1

2

3

4

5

6

7

8

9

10    **Genome-wide protein phylogenies for four African cichlid species**

11

12

13

14    Ajay Ramakrishnan Varadarajan, Rohini Mopuri, J. Todd Streelman, and Patrick T. McGrath[1]

15    [1]Department of Biological Sciences, Georgia Institute of Technology, Atlanta, GA 30332, USA

16

1

1   **ABSTRACT**

2   Background

3   The thousands of species of closely related cichlid fishes in the great lakes of East Africa are a

4   powerful model for understanding speciation and the genetic basis of trait variation. Recently,

5   the genomes of five species of African cichlids representing five distinct lineages were

6   sequenced and used to predict protein products at a genome-wide level. Here we characterize

7   the evolutionary relationship of each cichlid protein to previously sequenced animal species.

8   Results

9   We used the Treefam database, a set of preexisting protein phylogenies built using 109

10   previously sequenced genomes, to identify Treefam families for each protein annotated from

11   four cichlid species: *Metriaclima zebra*, *Astatotilapia burtoni*, *Pundamilia nyererei* and

12   *Neolamporologus brichardi*. For each of these Treefam families, we built new protein

13   phylogenies containing each of the cichlid protein hits. Using these new phylogenies we

14   identified the evolutionary relationship of each cichlid protein to its nearest human and zebrafish

15   protein. This data is available either through download or through a webserver we have

16   implemented.

17   Conclusion

18   These phylogenies will be useful for any cichlid researchers trying to predict biological and

19   protein function for a given cichlid gene, understanding the evolutionary history of a given cichlid

20   gene, identifying recently duplicated cichlid genes, or performing genome-wide analysis in

21   cichlids that relies on using databases generated from other species.

## 1   BACKGROUND

2   The rapid decrease in sequencing costs and the development of broadly applicable genetic

3   tools like TALENs and CRISPR/Cas9 has facilitated the development of a large number of new

4   species as model organisms (Joung and Sander 2013, Doudna and Charpentier 2014, Hsu,

5   Lander et al. 2014, Nemudryi, Valetdinova et al. 2014). For evolutionary biologists, this has

6   been especially fruitful – species with unique evolutionary traits can now be used as model

7   organisms to identify and understand the underlying genetic and cellular mechanisms

8   responsible for trait changes (Goldstein and King 2016). For example, threespine sticklebacks

9   have long fascinated evolutionary biologists for their coexisting phenotypically divergent forms

10   including freshwater/anadromous pairs (Hagen 1967, Mcphail 1969, McKinnon and Rundle

11   2002). Freshwater lakes created after the retreat of Pleistocene glaciers have been populated

12   by marine sticklebacks, evolving repeated changes in a number of traits. These adaptations

13   include morphological changes to body shape, pigmentation changes, salt handling, and

14   reproductive related behaviors (Bell and Foster 1994, McKinnon and Rundle 2002). A

15   combination of quantitative genetics and resequencing of individuals isolated from freshwater

16   and saltwater habitats identified a large number of loci putatively responsible for evolution of

17   marine-freshwater ecotypes (Colosimo, Peichel et al. 2004, Chan, Marks et al. 2010, Jones,

18   Grabherr et al. 2012). An important conclusion from this research, and a number of other

19   individual examples (Martin and Orgogozo 2013), is that despite the large number of genes that

20   control a trait, natural selection can act in predictable ways, isolating genetic changes in

21   preferred genes in response to specific environmental shifts. An important goal now is to identify

22   additional examples of repeated evolution, and understand why particular genes are repeatedly

23   selected.

24   Cichlid fishes offer an attractive avenue for this type of research. Cichlids are well-known for

25   their adaptive radiations in the Great Lakes of East Africa. The three largest radiations in Lakes

26   Victoria, Lake Malawi, and Lake Tanganyika have generated between 250 – 500 species per

27   lake in a period of time that ranges from 100,000 to 12 million years (Kocher 2004, Brawand,

28   Wagner et al. 2014). These radiations resulted in exceptional phenotypic diversity in behavior,

29   neurodevelopment, body shape, sexual traits, and ecological specialization. However, due to

30   the speed of evolution, nucleotide diversity between these species is on the order of nucleotide

31   diversity within the human population (Loh, Bezault et al. 2013, Brawand, Wagner et al. 2014).

32   Further, genetic barriers have not formed in this short period, allowing for genetics -

33   phenotypically-divergent species can still interbreed. These peculiarities of the cichlid family

1  make genomics and quantitative genetics approaches particularly attractive. Genes responsible

2  for phenotypic diversity can be identified using quantitative mapping approaches in progeny of

3  intercrossed species, association mapping in outbred animals, or tissue-specific transcriptomics

4  in behaving animals. To facilitate these approaches, high-quality genomes for five cichlid fishes

5  were generated (Brawand, Wagner et al. 2014). It is anticipated that genetic variants and genes

6  responsible for a variety of interesting trait differences will be identified in the coming years.

7  Due to the difficulty of experimental study of cichlids in the laboratory, assignment of molecular

8  and biological function to genes relies almost exclusively on homology to proteins characterized

9  biochemically, or in model organisms such as *Caenorhabditis elegans, Drosophila*

10  *melanogaster, Danio rerio,* or *Mus musculus*. Homologous proteins share a common

11  evolutionary ancestry (Fitch 1970), suggesting shared biochemical and/or biological role,

12  justifying the use of homology to assign function to genes identified in cichlid fish. Proteins with

13  shared homology can be characterized as orthologs (which diverged from a common ancestor

14  due to speciation) or paralogs (which diverged from a common ancestor due to a gene

15  duplication event). In general, orthologs are expected to retain similar (if not identical) function

16  with each other. Paralogs are expected to acquire novel function and/or biological roles. For

17  cichlids, paralogs are thought to be especially relevant to their evolution - the cichlid lineage has

18  undergone an increased rate of gene duplication, suggesting that these novel genes could

19  serve important roles in the cichlid's adaptive radiations (Lynch and Conery 2000, Brawand,

20  Wagner et al. 2014). Cichlids also belong to the teleost infraclass of fish, whose ancestors have

21  undergone a genome-wide duplication event resulting in the duplication of a large number of

22  genes (Taylor, Braasch et al. 2003). Gene duplication can allow resolution of adaptive conflict

23  by allowing a bifunctional ancestral gene to resolve into two specialized genes (Lynch and Force

24  2000). These gene duplicates have been proposed to play a role in the evolutionary success of

25  the teleost fish, which make up ~96% of all fish. Phylogenetic relationships could potentially be

26  used to identify the cichlid genes that have undergone subfunctionalization. For all of these

27  reasons, it would be helpful to place each cichlid protein into a phylogeny to aid in predicting the

28  gene function for a given cichlid gene.

29  In this report, we utilized the TreeFam database of protein phylogenies to create protein

30  phylogenies for all completely sequenced cichlid genomes. We analyzed these phylogenies to

31  determine evolutionary relationships for each of these cichlid genes. This data is available for

32  download or searching on a web server, and should be useful to any researchers studying

33  cichlid fish.

4

**METHODS**

**Overview**

We employed a phylogeny-based approach to study the function and evolution of <genes of interest> taken from four East African cichlid species. Our aim was to assign each cichlid gene to a pre-defined gene family to identify homologous proteins and their evolutionary relationship. To accomplish this, we used TreeFam, a database of phylogenetic trees drawn from 109 animal genomes (Li, Coghlan et al. 2006, Ruan, Li et al. 2008, Schreiber, Patricio et al. 2014). A webserver implementing the Treefam pipeline is provided (www.treefam.org) to add new proteins of interest to existing TreeFam trees. We implemented this pipeline locally to perform this on a genomic basis.

**Datasets and TreeFam analysis**

Protein coding sequences and annotation files for four cichlid species, *A. burtoni*, *M. zebra*, *N. brichardi*, and *P. nyererei,* were obtained from the supplemental dataset from the genome sequencing paper (Brawand, Wagner et al. 2014). An improved genome for *M. zebra* was also recently published; protein coding sequence and genome annotation files from this paper were downloaded from NCBI (Conte and Kocher 2015). Annotation files were parsed using custom Python scripts and used to identify the longest protein isoform and amino acid sequence for each gene. This was done to limit the phylogeny to one representative protein isoform for each gene. To assign each of these proteins to a single TreeFam family, we utilized the ıêÉÉÑ~ãëÅ~åKéä script provided as part of the TreeFam API (Schreiber, Patricio et al. 2014). This script uses the program HMMER to identify matches using hidden Markov model profiles generated for each of the TreeFam families (Eddy 1998). After this had run on all of the proteins, we collected all of the protein sequences that best matched a given TreeFam to add these to the preexisting phylogeny. Multiple sequence alignments and phylogenies for each TreeFam were retrieved from a locally cloned SQL database with API utilities provided by TreeFam. We used MAFFT (version 7.221) to add the new cichlid proteins to the retrieved multiple sequence alignment using the - ~ÇÇ, - êÉçêÇÉê, and - ~åóëóóãÄÇä options (Katoh, Misawa et al. 2002, Katoh and Standley 2013). The aligned output was then used to add the new proteins to the retrieved phylogeny file using RaXML (version 8.1.15) using the GAMMA model for rate heterogeneity with the WAG substitution matrix (Stamatakis, Ludwig et al. 2005, Stamatakis 2006).

**Identification of closest relationships to human and zebrafish proteins**

1    For each cichlid protein, we used custom Python scripts to identify the closest human and

2    zebrafish protein using the phylogenetic tree produced by RAxML. The structures of each tree

3    were analyzed using the ETE toolkit, which provides a Python framework for analysis and

4    visualization of protein trees (Huerta-Cepas, Serra et al. 2016). Trees were rooted using a

5    midpoint outgroup method implemented by the ÖÉí|ãáÇéçáåí|çìí Öèçìé function. To find the

6    closest human protein and its evolutionary relationship with a cichlid protein of interest, the trees

7    were then traversed to identify the smallest subtree containing the cichlid protein and one or

8    more human protein. If such a subtree could not be found (i.e. there was no human protein in

9    the phylogeny), the relationship was defined as kçeçãçäçÖ. If the subtree contained a single

10    human protein and a single cichlid protein from the cichlid species, the relationship was defined

11    as lãíÜçäçÖ. If the subtree contained a single human protein and exactly two cichlid proteins

12    from the cichlid species, this relationship was defined as a e~äÑãíÜçäçÖ with the human protein.

13    Finally, if the subtree contained multiple human proteins, or more than two cichlid proteins from

14    the cichlid species, this relationship was defined as a m~ê~äçÖ. The closest human protein was

15    identified using the shortest branch length. To convert the Ensembl protein ID's of the human

16    proteins to HGNC identifiers (Gray, Yates et al. 2015), we downloaded mapping data from

17    Ensembl BioMart (Aken, Ayling et al. 2016). An essentially identical process was also performed

18    between all cichlid proteins with zebrafish proteins. An excel spreadsheet (one per species) was

19    then created for each cichlid gene for this information.

20    PDFs of the resulting phylogenies were rendered using the ETE toolkit. A full size version of

21    each TreeFam phylogeny was created using all species. In addition, a smaller PDF was created

22    from a pruned tree containing a limited number of well-characterized species (human (*H.*

23    *sapiens*), mouse (*M. musculus*), zebrafish (*D. rerio*), fruit fly (*D. melanogaster*), and nematode

24    (*C. elegans*)), the closely related Nile tilapia (*O. niloticus*), and the four new cichlid species.

25    **RESULTS AND DISCUSSION**

26    **Identification of human and zebrafish relationships for each cichlid gene**

27    The cichlids species of East Africa have become a popular genomic model to understand the

28    evolution of a number of traits, including differences in morphology, coloration and behavior. To

29    broaden our understanding of the function and evolutionary history of the genes that are

30    encoded in the genomes of four recently-sequenced cichlid species, we performed phylogenetic

31    analysis using the previously published TreeFam pipeline to add the new cichlid proteins to

32    preexisting protein phylogenies generated from a large number of animal species (**Figure 1**).

33    The most current version of the TreeFam database (Schreiber, Patricio et al. 2014), which

1   contains 15,736 phylogenetic trees generated from 109 animal genomes covering ~2.2 million

2   sequences, can be used to study evolutionary relationships between homologous proteins.

3   While this database already includes the African cichlid *O. niloticus* (Nile tilapia), it does not

4   contain four recently sequenced African cichlids: *M. zebra* from Lake Malawi, *P. nyererei* from

5   Lake Victoria, *N. brichardi* from Lake Tanganyika, and *A. burtoni* found in a variety of African

6   lakes and rivers. For all four cichlid species, the majority of cichlid genes, 82.2% – 84.7%,

7   contained a hit to a preexisting TreeFam family (**Figure 2**). Using the resulting phylogenies, we

8   identified the closest human and zebrafish gene along with the evolutionary relationship to the

9   cichlid. These included traditional evolutionary relationships (Ortholog and Paralog) and also a

10  novel evolutionary definition we call HalfOrtholog, to account for the large number of cichlid

11  genes that duplicated in the ancestral teleost lineage and are retained in the extant species.

12  **Data accessibility**

13  This data is intended as a resource for the cichlid community. We have provided access to this

14  data in three ways. 1. Two PDF files for each TreeFam were generated for the purposes of

15  human inspection. One PDF contains a phylogeny for a TreeFam from a limited number of

16  species: humans, four well-characterized model organisms (*C. elegans, D. melanogaster, D.*

17  *rerio,* and *M. musculus*), Nile tilapia (*O. niloticus*), and the four recently studied cichlid species.

18  The second PDF contains a phylogeny of all 108 species used in the analysis. While the second

19  phylogeny is the most complete, it is difficult to analyze due to the large number of species. This

20  data is hosted on a web server (`http://cichlids.biosci.gatech.edu/`) and can be

21  searched using cichlid gene names, TreeFam IDs, or human and zebrafish names. 2. Excel files

22  for each cichlid species that contain each gene, its best hit to a human and zebra fish gene, and

23  its evolutionary relationship to that gene. We anticipate this data will be useful for genomic scale

24  analysis. For example, the excel file can be loaded into scripts to automatically map cichlid

25  genes to human or zebrafish homologs. This could be useful for the purposes of pathway

26  analysis (such as gene ontology), which often are limited to human genes. 3. Finally,

27  phylogenies of each TreeFam are available for download in enhanced Newick tree format.

28  These will be useful for any researchers interested in automated analysis of the phylogenies for

29  the purpose of enhancing the evolutionary relationships that we have reported here. For

30  example, researchers could use this dataset to identify genes whose protein phylogenies

31  contradict the species phylogenies.

32  **Example phylogeny generated from a tree containing members of the TGFβ superfamily**

1    To illustrate these evolutionary relationships as well as common issues users should be aware
2    of in using these trees, we have included two figures of new phylogenies generated in this
3    analysis. **Figure 3** shows a subtree of TF351789, which includes members of the TGFβ-
4    superfamily of proteins including BMP2 and BMP4. These proteins are ubiquitous throughout
5    metazoans, and control proliferation and differentiation of cells throughout development
6    (Salazar, Gamer et al. 2016). This tree includes both ortholog and paralog relationships. For
7    example, the subtree indicated by **a** in **Figure 3** shows ortholog relationships between the
8    cichlid proteins and human BMP4. These genes likely play similar biological roles in cichlids.
9    Similarly, subtree **b** contains cichlid orthologs to human BMP2 (with the exception of *M. zebra*,
10   which will be discussed below), suggesting these genes play similar biological roles as the
11   orthologs play in other species. There is also a cichlid-specific set of paralogs to BMP2 and
12   BMP4 not present in *D. rerio* (subtree **c**) suggesting that there was a duplication of BMP2 or
13   BMP4 in a recent common ancestor of all cichlid species following separation from the zebrafish
14   lineage. It is not obvious from the phylogeny what biological role these genes might play. This
15   clade of genes is potentially of interest to cichlid biologists, as they could play a role in the
16   extensive morphological diversity observed among cichlid species. However, analysis of the full
17   tree indicated that this clade contains genes from a large number of additional teleost fish along
18   with a coelacanth fish (*L. chalumnae*) and an anole lizard (*A. carolinensis*) (**Figure S1**). Further,
19   blasting the protein sequence encoded by the ab.gene.s112.4 from *A. burtoni* to the *D. rerio*
20   genome identified a match to a known protein annotated as BMP16 (Feiner, Begemann et al.
21   2009). BMP16 does not appear to be present in the Treefam database, which explains why it
22   was not present in the phylogeny. This set of BMP2/BMP4 paralogs thus seems to be a
23   duplication that occurred in an ancient vertebrate ancestor of these fish (preceding the teleost
24   ancestor) and lost in most tetrapod lineages as proposed by Marques et al (Marques,
25   Fernandez et al. 2016).

26   We observed a similar issue in the phylogeny surrounding the human IRX1 gene (TF319371)
27   (**Figure S2**). The Treefam phylogeny suggests that *D. rerio* contained a single ortholog to this
28   gene while each of the cichlid species contained two copies of this gene. However, previous
29   publications demonstrate that there are also two versions of IRX1 in *D. rerio* (called *irx1a* and
30   *irx1b*) (Dildrop and Ruther 2004, Feijoo, Manzanares et al. 2004). Inspection of the Treefam
31   data indicates that *irx1a* isn't present in the starting dataset. These examples illustrate a
32   common issue to most genomic analysis. Since Treefam relies on genomic-scale predictions,
33   there are likely errors within the resulting phylogenies. Users would do well to manually verify or
34   repeat any of these phylogenies for genes they are especially interested in.

8

1   We also were curious about the lack of a clear ortholog to BMP2 in *M. zebra* (**Figure 3**). It

2   seemed unlikely that this species could lose this protein entirely due to its essential function in

3   bone development. We were able to track down this discrepancy to an error in the annotation

4   file for *M. zebra*. Through blastp, we were able to identify mz.gene.s5.238 as a gene containing

5   a strong match to BMP2. mz.gene.s5.238, however, was assigned to the TF314677 family, and

6   predicted to be an ortholog to the human protein FERMT1. When we invested the protein

7   sequence more closely, it became clear that mz.gene.s5.238 appeared to contain a fusion of

8   two genes: an ortholog to BMP2 and an ortholog to FERMT1. Due to the longer length of

9   FERMT1, mz.gene.s5.238 was assigned to the TreeFam containing FERMT1. This is unlikely to

10  represent a real gene fusion, and the improved version of the *M. zebra* genome predicts

11  separate gene products consistent with other species(Conte and Kocher 2015). We observed a

12  similar potential error with the PTGFR prostaglandin receptor. An ortholog of PTGFR has

13  recently been shown to control female reproductive behaviors in the cichlid *A. burton(Juntti,*

14  *Hilliard et al. 2016)i*, however, the Treefam containing the human PTGFR gene (TF324982), did

15  not contain an ortholog of this gene in *A. burtoni*. Again, this seems to be due to an error in

16  annotation incorrectly predicting a fusion between two genes. The best blastp match

17  ab.gene.s495.12 contains a fusion between two genes, an ortholog to PTGFR and an ortholog

18  to the ZFYVE9. Due to the longer length of the ZFYVE9 protein, the ab.gene.s495.12 gene is

19  assigned to the Treefam containing the human ZFYVE9. Again, this is unlikely to represent a

20  real fusion, and it since has been corrected in new annotations. These two examples illustrate

21  how errors in the gene annotation can lead to incorrect phylogenies.

22  **Example phylogeny generated from a tree containing arginine vasopressin receptors**

23  **Figure 4** shows a subtree of the phylogeny for TF106499, which contains a number of receptors

24  for the arginine vasopressin and oxytocin neuropeptides that are thought to play a role in social

25  behavior and sexual motivation (Hammock, Lim et al. 2005, Insel 2010). We have limited this

26  phylogeny to the clade containing the AVPR1A and AVPR1B human proteins. The clade

27  indicated by **a** demonstrates the HalfOrtholog relationship (**Figure 4**). All of the sequenced

28  cichlid species (along with zebrafish and other teleost fish) contain two genes that fall within this

29  clade. This phylogeny suggests that the function of the ancestral AVPR1A gene bifurcated into

30  two genes in an ancestor to the teleost lineage. While the phylogeny suggests that both of these

31  receptors should retain a molecular role in arginine vasopressin/oxytocin signaling, the

32  biological function of AVPR1A should not be assigned to either of the two genes in each cichlid

33  species. Rather, experiment will be necessary to parse out the biological function of each of

1    these two half orthologs. A recent paper characterizing the expression pattern of these two

2    receptors in zebrafish demonstrated that these two genes are expressed in similar but non-

3    overlapping cell types(Iwasaki, Taguchi et al. 2013). This phylogeny also contains the human

4    AVPR1B protein. While mouse contains a clear ortholog to this gene, none of the cichlid species

5    nor zebrafish contain an ortholog to this gene. Analysis of the full phylogeny suggests that

6    AVPR1B was lost in the teleost fish completely. Thus, the phylogeny indicates that the biological

7    functions assigned to AVPR1B through the study of mouse and other mammals should not be

8    directly assigned to any of the cichlid homologs without experimental study.

9    **CONCLUSION**

10    This study reports a set of protein phylogenies generated for four recently sequenced African

11    cichlids. We hope that these phylogenies will be useful for cichlid researchers for the purpose of

12    inferring biological and molecular function of cichlid genes.

13    **ACKNOWLEDGEMENTS**

17    **FIGURE LEGENDS**

18    **Figure 1.** Pipeline for adding cichlid proteins to preexisting Treefam phylogenies.

19    **Figure 2.** Summary of the human evolutionary relationships found in each species. HalfOrtholog

20    is a non-standard relationship indicating a gene potentially duplicated and retained in an ancient

21    teleost ancestor.

22    **Figure 3**. Subtree from the TF351789 family from a limited number of cichlid species and well-

23    studied model organisms. This family contains a number of BMP growth factors belonging to the

24    transforming growth factor beta family. Letters indicate additional subtrees discussed in the text.

25    **Figure 4**. Subtree from the TF106499 family from a limited number of cichlid species and well-

26    studied model organisms. This family contains a number of G-protein receptors for the arginine

27    vasopressin and oxytocin nonapeptide hormones. Letters indicate additional subtrees discussed

28    in the text.

29    **Figure S1**. Full tree for the TF351789 family from all 109 species included in the Treefam

30    database. This family contains a number of BMP growth factors belonging to the transforming

31    growth factor beta family.

1  **Figure S2**. Subtree from the TF319371 family from a limited number of cichlid species and well-

2  studied model organisms. This family contains a number of Iroquois-family of homeodomain

3  transcription factors involved in patterning and other development processes.

4  **REFERENCES**

5  Aken, B. L., S. Ayling, D. Barrell, L. Clarke, V. Curwen, S. Fairley, J. Fernandez Banet, K. Billis,

6  C. Garcia Giron, T. Hourlier, K. Howe, A. Kahari, F. Kokocinski, F. J. Martin, D. N. Murphy, R.

7  Nag, M. Ruffier, M. Schuster, Y. A. Tang, J. H. Vogel, S. White, A. Zadissa, P. Flicek and S. M.

8  Searle (2016). "The Ensembl gene annotation system." Database (Oxford) **2016**.

9  Bell, M. A. and S. A. Foster (1994). The evolutionary biology of the threespine stickleback.

10  Oxford ; New York, Oxford University Press.

11  Brawand, D., C. E. Wagner, Y. I. Li, M. Malinsky, I. Keller, S. Fan, O. Simakov, A. Y. Ng, Z. W.

12  Lim, E. Bezault, J. Turner-Maier, J. Johnson, R. Alcazar, H. J. Noh, P. Russell, B. Aken, J. Alfoldi,

13  C. Amemiya, N. Azzouzi, J. F. Baroiller, F. Barloy-Hubler, A. Berlin, R. Bloomquist, K. L.

14  Carleton, M. A. Conte, H. D'Cotta, O. Eshel, L. Gaffney, F. Galibert, H. F. Gante, S. Gnerre, L.

15  Greuter, R. Guyon, N. S. Haddad, W. Haerty, R. M. Harris, H. A. Hofmann, T. Hourlier, G.

16  Hulata, D. B. Jaffe, M. Lara, A. P. Lee, I. MacCallum, S. Mwaiko, M. Nikaido, H. Nishihara, C.

17  Ozouf-Costaz, D. J. Penman, D. Przybylski, M. Rakotomanga, S. C. Renn, F. J. Ribeiro, M. Ron,

18  W. Salzburger, L. Sanchez-Pulido, M. E. Santos, S. Searle, T. Sharpe, R. Swofford, F. J. Tan, L.

19  Williams, S. Young, S. Yin, N. Okada, T. D. Kocher, E. A. Miska, E. S. Lander, B. Venkatesh, R.

20  D. Fernald, A. Meyer, C. P. Ponting, J. T. Streelman, K. Lindblad-Toh, O. Seehausen and F. Di

21  Palma (2014). "The genomic substrate for adaptive radiation in African cichlid fish." Nature

22  **513**(7518): 375-381.

23  Chan, Y. F., M. E. Marks, F. C. Jones, G. Villarreal, Jr., M. D. Shapiro, S. D. Brady, A. M.

24  Southwick, D. M. Absher, J. Grimwood, J. Schmutz, R. M. Myers, D. Petrov, B. Jonsson, D.

25  Schluter, M. A. Bell and D. M. Kingsley (2010). "Adaptive evolution of pelvic reduction in

26  sticklebacks by recurrent deletion of a Pitx1 enhancer." Science **327**(5963): 302-305.

27  Colosimo, P. F., C. L. Peichel, K. Nereng, B. K. Blackman, M. D. Shapiro, D. Schluter and D. M.

28  Kingsley (2004). "The genetic architecture of parallel armor plate reduction in threespine

29  sticklebacks." PLoS Biol **2**(5): E109.

30  Conte, M. A. and T. D. Kocher (2015). "An improved genome reference for the African

31  cichlid, Metriaclima zebra." BMC Genomics **16**: 724.

32  Dildrop, R. and U. Ruther (2004). "Organization of Iroquois genes in fish." Dev Genes Evol

33  **214**(6): 267-276.

34  Doudna, J. A. and E. Charpentier (2014). "Genome editing. The new frontier of genome

35  engineering with CRISPR-Cas9." Science **346**(6213): 1258096.

36  Eddy, S. R. (1998). "Profile hidden Markov models." Bioinformatics **14**(9): 755-763.

37  Feijoo, C. G., M. Manzanares, E. de la Calle-Mustienes, J. L. Gomez-Skarmeta and M. L.

38  Allende (2004). "The Irx gene family in zebrafish: genomic structure, evolution and initial

39  characterization of irx5b." Dev Genes Evol **214**(6): 277-284.

40  Feiner, N., G. Begemann, A. J. Renz, A. Meyer and S. Kuraku (2009). "The origin of bmp16, a

41  novel Bmp2/4 relative, retained in teleost fish genomes." BMC Evol Biol **9**: 277.

42  Fitch, W. M. (1970). "Distinguishing homologous from analogous proteins." Syst Zool **19**(2):

43  99-113.

Goldstein, B. and N. King (2016). "The Future of Cell Biology: Emerging Model Organisms." Trends Cell Biol **26**(11): 818-824.

Gray, K. A., B. Yates, R. L. Seal, M. W. Wright and E. A. Bruford (2015). "Genenames.org: the HGNC resources in 2015." Nucleic Acids Res **43**(Database issue): D1079-1085.

Hagen, D. W. (1967). "Isolating Mechanism in Threespine Sticklebacks (Gasterosteus)." Journal of the Fisheries Research Board of Canada **24**(8): 1637-&.

Hammock, E. A., M. M. Lim, H. P. Nair and L. J. Young (2005). "Association of vasopressin 1a receptor levels with a regulatory microsatellite and behavior." Genes Brain Behav **4**(5): 289-301.

Hsu, P. D., E. S. Lander and F. Zhang (2014). "Development and applications of CRISPR-Cas9 for genome engineering." Cell **157**(6): 1262-1278.

Huerta-Cepas, J., F. Serra and P. Bork (2016). "ETE 3: Reconstruction, Analysis, and Visualization of Phylogenomic Data." Mol Biol Evol **33**(6): 1635-1638.

Insel, T. R. (2010). "The challenge of translation in social neuroscience: a review of oxytocin, vasopressin, and affiliative behavior." Neuron **65**(6): 768-779.

Iwasaki, K., M. Taguchi, J. L. Bonkowsky and J. Y. Kuwada (2013). "Expression of arginine vasotocin receptors in the developing zebrafish CNS." Gene Expr Patterns **13**(8): 335-342.

Jones, F. C., M. G. Grabherr, Y. F. Chan, P. Russell, E. Mauceli, J. Johnson, R. Swofford, M. Pirun, M. C. Zody, S. White, E. Birney, S. Searle, J. Schmutz, J. Grimwood, M. C. Dickson, R. M. Myers, C. T. Miller, B. R. Summers, A. K. Knecht, S. D. Brady, H. Zhang, A. A. Pollen, T. Howes, C. Amemiya, P. Broad Institute Genome Sequencing, T. Whole Genome Assembly, J. Baldwin, T. Bloom, D. B. Jaffe, R. Nicol, J. Wilkinson, E. S. Lander, F. Di Palma, K. Lindblad-Toh and D. M. Kingsley (2012). "The genomic basis of adaptive evolution in threespine sticklebacks." Nature **484**(7392): 55-61.

Joung, J. K. and J. D. Sander (2013). "TALENs: a widely applicable technology for targeted genome editing." Nat Rev Mol Cell Biol **14**(1): 49-55.

Juntti, S. A., A. T. Hilliard, K. R. Kent, A. Kumar, A. Nguyen, M. A. Jimenez, J. L. Loveland, P. Mourrain and R. D. Fernald (2016). "A Neural Basis for Control of Cichlid Female Reproductive Behavior by Prostaglandin F2alpha." Curr Biol **26**(7): 943-949.

Katoh, K., K. Misawa, K. Kuma and T. Miyata (2002). "MAFFT: a novel method for rapid multiple sequence alignment based on fast Fourier transform." Nucleic Acids Res **30**(14): 3059-3066.

Katoh, K. and D. M. Standley (2013). "MAFFT multiple sequence alignment software version 7: improvements in performance and usability." Mol Biol Evol **30**(4): 772-780.

Kocher, T. D. (2004). "Adaptive evolution and explosive speciation: the cichlid fish model." Nat Rev Genet **5**(4): 288-298.

Li, H., A. Coghlan, J. Ruan, L. J. Coin, J. K. Heriche, L. Osmotherly, R. Li, T. Liu, Z. Zhang, L. Bolund, G. K. Wong, W. Zheng, P. Dehal, J. Wang and R. Durbin (2006). "TreeFam: a curated database of phylogenetic trees of animal gene families." Nucleic Acids Res **34**(Database issue): D572-580.

Loh, Y. H., E. Bezault, F. M. Muenzel, R. B. Roberts, R. Swofford, M. Barluenga, C. E. Kidd, A. E. Howe, F. Di Palma, K. Lindblad-Toh, J. Hey, O. Seehausen, W. Salzburger, T. D. Kocher and J. T. Streelman (2013). "Origins of shared genetic variation in African cichlids." Mol Biol Evol **30**(4): 906-917.

Lynch, M. and J. S. Conery (2000). "The evolutionary fate and consequences of duplicate genes." Science **290**(5494): 1151-1155.

1 Lynch, M. and A. Force (2000). "The probability of duplicate gene preservation by
2 subfunctionalization." Genetics **154**(1): 459-473.
3 Marques, C. L., I. Fernandez, M. N. Viegas, C. J. Cox, P. Martel, J. Rosa, M. L. Cancela and V.
4 Laize (2016). "Comparative analysis of zebrafish bone morphogenetic proteins 2, 4 and 16:
5 molecular and evolutionary perspectives." Cell Mol Life Sci **73**(4): 841-857.
6 Martin, A. and V. Orgogozo (2013). "The Loci of repeated evolution: a catalog of genetic
7 hotspots of phenotypic variation." Evolution **67**(5): 1235-1250.
8 McKinnon, J. S. and H. D. Rundle (2002). "Speciation in nature: the threespine stickleback
9 model systems." Trends in Ecology & Evolution **17**(10): 480-488.
10 Mcphail, J. D. (1969). "Predation and Evolution of a Stickleback (Gasterosteus)." Journal of
11 the Fisheries Research Board of Canada **26**(12): 3183-&.
12 Nemudryi, A. A., K. R. Valetdinova, S. P. Medvedev and S. M. Zakian (2014). "TALEN and
13 CRISPR/Cas Genome Editing Systems: Tools of Discovery." Acta Naturae **6**(3): 19-40.
14 Ruan, J., H. Li, Z. Chen, A. Coghlan, L. J. Coin, Y. Guo, J. K. Heriche, Y. Hu, K. Kristiansen, R. Li,
15 T. Liu, A. Moses, J. Qin, S. Vang, A. J. Vilella, A. Ureta-Vidal, L. Bolund, J. Wang and R. Durbin
16 (2008). "TreeFam: 2008 Update." Nucleic Acids Res **36**(Database issue): D735-740.
17 Salazar, V. S., L. W. Gamer and V. Rosen (2016). "BMP signalling in skeletal development,
18 disease and repair." Nat Rev Endocrinol **12**(4): 203-221.
19 Schreiber, F., M. Patricio, M. Muffato, M. Pignatelli and A. Bateman (2014). "TreeFam v9: a
20 new website, more species and orthology-on-the-fly." Nucleic Acids Res **42**(Database
21 issue): D922-925.
22 Stamatakis, A. (2006). "RAxML-VI-HPC: maximum likelihood-based phylogenetic analyses
23 with thousands of taxa and mixed models." Bioinformatics **22**(21): 2688-2690.
24 Stamatakis, A., T. Ludwig and H. Meier (2005). "RAxML-III: a fast program for maximum
25 likelihood-based inference of large phylogenetic trees." Bioinformatics **21**(4): 456-463.
26 Taylor, J. S., I. Braasch, T. Frickey, A. Meyer and Y. Van de Peer (2003). "Genome duplication,
27 a trait shared by 22000 species of ray-finned fish." Genome Res **13**(3): 382-390.
28

```
┌─────────────────────────┐      ┌─────────────────────────┐
│  Protein files downloaded│      │   TreeFam files and API  │
│    for 4 cichlid species │      │        downloaded        │
└─────────────────────────┘      └─────────────────────────┘
              ╲                          ╱
               ╲                        ╱
                ◇ treefamscan ◇
                      │
                      ▼
         ┌─────────────────────────┐
         │     TreeFam identified  │
         │      for each protein   │
         └─────────────────────────┘
                      │
                      ▼
                ◇  MAAFT
                   RAxML ◇
                      │
                      ▼
         ┌─────────────────────────┐
         │     Proteins added to   │
         │     TreeFam phylogeny   │
         └─────────────────────────┘
                      │
                      ▼
                ◇  ete3
                   Custom ◇
                      │
                      ▼
         ┌─────────────────────────┐
         │    Human/zebrafish      │
         │  relationships identified│
         └─────────────────────────┘
```

*D. rerio:* 895
*M. zebra:* mz.gene.s93.67
*A. burtoni:* ab.gene.s89.43
*P. nyererei:* pn.gene.s26.157
*N. brichardi:* nb.gene.s3.305
*O. niloticus:* ENSONIP00000011152

AVPR1A HalfOrthlog

*D. rerio:* ZDB-GENE-041210-105
*M. zebra:* mz.gene.s119.8
*P. nyererei:* pn.gene.s290.7
*A. burtoni:* ab.gene.s39.34
*N. brichardi:* nb.gene.s1.292
*O. niloticus:* ENSONIP00000022654

AVPR1A HalfOrthlog

*M. musculus:* ENSMUSP00000020323
*H. sapiens:* AVPR1A

**a**

*H. sapiens:* AVPR1B
*M. musculus:* ENSMUSP00000027690

**b**

-
1.90