

Lawrence Berkeley National Laboratory

Recent Work

Title

Metazome

Permalink

<https://escholarship.org/uc/item/5qx2z2r9>

Authors

Dirks, Bill
Goodstein, David
Hellsten, Uffe
et al.

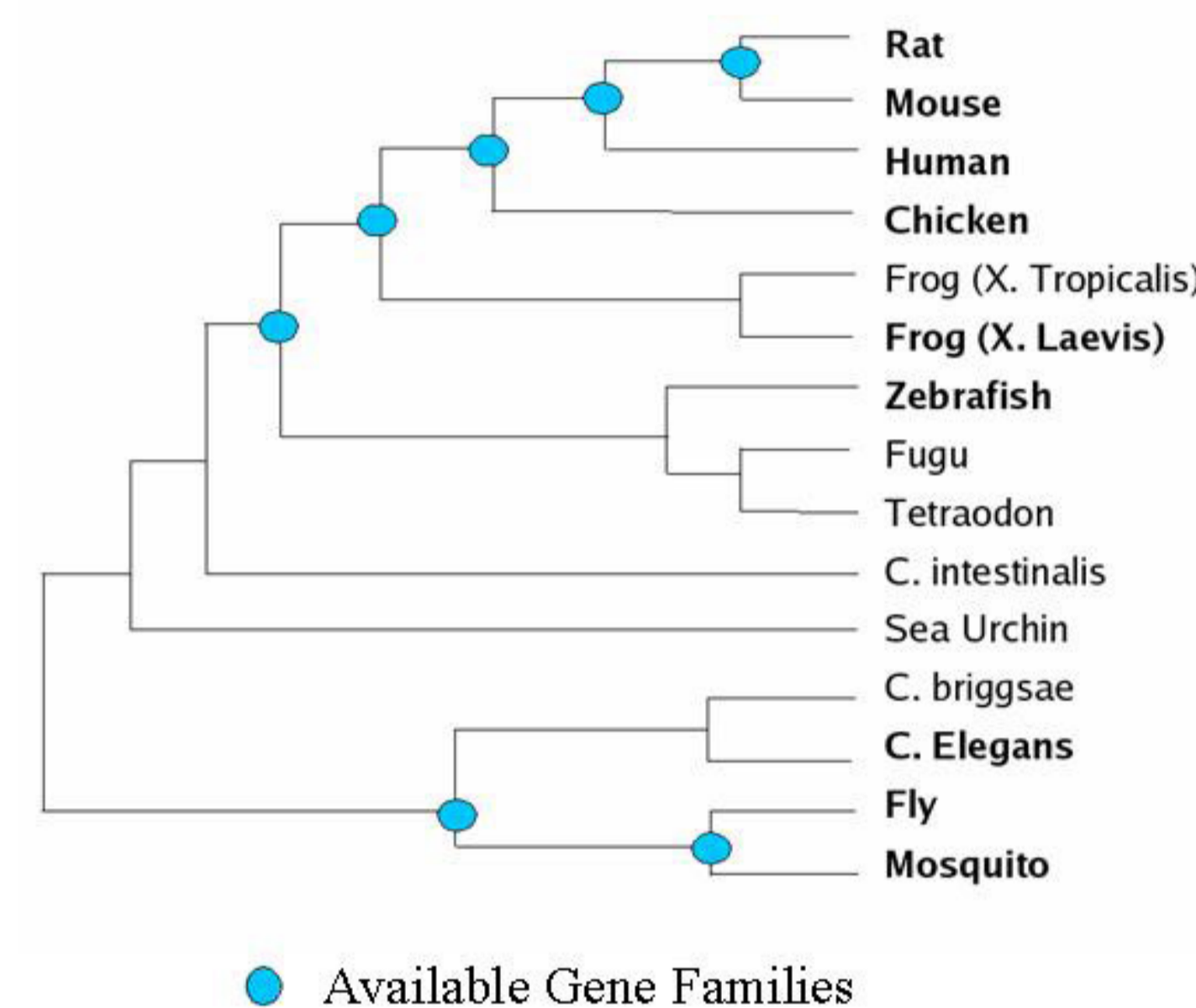
Publication Date

2005-05-10

Abstract

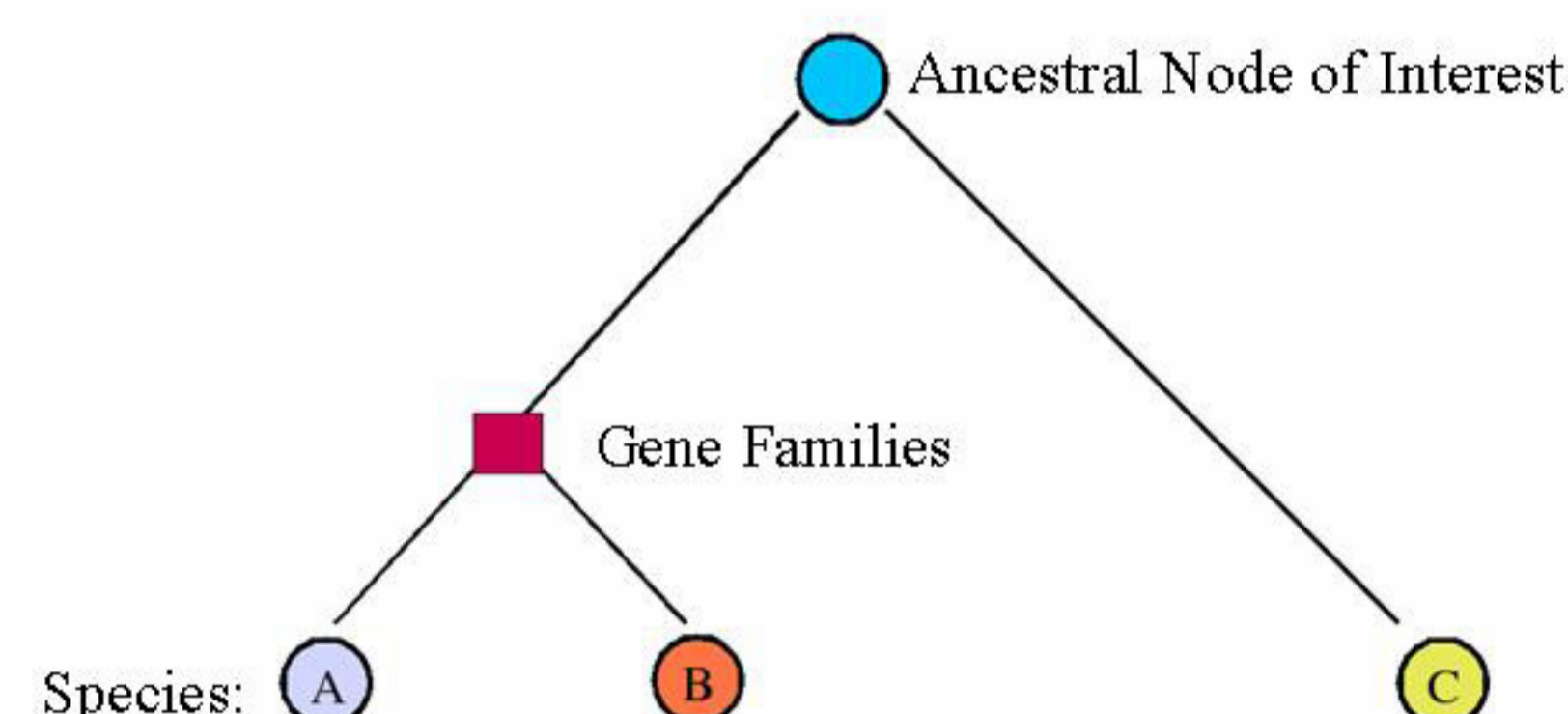
Metazome is a database and graphical user interface enabling comparative genomic studies within the Metazoa. This system allows uniform access to a high quality set of protein-coding genes from all annotated metazoan genomes. Each gene has been assigned with PFAM, KOG, and PANTHER annotations, and publicly available annotations from RefSeq, SwissProt, Ensembl, and JGI are hyper-linked and searchable. In addition, these genes are organized into orthologous gene families each representing the modern descendents of the ancestral gene at a key phylogenetic node such as the ancestral mammal, tetrapod, vertebrate, etc.

Current Genomes and Gene Families

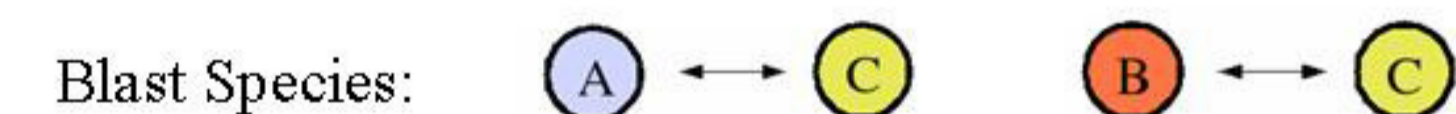


Gene Family Creation Algorithm

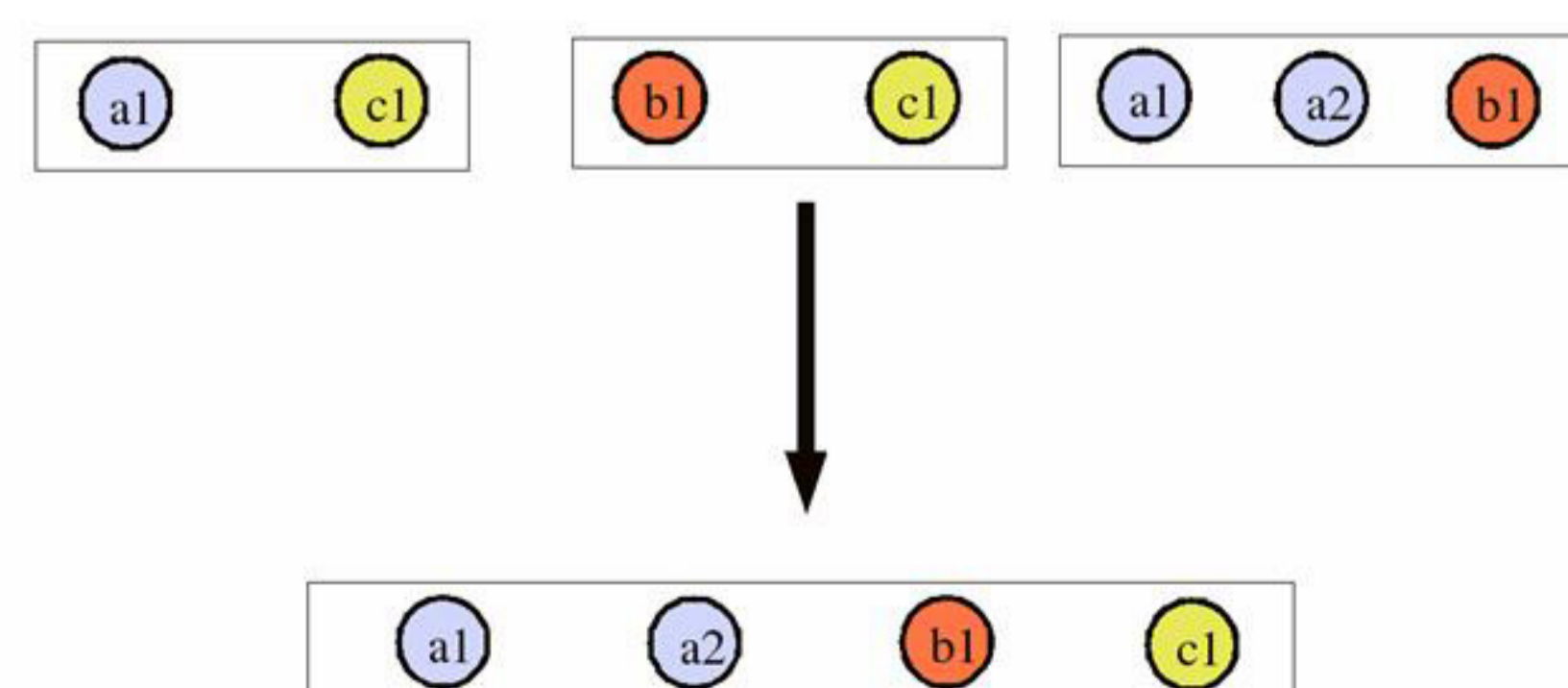
Step 1 - Making a Backbone of high quality putative orthologs



First, we find mutual best hits across the ancestral node of interest. This is done using BLAST, stringing together the highest scoring set of nonoverlapping hsp. This gives us pseudo-global alignments.

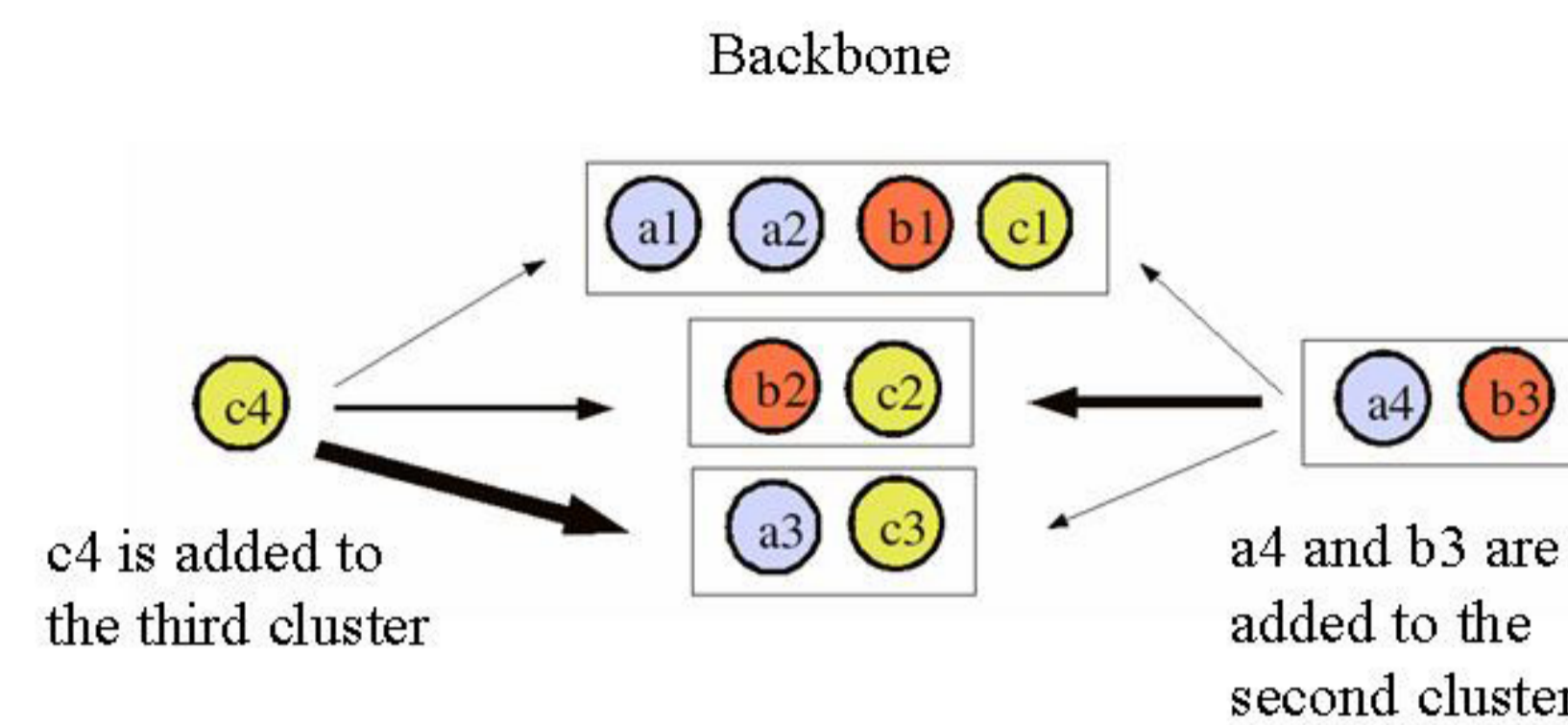


Next, we combine all the pairwise best hits and previous families if they share common genes.



This becomes a backbone element of orthologous genes.

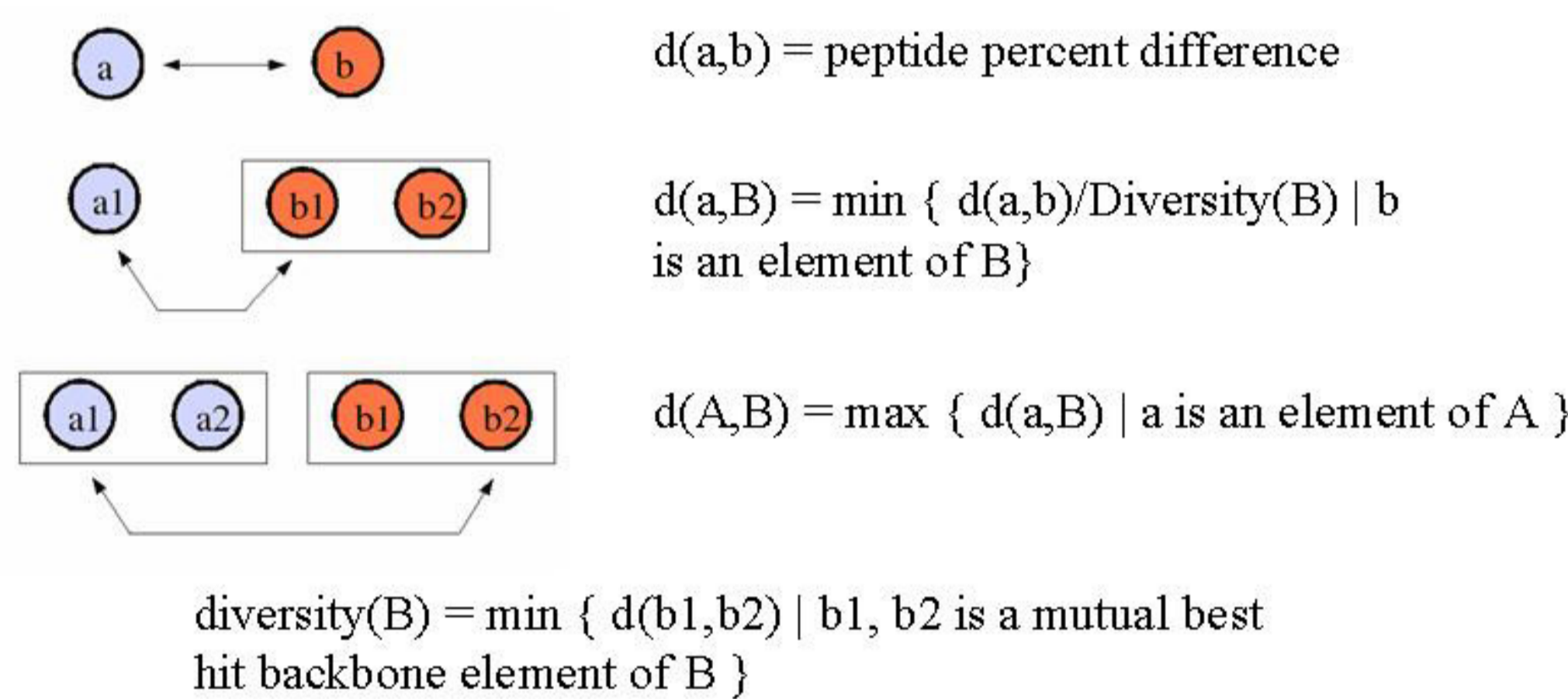
Step 2 - The backbones are grown out into gene families



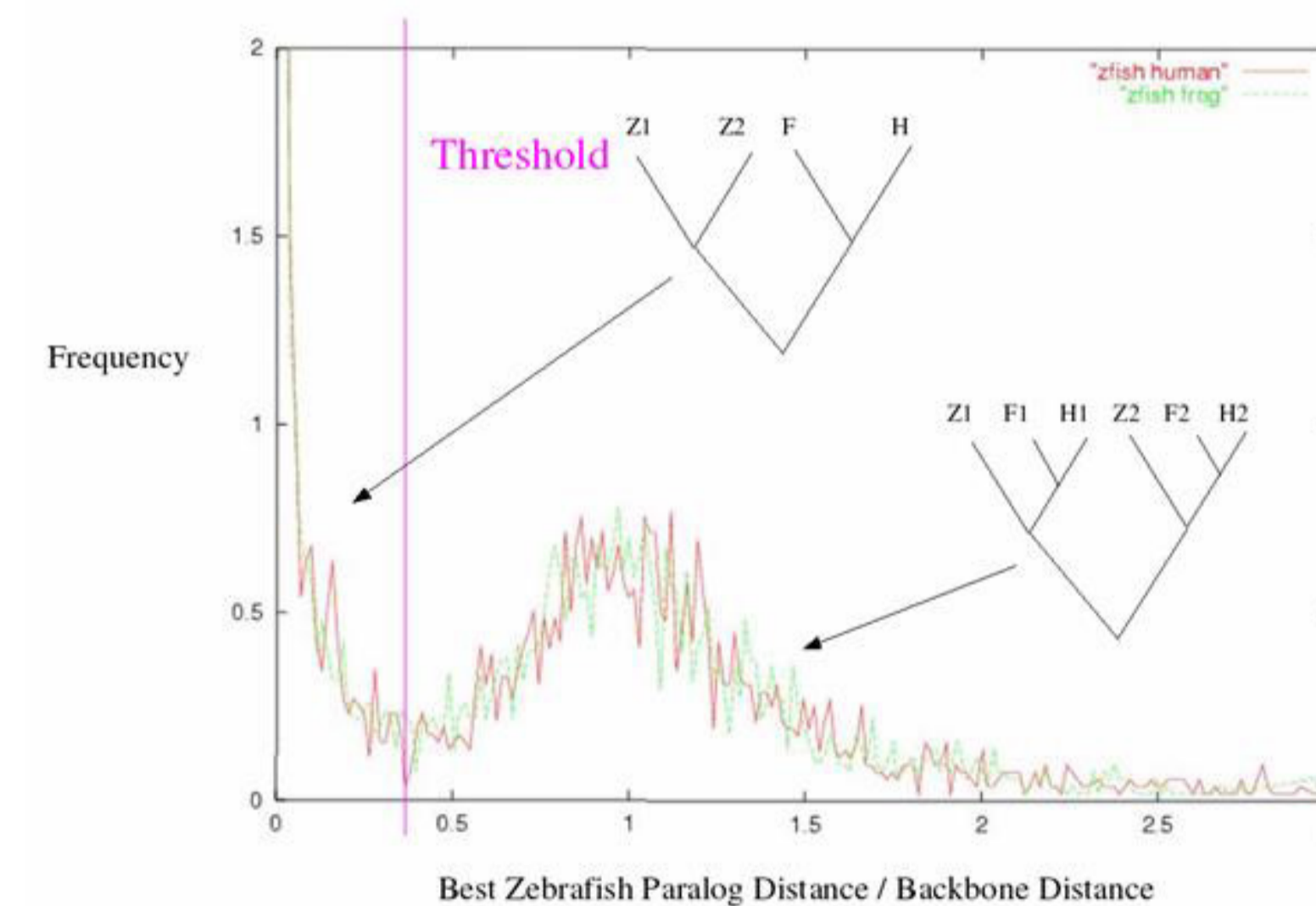
Genes and previously made gene families are added to a backbone element if we believe they are paralogs to this backbone element at this depth of evolution. We assert a gene or gene family is paralogous to its closest backbone element if the distance is better than a pre-computed empirical threshold (see below).

Distances

Our distances respect the variable mutability of different gene families.



Empirical Paralogy Threshold for Constructing Gene Families



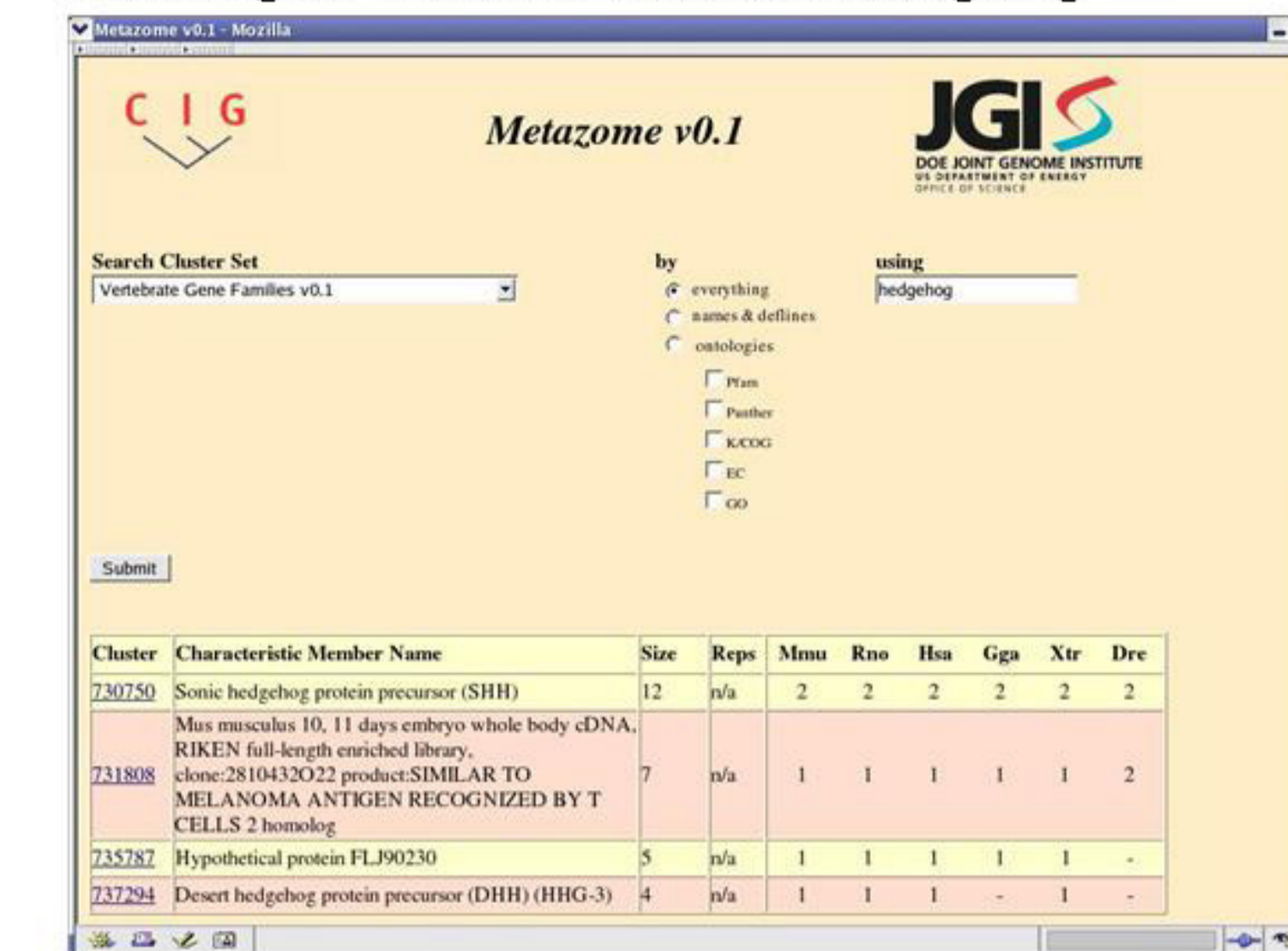
When choosing a threshold for constructing gene families we look at triplets of best hit gene pairs within a species and between species. This lets us distinguish old paralogs from new ones.

Gene Family Statistics

Evolutionary Depth	N. of Organisms	N. of Gene Families	Average Standard Size	Family Size Deviation	N. of Genes in Families	N. of Genes Total
Rodent	2	16569	2.051	0.003	33893	47127
Mammal	3	15645	3.082	0.009	48225	69216
Amniote	4	11369	4.398	0.003	49435	86611
Tetrapod	5	10787	5.539	0.034	76570	143484
Vertebrate	6	10684	7.167	0.050	76570	143484
Fly-Mosquito	2	7297	2.088	0.001	15235	27418
Ecdysozoa	3	4926	3.095	0.016	15247	47238

Screenshots

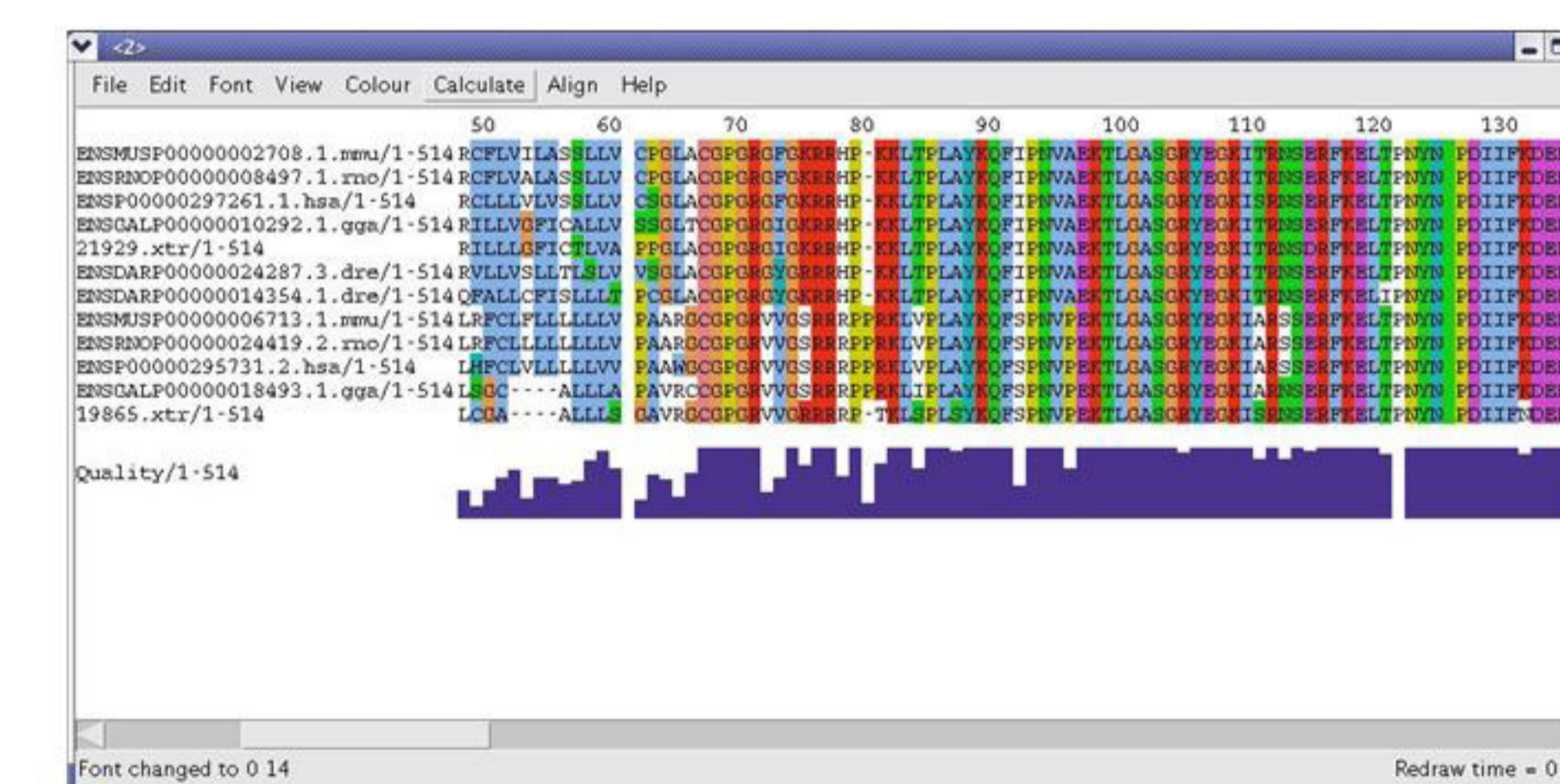
Searching the vertebrate families for hedgehog



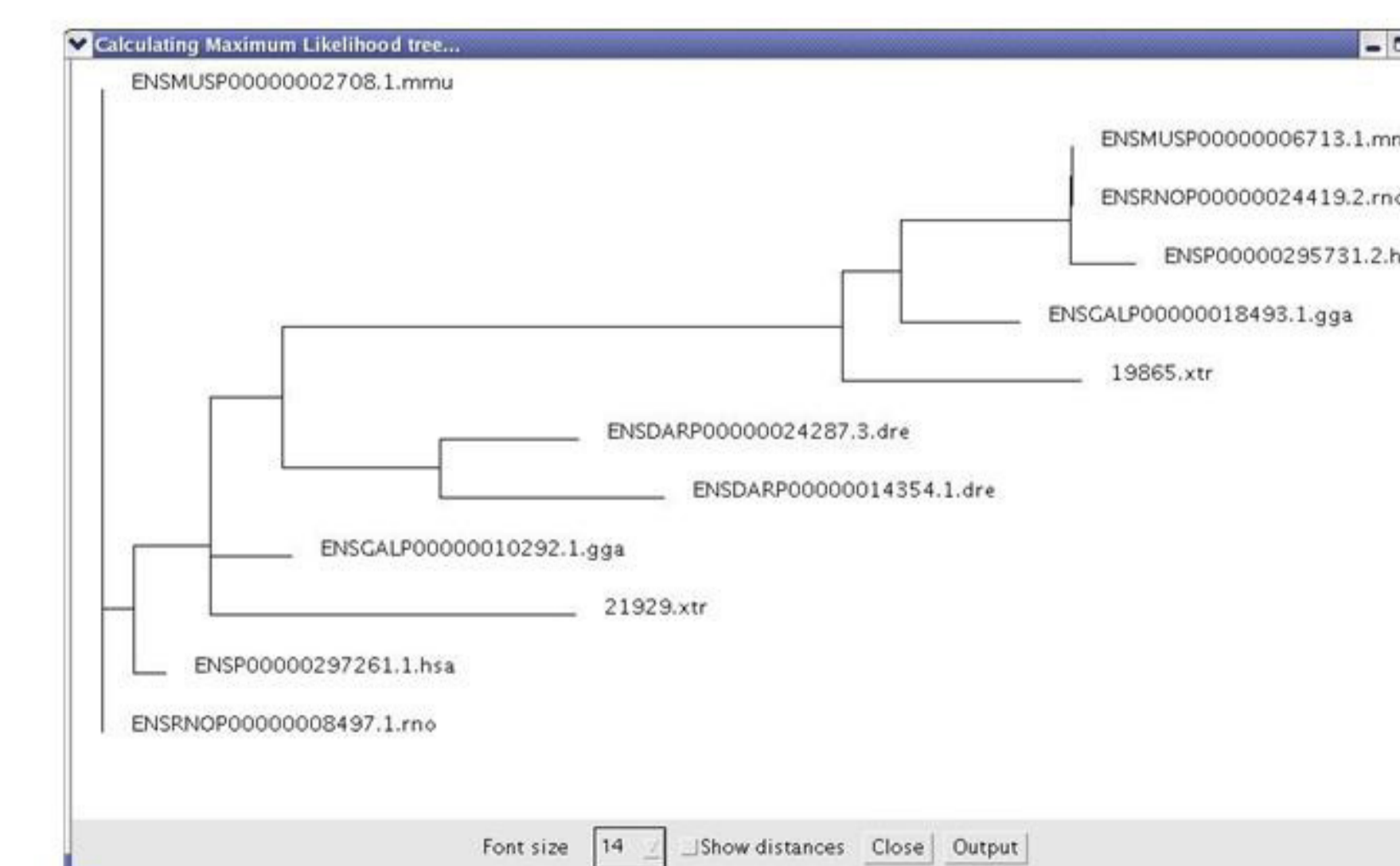
Cluster Detail Information



Examining a pre-computed multiple sequence alignment with JalView



Retrieving a tree precomputed using Tree-Puzzle



Coming Features

- 1) additional organisms and gene families
- 2) Searching with sequence (BLAST)
- 2) Retrieving nearby genomic sequence to facilitate cis-regulation studies
- 3) Manual curation of family names for core families at key nodes
- 4) Incorporation of expression data and protein-protein interactions