toolshed.g2.bx.psu.edu/view/earlhaminst/treebest_best/. Accessed 20 Dec 2016.
5
6
7 29. Gene Align and Family Aggregator (GAFA) : Galaxy Tool Shed.
8 <https://toolshed.g2.bx.psu.edu/view/earlhaminst/gafa/>. Accessed 21 Dec 2016.
9
10 30. uniprot_rest_interface : Galaxy Tool Shed.
11 https://toolshed.g2.bx.psu.edu/view/bgruening/uniprot_rest_interface/. Accessed 20 Mar 2017.
12
13 31. Yates A, Beal K, Keenan S, McLaren W, Pignatelli M, Ritchie GRS, et al. The Ensembl
14 REST API: Ensembl Data for Any Language. *Bioinformatics*. 2015;31:143–5.
15
16 32. Representational State Transfer. <http://www.peej.co.uk/articles/rest.html>. Accessed 4 Feb
17 2016.
18
19 33. Huerta-Cepas J, Serra F, Bork P. ETE 3: Reconstruction, Analysis, and Visualization of
20 Phylogenomic Data. *Mol Biol Evol*. 2016. doi:10.1093/molbev/msw046.
21
22 34. GFF3 - GMOD. <http://gmod.org/wiki/GFF3>. Accessed 4 Feb 2016.
23
24 35. JSON. <http://www.json.org>. Accessed 4 Feb 2016.
25
26 36. Rice P, Longden I, Bleasby A. EMBOSS: the European Molecular Biology Open Software
27 Suite. *Trends Genet*. 2000;16:276–7.
28
29 37. Cock PJA, Chilton JM, Grüning B, Johnson JE, Soranzo N. NCBI BLAST+ integrated into
30 Galaxy. *Gigascience*. 2015;4:39.
31
32 38. National Center for Biotechnology Information (U.S.), Camacho C. BLAST(r) Command Line
33 Applications User Manual. 2008.
34
35 39. “Newick’s 8:45” Tree Format Standard.
36 http://evolution.genetics.washington.edu/phylip/newick_doc.html. Accessed 8 Apr 2016.
37
38 40. Sequence Alignment/Map Format Specification. [http://samtools.github.io/hts-](http://samtools.github.io/hts-specs/SAMv1.pdf)
39 [specs/SAMv1.pdf](http://samtools.github.io/hts-specs/SAMv1.pdf). Accessed 20 Dec 2016.
40
41 41. TGAC. TGAC/earlham-galaxytools. GitHub. <https://github.com/TGAC/earlham-galaxytools>.
42 Accessed 21 Mar 2016.
43
44 42. Gene: BRAT1 (ENSG00000106009) - Gene tree - Homo sapiens - Ensembl genome
45 browser 87.
46 [http://dec2016.archive.ensembl.org/Homo_sapiens/Gene/Compara_Tree?g=ENSG0000010600](http://dec2016.archive.ensembl.org/Homo_sapiens/Gene/Compara_Tree?g=ENSG00000106009;r=7:2537877-2555727)
47 [9;r=7:2537877-2555727](http://dec2016.archive.ensembl.org/Homo_sapiens/Gene/Compara_Tree?g=ENSG00000106009;r=7:2537877-2555727); Accessed 23 Dec 2016.
48
49 43. Gene: INSR (ENSG00000171105) - Gene tree - Homo sapiens - Ensembl genome browser
50 87.
51 [http://dec2016.archive.ensembl.org/Homo_sapiens/Gene/Compara_Tree?g=ENSG0000017110](http://dec2016.archive.ensembl.org/Homo_sapiens/Gene/Compara_Tree?g=ENSG00000171105;r=19:7112255-7294034)
52 [5;r=19:7112255-7294034](http://dec2016.archive.ensembl.org/Homo_sapiens/Gene/Compara_Tree?g=ENSG00000171105;r=19:7112255-7294034); Accessed 23 Dec 2016.
53
54 44. Gene: MAOA (ENSG00000189221) - Gene tree - Homo sapiens - Ensembl genome
55 browser 87.
56 [http://dec2016.archive.ensembl.org/Homo_sapiens/Gene/Compara_Tree?g=ENSG0000018922](http://dec2016.archive.ensembl.org/Homo_sapiens/Gene/Compara_Tree?g=ENSG00000189221;r=X:43654907-43746824)
57 [1;r=X:43654907-43746824](http://dec2016.archive.ensembl.org/Homo_sapiens/Gene/Compara_Tree?g=ENSG00000189221;r=X:43654907-43746824); Accessed 23 Dec 2016.
58
59
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

45. Gene: MAOB (ENSG00000069535) - Gene tree - Homo sapiens - Ensembl genome browser 87.
http://dec2016.archive.ensembl.org/Homo_sapiens/Gene/Compara_Tree?g=ENSG00000069535;r=X:43766611-43882447; Accessed 23 Dec 2016.

46. Thanki AS, Soranzo N, Haerty W, Davey R. GeneSeqToFamily.zip. 2017.
doi:10.6084/m9.figshare.4484141.v3.

47. Mi H, Lazareva-Ulitsky B, Loo R, Kejariwal A, Vandergriff J, Rabkin S, et al. The PANTHER database of protein families, subfamilies, functions and pathways. *Nucleic Acids Res.* 2005;33 Database issue:D284–8.





Sequence1: NLYIQWLKDGGPSSGRPPPS

Sequence2: NLYIQWLKDQGPSSGRPPPS

Sequence3: GDAYAQWLADGGPSSGRPPPSG

Sequence1: -NLYIQWLKDGGPSSGRPPP-S

Sequence2: -NLYIQWLKDQGPSSGRPPP-S

Sequence3: GDAYAQWLADGGPSSGRPPPSG

CIGAR1: D19MDM

CIGAR2: D19MDM

CIGAR3: 22M











