

# UCLA

## UCLA Previously Published Works

### Title

Genome-wide ultraconserved elements exhibit higher phylogenetic informativeness than traditional gene markers in percomorph fishes

### Permalink

<https://escholarship.org/uc/item/95m48485>

### Journal

MOLECULAR PHYLOGENETICS AND EVOLUTION, 92(Nov 2015)

### ISSN

1055-7903

### Authors

Gilbert, Princess S  
Chang, Jonathan  
Pan, Calvin  
et al.

### Publication Date

2015

### DOI

10.1016/j.ympcv.2015.05.027

Peer reviewed



# Genome-wide ultraconserved elements exhibit higher phylogenetic informativeness than traditional gene markers in percomorph fishes <sup>☆</sup>



Princess S. Gilbert<sup>a,\*</sup>, Jonathan Chang<sup>a</sup>, Calvin Pan<sup>b</sup>, Eric M. Sobel<sup>d</sup>, Janet S. Sinsheimer<sup>c,d,e</sup>, Brant C. Faircloth<sup>f</sup>, Michael E. Alfaro<sup>a,\*</sup>

<sup>a</sup> Department of Ecology & Evolutionary Biology, University of California, Los Angeles, CA, USA

<sup>b</sup> Department of Medicine, University of California, Los Angeles, CA, USA

<sup>c</sup> Department of Biomathematics, University of California, Los Angeles, CA, USA

<sup>d</sup> Department of Human Genetics, University of California, Los Angeles, CA, USA

<sup>e</sup> Department of Biostatistics, University of California, Los Angeles, CA, USA

<sup>f</sup> Department of Biological Sciences and Museum of Natural Science, Louisiana State University, Baton Rouge, LA, USA

## ARTICLE INFO

### Article history:

Received 14 February 2015

Revised 13 May 2015

Accepted 26 May 2015

Available online 12 June 2015

### Keywords:

Ultraconserved elements  
Next-generation sequencing  
Non-coding DNA  
Phylogenomics  
Molecular evolution  
Phylogenetic informativeness

## ABSTRACT

Ultraconserved elements (UCEs) have become popular markers in phylogenomic studies because of their cost effectiveness and their potential to resolve problematic phylogenetic relationships. Although UCE datasets typically contain a much larger number of loci and sites than more traditional datasets of PCR-amplified, single-copy, protein coding genes, a fraction of UCE sites are expected to be part of a nearly invariant core, and the relative performance of UCE datasets versus protein coding gene datasets is poorly understood. Here we use phylogenetic informativeness (PI) to compare the resolving power of multi-locus and UCE datasets in a sample of percomorph fishes with sequenced genomes (genome-enabled). We compare three data sets: UCE core regions, flanking sequence adjacent to the UCE core and a set of ten protein coding genes commonly used in fish systematics. We found the net informativeness of UCE core and flank regions to be roughly ten-fold and 100-fold more informative than that of the protein coding genes. On a per locus basis UCEs and protein coding genes exhibited similar levels of phylogenetic informativeness. Our results suggest that UCEs offer enormous potential for resolving relationships across the percomorph tree of life.

© 2015 Elsevier Inc. All rights reserved.

## 1. Introduction

Ultraconserved elements (UCEs) have become increasingly popular in recent phylogenomic studies. They have been used to reconstruct phylogenies for clades as divergent as the mammals, fish, birds, turtles, and arthropods (Bejerano et al., 2004; Faircloth et al., 2014, 2013; McCormack et al., 2013; Smith et al., 2014; Sun et al., 2014). The utility of UCEs for sequence-capture approaches has been well justified on practical grounds. They are shared loci found among most, if not all vertebrate genomes (Bejerano et al., 2004; Siepel et al., 2005) and researchers can easily detect and align UCEs from divergent taxonomic groups (Miller et al., 2007). UCEs do not intersect paralogous genes (Derti et al., 2006) or have retroelement insertions (Simons et al., 2006). Stephen et al. (2008) found that most eutherian UCEs were intergenic with only 3% falling

within protein coding exons and suggested splicing regulation as one of their functions. One of the most compelling phylogenetic characteristics of UCEs is that the flanking regions increase in variant sites as the distance from the UCE center increases, allowing for better resolution of nodes across a range of evolutionary timescales in a given phylogeny (Faircloth et al., 2012b). This aspect potentially allows phylogeneticists to tailor their use of UCEs by choosing those with similar evolutionary rates or selecting a subsample of UCE regions whose flanking regions optimize their analyses. However, the relative performance of UCEs compared to traditional molecular markers remains poorly understood.

Traditional markers might be expected to exhibit better phylogenetic performance than UCEs because traditional markers have been highly selected for their potential ability to resolve polytomies and they have been well curated and validated. Sets of traditional markers that yield reasonable phylogenetic results have been identified for many major sections of the tree of life. In fishes for example, Li et al. (2007) identified a cohort of 10 genes from a pool of 154 that have become widely used at various phylogenetic scales (Betancur-R et al., 2013; Li et al., 2009, 2008; Near et al.,

<sup>☆</sup> This paper was edited by the Associate Editor Giacomo Bernardi.

\* Corresponding authors at: Department of Ecology & Evolutionary Biology, 621 Charles E. Young Drive South, University of California, Los Angeles, CA 90095, USA.

E-mail addresses: [ps.gilbert@ucla.edu](mailto:ps.gilbert@ucla.edu) (P.S. Gilbert), [michaelalfaro@ucla.edu](mailto:michaelalfaro@ucla.edu) (M.E. Alfaro).

2012; Wainwright et al., 2012). These protein coding genes were carefully selected and validated for the purpose of reconstructing the ray-finned fish phylogeny (Li et al., 2007). In contrast, UCEs are identified by the presence of nearly invariant core regions. UCE cores are thus expected to have very low to no phylogenetic resolving power. The flanking regions of the UCE are, by definition, not invariant and should thus provide more resolving power than the core. However individual UCE loci have not generally been subjected to the same degree of scrutiny as the phylogenetic workhorse, PCR-amplified, single copy protein coding genes, and thus, on average, might be expected to perform more poorly at resolving phylogenetic problems. One resolution of this paradox would be that the greater degree of resolution obtained in recent UCE studies (Crawford et al., 2012; McCormack et al., 2012) is largely due to the sheer number of sites that are captured through high-throughput sequencing methods, as on a per locus basis the ability of UCEs to resolve polytomies is thought to be relatively poor.

UCE cores are highly conserved throughout the genome, which suggests there may be little phylogenetic informativeness in these regions. More specifically, we ask the question, what is the impact of UCE core conservation on overall phylogenetic informativeness and on the UCEs' ability to resolve hypothetical polytomies?

To better understand the utility of UCEs in a phylogenetic context, we characterize their phylogenetic informativeness (Townsend, 2007) by analyzing a dataset comprised of 1201 UCEs and 10 protein coding genes collected from eight species of percomorphs with fully sequenced genomes (genome-enabled), *Gasterosteus aculeatus*, *Oryzias latipes*, *Takifugu rubripes*, *Tetraodon nigroviridis*, *Oreochromis niloticus*, *Neolamprologus brichardi*, *Pundamilia nyererei* and *Haplochromis burtoni*. We chose to examine the percomorphs because recent studies have demonstrated that this large clade has undergone recent radiations and many relationships remain unresolved, which heavily impact age estimations in the clade (Betancur-R et al., 2013; Broughton et al., 2013; Smith et al., 2007; Wainwright et al., 2012). Li et al. (2007) demonstrated that a carefully chosen set of 10 protein coding genes can successfully resolve many groups within the percomorphs. Faircloth et al. (2012b) demonstrated that UCEs successfully resolve older lineage relationships in the euteleost tree of life but they did not specifically focus on resolving polytomies within sub-clades of the percomorphs, for example the order Perciformes, and it is yet untested whether more recent radiations within the Euteleosts can be resolved using UCEs.

We chose phylogenetic informativeness (PI) to make our comparison. PI estimates the probability that a character resolves a hypothetical polytomy in a four-taxon phylogeny and then remains unchanged along the peripheral branches (Townsend, 2007). PI is a function of the rate of evolutionary change and the time to most recent common ancestor among the taxa under analysis, and it provides one estimate of the amount of phylogenetic signal relative to noise across a specified time period. Marker sets for more than four taxa can be compared using PI if a consistent topology is used across the markers. Calculation of the PI per nucleotide allows estimation of the cost-effectiveness of character sampling. Thus our study seeks to address which dataset, the UCEs or the protein coding genes, has the greatest PI so that researchers interested in clades within the percomorphs can focus on the appropriate data to best resolve the remaining polytomies.

## 2. Materials and methods

### 2.1. UCE core region design pipeline

We identified 1201 UCEs found in the eight percomorphs whose genomes were available at the start of our study, one three-spined

stickleback, *G. aculeatus*, one medaka, *O. latipes*, two puffers, *T. rubripes*, *T. nigroviridis*, and four cichlids, *O. niloticus*, *N. brichardi*, *P. nyererei* and *H. burtoni*. Following Faircloth et al. (2013), we: (1) located nuclear DNA regions of  $180 \pm 10$  base pairs (bp) where there were at least 80 contiguous bp with 100% conservation and the remainder with >80% conservation between *G. aculeatus* and *O. latipes*; (2) aligned these sequences to the genomes of the remaining six fishes (*T. rubripes*, *T. nigroviridis*, *O. niloticus*, *N. brichardi*, *P. nyererei* and *H. burtoni*) using LASTZ (Harris, 2007); and (3) required >80% sequence identity across all eight species. We defined the core as the contiguous region of the aligned sequence, which corresponds to the original 180 bp from *G. aculeatus* and *O. latipes*, and flank as all the remaining sequence 5' or 3' of the core. To ensure that PI is accurately calculated, we limited our analysis to UCEs with at least 50 bp flanking the 5' or 3' end of the core. This reduced the final count used for all further analysis to 988 UCE loci with cores of aligned lengths of 171 bp to 219 bp and flanks of aligned lengths of 144 bp to 1626 bp.

### 2.2. Protein coding genes

We compared the UCEs recovered in this study to ten protein coding genes identified by Li et al. (2007) (see Supplemental Table S1). We downloaded individual gene data for each of these loci across the eight genome-enabled percomorph species from the ENSEMBL Genome Browser (Hubbard et al., 2007), the UCSC genome browser (Kent et al., 2002), and NCBI GenBank (Benson et al., 2005). We translated the nucleotide sequences of the ten loci into amino acid sequences using TranslatorX (Abascal et al., 2010) and aligned amino acids using MUSCLE (Edgar, 2004). We used the DNA version of these alignments when calculating PI.

### 2.3. In silico phylogeny design for the PI guide tree

We constructed a time-calibrated phylogenetic framework needed for calculation of PI using divergence times from recently published phylogenetic studies to date node splits for the eight taxon tree of genome-enabled percomorph fishes (Betancur-R et al., 2013; Broughton et al., 2013; Santini et al., 2009; Wainwright et al., 2012). We provide the time-calibrated phylogeny for the eight genome-enabled species used in this study (Supplemental Fig. 1).

### 2.4. PI Calculations

We used the software package TAPIR (<http://faircloth-lab.github.com/tapir/>) to measure the PI of the UCE core regions, the flanking regions of the UCE cores and the set of ten protein coding genes. TAPIR employs a similar pipeline for estimating PI to that used in PhyDesign (Lopez-Giraldez and Townsend, 2011) although the PI computation is parallelized to work across large genomic datasets (Faircloth et al., 2012a; Pond et al., 2005). TAPIR calculates substitution rates from sequence alignment files and then uses those substitution rates to estimate the PI profile of each locus. We calculated net PI for each dataset, PI per locus per dataset, and PI per nucleotide per locus per dataset. The net PI is the sum of the individual PI's for each nucleotide across all loci in a dataset. Thus, net PI is additive and the length of each dataset contributes to its respective net PI curve. When displaying or analyzing the time of maximum PI, we removed seven UCEs whose cores were invariant across all taxa and thus had PI = 0 across the entire time-calibrated phylogeny.

## 2.5. Statistical analysis

We conducted statistical analyses using the R package (<http://www.r-project.org/>) and TAPIR (<http://faircloth-lab.github.com/tapir/>). We calculated the distribution of the average per nucleotide PI, the maximum nucleotide PI, and the time in millions of years (Ma) of maximum PI using plyr, gtools, and xtable libraries in R and ggplot2 (Harrell and Dupont, 2014; Team, 2014; Warnes et al., 2014; Wickham, 2009, 2011). We performed regression analyses using the lm function of R.

## 2.6. Verification of the percomorph phylogeny

To verify that both the UCE dataset and the protein coding gene dataset produced the expected phylogeny (Dornburg et al., 2014; Faircloth et al., 2013; Near et al., 2013; Wainwright et al., 2012)) we reconstructed the phylogeny for the eight genome-enabled species (Supplemental Table S2). We prepared our data for phylogenetic reconstruction using phyluce (<https://github.com/faircloth-lab/phyluce>). To estimate the best fitting locus-specific site rate substitution models we used Cloudforest (Crawford and Faircloth, 2014) and partitioned the UCEs by their best-fitting substitution models. Bayesian methods were used for phylogenetic inference as implemented in MrBayes 3.1 (Huelsenbeck and Ronquist, 2001; Ronquist and Huelsenbeck, 2003; Ronquist et al., 2012) thus over 5,000,000 iterations we sampled trees every 500 iterations to yield 10,000 trees. Convergence was confirmed by checking Effective Sampling Size values >200 in TRACER (Rambaut et al., 2014).

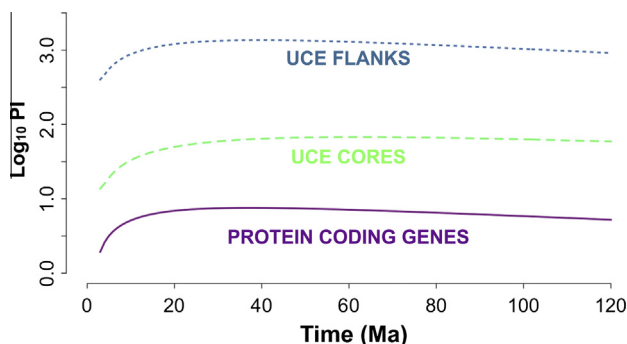
## 3. Results

### 3.1. Net phylogenetic informativeness of each dataset

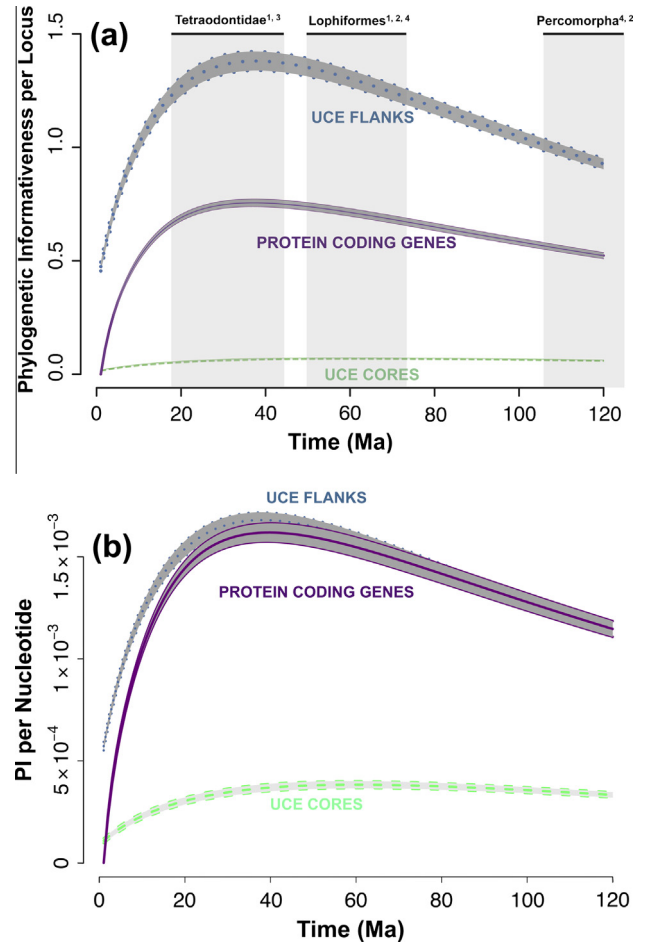
The UCE flanking regions outperformed the UCE core regions, which outperformed the protein coding genes, for estimates of net PI across all times scales (presented as the  $\log_{10}$  of PI versus time in Ma in Fig. 1). PI for the UCE flanks rose rapidly, reached a maximum at 43 Ma and then slowly tapered off. We observed similar behavior for the PI of the UCE cores and the PI of the protein coding genes (Fig. 1).

### 3.2. PI per locus in each data set

The average and 95% confidence interval (CI) for the per locus PI of the UCE flanking regions, the UCE core regions, and the ten protein coding genes are shown versus time in Ma (Fig. 2a). UCE



**Fig. 1.** The  $\log_{10}$  of net phylogenetic informativeness plotted against time for each data type. The blue short dashed line shows UCE flanking regions, the green long dashed line shows the UCE core, and the purple line shows the protein coding genes chosen from Li et al. (2007).



**Fig. 2.** (a and b) The 95% confidence interval for phylogenetic informativeness (PI) per locus (a) and per nucleotide (b) across time. Flanking regions (dotted, blue), UCE core regions (dashed, green) and protein coding genes (solid, purple) overlay a shaded gray region illustrating the average  $\pm 2$  std. errors. The central line is the average PI across all UCEs or loci for each time point. The estimate for the age of the most recent common ancestor (MRCA) of Tetraodontidae, Lophiformes and Percomorpha is plotted on the x-axis of (a) with grey shading. Chen et al., 2014<sup>1</sup>; Near et al., 2013<sup>2</sup>; Santini et al., 2013<sup>3</sup>; Betancur-R et al., 2013<sup>4</sup>. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

flanking regions had the highest PI per locus, surpassing both the UCE core regions and protein coding genes. The UCE core had the lowest per locus PI, reflecting that region's relative invariance.

The ability of UCEs to resolve polytomies depends on the time of divergence from the most recent common ancestor (MRCA) of the polytomy, thus we calculated the time in Ma at which PI is maximized. Based on the average and 95% CI, we observed that the UCE flanking region PI reached its maximum at  $39 \pm 20$  Ma (Fig. 2a), which was similar to that of the protein coding genes, suggesting UCE loci should be suitable for resolving the same polytomies as protein coding genes. Similarly, the maximum PI for UCE cores occurred at  $61 \pm 20$  Ma (Fig. 2a), suggesting these data are suitable for resolving polytomies occurring deeper in time.

To illustrate how these maxima correspond to the age of the MRCA of the percomorphs and two key clades within the percomorphs, we included in Fig. 2a the estimates of the ages of these clades. We use the results of four previously time calibrated phylogenetic reconstructions. The estimates for the MRCA of the Tetraodontidae span from 18 Ma to 44 Ma (Chen et al., 2014; Santini et al., 2013). The maximum PI for the UCE flanking region and for the protein coding genes fall within this range therefore PIs are still driven far more by signal than noise (Townsend,

2007) (Fig. 2a). The estimates for the age of the MRCA of Lophiformes span from 50 Ma to 73 Ma (Betancur-R et al., 2013; Chen et al., 2014). At ~60 Ma, the UCE flanking region PI and the protein coding gene PI have decayed to less than 10% from their maxima indicating again that these loci are still within optimal signal for this clade. The estimates for the age of the MRCA for the percomorphs span from 106 Ma to 133 Ma (Betancur-R et al., 2013; Near et al., 2013; Chen et al., 2014). At ~120 Ma, UCE flanking region PI and the protein coding gene PI have decayed less than 33% from their maxima. These comparisons illustrate that UCE flanking regions are appropriate for resolving polytomies within Tetraodontidae and Lophiformes as well as within Percomorpha.

### 3.3. PI per nucleotide in each data set

The average and 95% CI for the per nucleotide PI of the UCE flanking regions, the UCE core regions and the protein coding genes are shown versus time in Ma in Fig. 2b. The UCE flanking regions had PI values that are slightly higher but similar to the protein coding genes. The UCE core regions had the lowest PI at each time point which is likely a consequence of how UCEs are chosen and the different evolutionary pressures on the UCE cores relative to the UCE flanks or the protein coding genes.

### 3.4. Average PI, Max PI, and time at maximum PI for the UCE core, flank and protein coding datasets

The results shown thus far provide the average UCE behavior for each point in time. When comparing the individual UCEs versus the average behavior across the set, we found that the per nucleotide PI maxima and averages were higher for the flanking regions (mean of max PI =  $1.700 \times 10^{-3}$ , std. dev. of max PI =  $5.955 \times 10^{-4}$  and mean of average PI =  $1.437 \times 10^{-3}$ , std. dev. of average PI =  $4.636 \times 10^{-4}$ ) than for its corresponding core

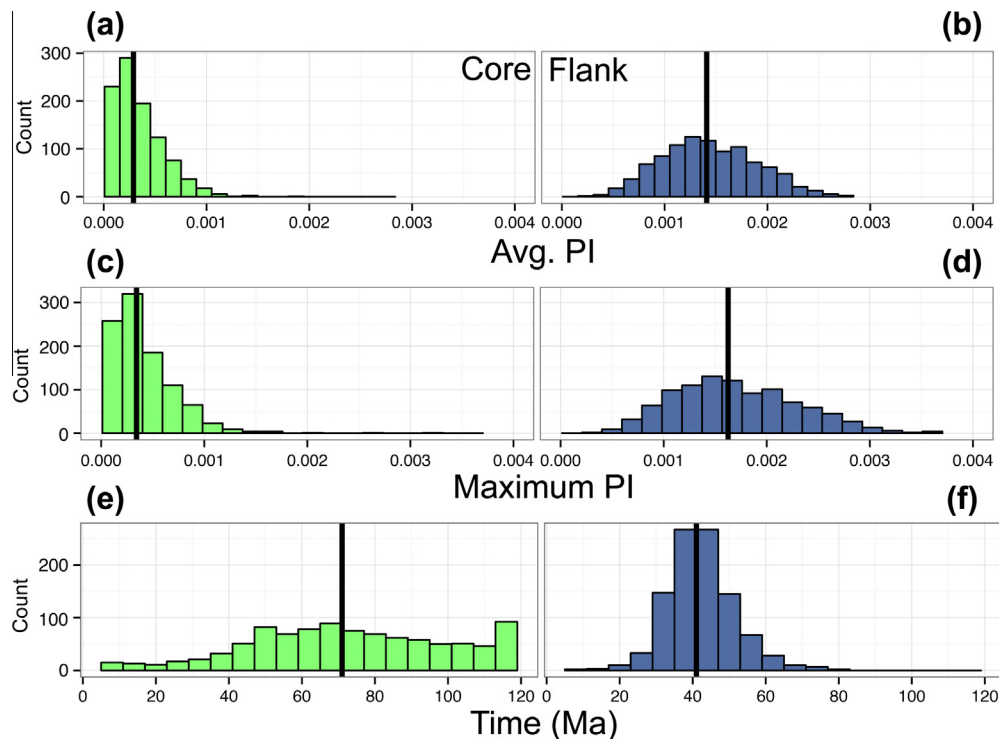
regions (mean of max PI =  $4.097 \times 10^{-4}$ , std. dev. of max PI =  $3.103 \times 10^{-4}$  and mean of average PI =  $2.899 \times 10^{-4}$ , std. dev. of average PI =  $2.443 \times 10^{-4}$ ) and were better approximated by normal distributions (Fig. 3a–d and Supplemental Table S1).

For UCE core regions, the median time of maximum PI was 71 Ma (interquartile range for core = 53 Ma, 94 Ma) but the distribution was quite wide with a number of UCE cores reaching maximum PI at 120 Ma, the oldest time point included in our analysis (Fig. 3e). For the UCE flanking regions, the median time of maximum PI was 41 Ma with an interquartile range for the flank of (36 Ma, 47 Ma, Fig. 3f). For the protein coding genes, the median time of maximum PI was 32 Ma (Table 1) with an interquartile range of (28.75 Ma, 44.25 Ma).

### 3.5. Determinants of PI – linear regression analyses

As expected, there was a strong correlation between average per nucleotide PI and the maximum per nucleotide PI for each locus in the UCE core ( $R^2 = 0.91$ ) and UCE flanking regions ( $R^2 = 0.99$ ) (Supplemental Fig. S3a and S3b). We thus only present results for the average per nucleotide PI. We found a significant but weak correlation between average PI per nucleotide for UCE flanking regions and the average PI per nucleotide for the UCE core regions,  $R^2 = 0.14$  (Fig. 4), indicating that if the UCE had an increased average PI for its core region, they also had an increased PI for its flanking region.

We plotted the average PI per upstream and downstream UCE flanking region against that region's length (Fig. 5). We observed an increasing trend in average PI per region as the flanking region's length increased, as would be expected as variation has been shown to increase with distance from the core (Faircloth et al., 2012b). Further, if we controlled for the average per nucleotide PI of the core, we found that total flank length was a significant predictor of average per nucleotide PI of the flank ( $p < 2.2 \times 10^{-16}$ , Table 2).

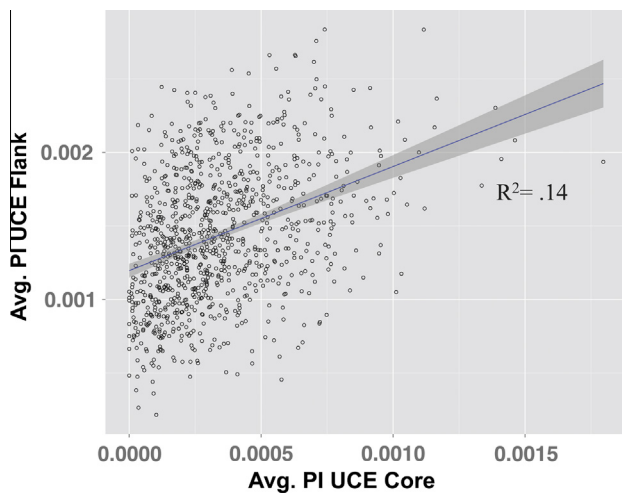


**Fig. 3.** (a–f) UCE core and flank dataset phylogenetic informativeness (PI) distributions. The left column of histograms shows the observed distributions of the core UCE regions. The right column of histograms shows the observed distributions of the flank UCE regions. Average PI per nucleotide for each dataset (a and b); Maximum PI per nucleotide for each dataset (c and d); The time point when PI reaches its maximum for each dataset (e and f). The black line marks the median of each histogram.

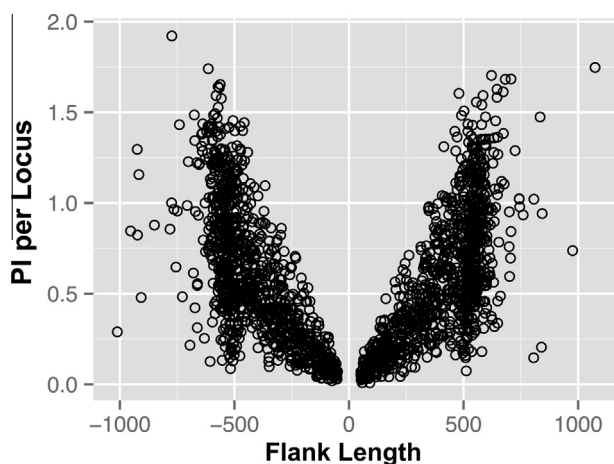
**Table 1**  
Summary statistics for average per nucleotide PI, maximum per nucleotide and time at maximum PI for the core, flank and protein coding genes.

	UCE core <i>avg. per nucleotide PI</i>	UCE flank <i>avg. per nucleotide PI</i>	Protein coding genes <i>avg. per nucleotide PI</i>
Median	$2.889 \times 10^{-4}$	$1.409 \times 10^{-3}$	$1.08 \times 10^{-3}$
Average	$3.406 \times 10^{-4}$	$1.437 \times 10^{-3}$	$1.34 \times 10^{-3}$
Std. deviation	$2.443 \times 10^{-4}$	$4.636 \times 10^{-4}$	$5.8 \times 10^{-4}$
	UCE core <i>max. per nucleotide PI</i>	UCE flank <i>max. per nucleotide PI</i>	Protein coding genes <i>max. per nucleotide PI</i>
Median	$3.430 \times 10^{-4}$	$1.625 \times 10^{-3}$	$1.37 \times 10^{-3}$
Average	$4.097 \times 10^{-4}$	$1.700 \times 10^{-3}$	$1.61 \times 10^{-3}$
Std. deviation	$3.103 \times 10^{-4}$	$5.955 \times 10^{-4}$	$6.8 \times 10^{-4}$
	UCE core <i>time at maximum PI</i>	UCE flank <i>time at maximum PI</i>	Protein coding genes <i>time at maximum PI</i>
Median	71 Ma	41 Ma	32 Ma
Average	72.71 Ma	41.84 Ma	34.9 Ma
Std. deviation	27.39 Ma	9.37 Ma	10.1 Ma

Note: See Section 2.5 and Fig. 3 for details.



**Fig. 4.** Average PI per nucleotide for the UCE flanking regions versus average PI per nucleotide for the UCE core regions. Linear regression results: adjusted  $R^2 = 0.14$ ;  $p$ -value =  $<2.2 \times 10^{-16}$ ; slope = 0.7081; and Y-intercept =  $1.196 \times 10^{-3}$ .



**Fig. 5.** Average PI for each UCE plotted against upstream and downstream flank length.

### 3.6. Verification of the Phylogeny

We recovered the relationships supported in the current literature (Faircloth et al., 2013; Li et al., 2007) with high posterior probabilities using either the protein coding genes or the 988 UCes (Supplemental Fig. S2).

## 4. Discussion

Molecular marker choice is arguably the most important decision made before one embarks on a phylogenetic analysis. Here we explore 3 datasets: UCE core regions, UCE flanking regions and protein coding gene regions, in order to understand PI patterns. UCE flanking and core regions have higher net PI than protein coding genes (Fig. 1). This outcome was expected as there were far more UCes than protein coding genes analyzed. Our analysis corroborates Faircloth et al. (2012b) by finding that the major source of PI for more recent splits is derived from the UCE flanking region and not its core (Faircloth et al., 2012b). Furthermore as the flanking region length increased the average per locus PI for that region increased (Fig. 5). We believe this can be attributed to the fact that longer flanking regions had greater sequence diversity and thus higher PI than shorter regions.

A second important result is that on a per nucleotide scale, the UCE flanking regions have similar PI to protein coding genes (Fig. 2b). *A priori*, we suspected that the protein coding genes would have greater PI than the UCE flanking regions on a per-nucleotide and per-locus level because the protein coding genes we used were carefully selected and validated to be useful in reconstructing the ray-finned fish phylogeny (Li et al., 2007). UCE flanking regions show more variation than the UCE cores and yet are still readily aligned among a set of taxa such as the percomorphs chosen for our analysis. Although we suspect that our results extend beyond these eight taxa, it would be interesting to determine if they hold for a larger set of fishes, birds or mammals.

Despite the low PI of UCE core regions on a per-locus or per nucleotide basis (Fig. 2a and b), the net PI of the UCE cores exceeds that of the protein coding genes (Fig. 1). Although UCes are highly conserved, they still yield varying levels of PI. The explanation for UCE cores exceeding protein coding genes in net PI is sheer loci number. The median time when UCE cores reach its maximum PI is greater than the median time when the UCE flanks reach its maximum PI (Fig. 3 and Table 1), suggesting that UCE cores may be more useful for resolving phylogenetic relationships than previously thought, relationships that are more ancient than the radiation of the percomorphs. Therefore UCE core regions can and should be retained in a phylogenetic reconstruction along with the UCE flanking regions.

Our choice of phylogenetic informativeness as a measure of the suitability of a marker stems from a growing body of publications that demonstrate the comparative quality of PI (Lopez-Giraldez et al., 2013; Schoch et al., 2009; Townsend, 2007; Townsend and Leuenberger, 2011; Townsend et al., 2008). We believe PI holds the key to framing quantitative comparisons of marker types and gives researchers the ability to choose markers based on real data and not just hypothetical assumptions. However PI has garnered

**Table 2**  
Multiple linear regression analysis of PI per nucleotide for the flanking region.

Coefficients	Estimate	Std. Error	t-value	Pr (> t )
Y-intercept	$1.042 \times 10^{-3}$	$5.178 \times 10^{-5}$	20.114	$<2 \times 10^{-16}$
Average PI per nucleotide in the core region	$7.322 \times 10^{-1}$	$5.623 \times 10^{-2}$	13.022	$<2 \times 10^{-16}$
Total flank length	$1.795 \times 10^{-7}$	$5.361 \times 10^{-8}$	3.349	$8.41 \times 10^{-4}$

Note: Residual standard error of  $4.28 \times 10^{-4}$  on 985 degrees of freedom. Adjusted  $R^2$  of 0.149, F-statistic of 86.23 on 2 and 985 degrees of freedom. P-value  $<2.2 \times 10^{-16}$ .

criticism in regards to possible biases placed on fast evolving characters in a given sequence or gene and reduced applicability to real datasets with greater than four taxa (Klopfstein et al., 2010). Per Townsend and Leuenberger (2011), we limited our interpretation of PI profiles to details of the phylogeny on which we based our analyses. Detection of the phylogenetic signal, the subsequent loss of that signal and replacement with non-informative character states all depend upon the specific time epoch one is interested in studying.

In summary, our study provides preliminary evidence that the net phylogenetic informativeness of ultraconserved elements, at both flank and core regions, is superior to the phylogenetic informativeness of the set of protein coding genes recommended for resolving polytomies in the percomorphs. The improvement over the protein coding genes in net phylogenetic informativeness is made possible due to the large number of UCEs that can be detected and aligned among these taxa. It is also a novel finding of this study that UCE flanking regions and protein coding genes have similar levels of per nucleotide phylogenetic informativeness. Although a more comprehensive test with more taxa is required to insure that these results are not limited to the specific clades tested here, our results suggest that UCEs are likely to be an effective means for resolving relationships within percomorphs across a range of time scales.

## Acknowledgments

This work was supported by the National Institute of Health-Genomic Analysis Training Program (T32HG002536 to P.S.G.); the U.S. Department of Education Graduate Assistance in Areas of National Need (to P.S.G.); the Whitcome Research Fellowship (to J.C.); and the National Science Foundation (DMS 1264153 to J.S.S., DEB 6701648 and DEB 6861953 to M.E.A.

## Appendix A. Supplementary material

Supplementary data associated with this article can be found, in the online version, at <http://dx.doi.org/10.1016/j.ympev.2015.05.027>.

## References

- Abascal, F., Zardoya, R., Telford, M.J., 2010. TranslatorX: multiple alignment of nucleotide sequences guided by amino acid translations. *Nucl. Acids Res.* 38, W7–W13.
- Bejerano, G., Pheasant, M., Makunin, I., Stephen, S., Kent, W.J., Mattick, J.S., Haussler, D., 2004. Ultraconserved elements in the human genome. *Science* 304, 1321–1325.
- Benson, D.A., Karsch-Mizrachi, I., Lipman, D.J., Ostell, J., Wheeler, D.L., 2005. GenBank. *Nucl. Acids Res.* 33, D34–D38.
- Betancur-R, R., Broughton, R.E., Wiley, E.O., Carpenter, K., Lopez, J.A., Li, C., Holcroft, N.I., Arcila, D., Sanciangco, M., Cureton II, J.C., Zhang, F., Buser, T., Campbell, M.A., Ballesteros, J.A., Roa-Varon, A., Willis, S., Borden, W.C., Rowley, T., Reneau, P.C., Hough, D.J., Lu, G., Grande, T., Arratia, G., Orti, G., 2013. The tree of life and a new classification of bony fishes. *PLoS Curr.* 5. <http://dx.doi.org/10.1371/currents.tol.53ba26640df0ccaee75bb165c8c26288>.
- Broughton, R.E., Betancur-R, R., Li, C., Arratia, G., Orti, G., 2013. Multi-locus phylogenetic analysis reveals the pattern and tempo of bony fish evolution. *PLoS Curr.* 5. <http://dx.doi.org/10.1371/currents.tol.2ca8041495ffafdc92756e75247483e>.
- Chen, W.-J., Santini, F., Carnevale, G., Chen, J.-N., Liu, S.-H., Lavoué, S., Mayden, R.L., 2014. New insights on early evolution of spiny-rayed fishes (Teleostei: Acanthomorpha). *Front. Mar. Sci.* 1. <http://dx.doi.org/10.3389/fmars.2014.00053>.
- Crawford, N.G., Faircloth, B.C., 2014. Cloudforest: Code to Calculate Species Trees from Large Genomic Datasets. doi: <http://dx.doi.org/10.5281/zenodo.12259>.
- Crawford, N.G., Faircloth, B.C., McCormack, J.E., Brumfield, R.T., Winker, K., Glenn, T.C., 2012. More than 1000 ultraconserved elements provide evidence that turtles are the sister group of Archosaurs. *Biol. Lett.* 8, 783–786.
- Derti, A., Roth, F.P., Church, G.M., Wu, C.T., 2006. Mammalian ultraconserved elements are strongly depleted among segmental duplications and copy number variants. *Nat. Genet.* 38, 1216–1220.
- Dornburg, A., Townsend, J.P., Friedman, M., Near, T.J., 2014. Phylogenetic informativeness reconciles ray-finned fish molecular divergence times. *BMC Evol. Bio.* 14, 169. <http://dx.doi.org/10.1186/s12862-014-0169-0>.
- Edgar, R.C., 2004. MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucl. Acids Res.* 32, 1792–1797.
- Faircloth, B.C., Branstetter, M.G., White, N.D., Brady, S.G., 2014. Target enrichment of ultraconserved elements from arthropods provides a genomic perspective on relationships among Hymenoptera. *Mol. Ecol. Res.* 15, 489–501. <http://dx.doi.org/10.1111/1755-0998.12328>.
- Faircloth, B.C., Chang, J., Alfaro, M.E., 2012a. TAPIR Enables High-throughput Estimation and Comparison of Phylogenetic Informativeness using Locus-specific Substitution Models. arXiv preprint arXiv:12021215 2012, p. 1215.
- Faircloth, B.C., McCormack, J.E., Crawford, N.G., Harvey, M.G., Brumfield, R.T., Glenn, T.C., 2012b. Ultraconserved elements anchor thousands of genetic markers spanning multiple evolutionary timescales. *Syst. Biol.* 61, 717–726.
- Faircloth, B.C., Sorenson, L., Santini, F., Alfaro, M.E., 2013. A phylogenomic perspective on the radiation of ray-finned fishes based upon targeted sequencing of ultraconserved elements (UCEs). *PLoS One* 8, e65923.
- Harrell Jr., F.E., Dupont, M.C., 2014. R Package Hmisc. R Foundation for Statistical Computing, Vienna, Austria.
- Harris, R.S., 2007. Improved Pairwise Alignment of Genomic DNA. Computer Science and Engineering. The Pennsylvania State University, PA, USA.
- Hubbard, T.J., Aken, B.L., Beal, K., Ballester, B., Caccamo, M., Chen, Y., Clarke, L., Coates, G., Cunningham, F., Cutts, T., Down, T., Dyer, S.C., Fitzgerald, S., Fernandez-Banet, J., Graf, S., Haider, S., Hammond, M., Herrero, J., Holland, R., Howe, K., Howe, K., Johnson, N., Kahari, A., Keefe, D., Kokocinski, F., Kulesha, E., Lawson, D., Longden, I., Melsopp, C., Megy, K., Meidl, P., Ouverdin, B., Parker, A., Pricl, A., Rice, S., Rios, D., Schuster, M., Sealy, I., Severin, J., Slater, G., Smedley, D., Spudich, G., Trevanion, S., Vilella, A., Vogel, J., White, S., Wood, M., Cox, T., Curwen, V., Durbin, R., Fernandez-Suarez, X.M., Flicek, P., Kasprzyk, A., Proctor, G., Searle, S., Smith, J., Ureta-Vidal, A., Birney, E., 2007. Ensembl 2007. *Nucl. Acids Res.* 35, D610–617.
- Huelsenbeck, J.P., Ronquist, F., 2001. MRBAYES: Bayesian inference of phylogeny. *Bioinformatics* 17, 754–755.
- Kent, W.J., Sugnet, C.W., Furey, T.S., Roskin, K.M., Pringle, T.H., Zahler, A.M., Haussler, David., 2002. The human genome browser at UCSC. *Genome Res.* 12, 996–1006.
- Klopfstein, S., Kropf, C., Quicke, D.L., 2010. An evaluation of phylogenetic informativeness profiles and the molecular phylogeny of diplazontinae (Hymenoptera, Ichneumonidae). *Syst. Biol.* 59, 226–241.
- Li, B., Dettai, A., Cruaud, C., Couloux, A., Desoutter-Meniger, M., Lecointre, G., 2009. RNF213, a new nuclear marker for acanthomorph phylogeny. *Mol. Phylog. Evol.* 50, 345–363.
- Li, C., Lu, G., Orti, G., 2008. Optimal data partitioning and a test case for ray-finned fishes (Actinopterygii) based on ten nuclear loci. *Syst. Biol.* 57, 519–539.
- Li, C., Orti, G., Zhang, G., Lu, G., 2007. A practical approach to phylogenomics: the phylogeny of ray-finned fish (Actinopterygii) as a case study. *BMC Evol. Biol.* 7, 44.
- Lopez-Giraldez, F., Moeller, A.H., Townsend, J.P., 2013. Evaluating phylogenetic informativeness as a predictor of phylogenetic signal for metazoan, fungal, and mammalian phylogenomic data sets. *Biomed. Res. Int.* 2013, 621604. <http://dx.doi.org/10.1155/2013/621604>.
- Lopez-Giraldez, F., Townsend, J.P., 2011. PhyDesign: an online application for profiling phylogenetic informativeness. *BMC Evol. Biol.* 11, 152.
- McCormack, J.E., Faircloth, B.C., Crawford, N.G., Gowaty, P.A., Brumfield, R.T., Glenn, T.C., 2012. Ultraconserved elements are novel phylogenomic markers that resolve placental mammal phylogeny when combined with species-tree analysis. *Genome Res.* 22, 746–754.
- McCormack, J.E., Harvey, M.G., Faircloth, B.C., Crawford, N.G., Glenn, T.C., Brumfield, R.T., 2013. A phylogeny of birds based on over 1500 loci collected by target enrichment and high-throughput sequencing. *PLoS One* 8, e54848.

- Miller, W., Rosenbloom, K., Hardison, R.C., Hou, M., Taylor, J., Raney, B., Burhans, R., King, D.C., Baertsch, R., Blankenberg, D., Kosakovsky Pond, S.L., Nekrutenko, A., Giardine, B., Harris, R.S., Tyekucheva, S., Diekhans, M., Pringle, T.H., Murphy, W.J., Lesk, A., Weinstock, G.M., Lindblad-Toh, K., Gibbs, R.A., Lander, E.S., Siepel, A., Haussler, D., Kent, W.J., 2007. 28-Way vertebrate alignment and conservation track in the UCSC genome browser. *Genome Res.* 17, 1797–1808.
- Near, T.J., Dornburg, A., Eytan, R.I., Keck, B.P., Smith, W.L., Kuhn, K.L., Moore, J.A., Price, S.A., Burbrink, F.T., Friedman, M., Wainwright, P.C., 2013. Phylogeny and tempo of diversification in the superradiation of spiny-rayed fishes. *Proc. Natl. Acad. Sci. USA* 110, 12738–12743.
- Near, T.J., Dornburg, A., Kuhn, K.L., Eastman, J.T., Pennington, J.N., Patarnello, T., Zane, L., Fernández, D.A., Jones, C.D., 2012. Ancient climate change, antifreeze, and the evolutionary diversification of Antarctic fishes. *Proc. Natl. Acad. Sci.* 109, 3434–3439.
- Pond, S.L., Frost, S.D., Muse, S.V., 2005. HyPhy: hypothesis testing using phylogenies. *Bioinformatics* 21, 676–679.
- Rambaut, A., Suchard, M.A., Xie, D., Drummond, A.J., 2014. Tracer v.1.6- MCMC Trace Analysis Tool. <<http://beast.bio.ed.ac.uk/Tracer>>.
- Ronquist, F., Huelsenbeck, J.P., 2003. MrBayes 3: Bayesian phylogenetic inference under mixed models. *Bioinformatics* 19, 1572–1574.
- Ronquist, F., Teslenko, M., van der Mark, P., Ayres, D.L., Darling, A., Höhna, S., Larget, B., Liu, L., Suchard, M.A., Huelsenbeck, J.P., 2012. MrBayes 3.2: efficient Bayesian phylogenetic inference and model choice across a large model space. *Syst. Biol.* 61, 539–542.
- Santini, F., Harmon, L.J., Carnevale, G., Alfaro, M.E., 2009. Did genome duplication drive the origin of teleosts? A comparative study of diversification in ray-finned fishes. *BMC Evol. Biol.* 9, 194.
- Santini, F., Nguyen, M.T.T., Sorenson, L., Waltzek, T.B., Lynch Alfaro, J.W., Eastman, J.M., Alfaro, M.E., 2013. Do habitat shifts drive diversification in teleost fishes? An example from the pufferfishes (Tetraodontidae). *J. Evol. Biol.* 26, 1003–1018.
- Schoch, C.L., Sung, G.H., Lopez-Giraldez, F., Townsend, J.P., Miadlikowska, J., Hofstetter, V., Robbertse, B., Matheny, P.B., Kauff, F., Wang, Z., Gueidan, C., Andrie, R.M., Trippe, K., Ciuffetti, L.M., Wynns, A., Fraker, E., Hodkinson, B.P., Bonito, G., Groenewald, J.Z., Arzanlou, M., de Hoog, G.S., Crous, P.W., Hewitt, D., Pfister, D.H., Peterson, K., Gryzenhout, M., Wingfield, M.J., Aptroot, A., Suh, S.O., Blackwell, M., Hillis, D.M., Griffith, G.W., Castlebury, L.A., Rossman, A.Y., Lumbsch, H.T., Lücking, R., Budel, B., Rauhut, A., Diederich, P., Ertz, D., Geiser, D.M., Hosaka, K., Inderbitzin, P., Kohlmeyer, J., Volkmann-Kohlmeyer, B., Mostert, L., O'Donnell, K., Sipman, H., Rogers, J.D., Shoemaker, R.A., Sugiyama, J., Summerbell, R.C., Untereiner, W., Johnston, P.R., Stenroos, S., Zuccaro, A., Dyer, P.S., Crittenden, P.D., Cole, M.S., Hansen, K., Trappe, J.M., Yahr, R., Lutzoni, F., Spatafora, J.W., 2009. The Ascomycota tree of life: a phylum-wide phylogeny clarifies the origin and evolution of fundamental reproductive and ecological traits. *Syst. Biol.* 58, 224–239.
- Siepel, A., Bejerano, G., Pedersen, J.S., Hinrichs, A.S., Hou, M., Rosenbloom, K., Clawson, H., Spieth, J., Hillier, L.W., Richards, S., Weinstock, G.M., Wilson, R.K., Gibbs, R.A., Kent, W.J., Miller, W., Haussler, D., 2005. Evolutionarily conserved elements in vertebrate, insect, worm, and yeast genomes. *Genome Res.* 15, 1034–1050.
- Simons, C., Pheasant, M., Makunin, I.V., Mattick, J.S., 2006. Transposon-free regions in mammalian genomes. *Genome Res.* 16, 164–172.
- Smith, B.T., Harvey, M.G., Faircloth, B.C., Glenn, T.C., Brumfield, R.T., 2014. Target capture and massively parallel sequencing of ultraconserved elements for comparative studies at shallow evolutionary time scales. *Syst. Biol.* 63, 83–95.
- Smith, W.L., Craig, M.T., Quattro, J.M., 2007. Casting the Percomorph Net Widely: the importance of broad taxonomic sampling in the search for the placement of Serranid and Percid fishes. *Copeia* 1, 35–55.
- Stephen, S., Pheasant, M., Makunin, I.V., Mattick, J.S., 2008. Large-scale appearance of ultraconserved elements in tetrapod genomes and slowdown of the molecular clock. *Mol. Biol. Evol.* 25, 402–408.
- Sun, K., Meiklejohn, K.A., Faircloth, B.C., Glenn, T.C., Braun, E.L., Kimball, R.T., 2014. The evolution of peafowl and other taxa with ocelli (eyespot): a phylogenomic approach. *Proc. Roy. Soc. B* 281. <http://dx.doi.org/10.1098/rspb.2014.0823>.
- R Core Team, 2014. R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. <<http://www.R-project.org/>>.
- Townsend, J.P., 2007. Profiling phylogenetic informativeness. *Syst. Biol.* 56, 222–231.
- Townsend, J.P., Leuenberger, C., 2011. Taxon sampling and the optimal rates of evolution for phylogenetic inference. *Syst. Biol.* 60, 358–365.
- Townsend, J.P., Lopez-Giraldez, F., Friedman, R., 2008. The phylogenetic informativeness of nucleotide and amino acid sequences for reconstructing the vertebrate tree. *J. Mol. Evol.* 67, 437–447.
- Wainwright, P.C., Smith, W.L., Price, S.A., Tang, K.L., Sparks, J.S., Ferry, L.A., Kuhn, K.L., Eytan, R.I., Near, T.J., 2012. The evolution of pharyngognath: a phylogenetic and functional appraisal of the pharyngeal jaw key innovation in labroid fishes and beyond. *Syst. Biol.* 61, 1001–1027.
- Warnes, G.R., Bolker, B., Lumley T., 2014. Gtools: Various R Programming Tools. R Package Version 3.4.1. <<http://CRAN.R-project.org/package=gtools>>.
- Wickham, H., 2009. ggplot2: Elegant Graphics for Data Analysis. Springer, New York.
- Wickham, H., 2011. The split-apply-combine strategy for data analysis. *J. Stat. Softw.* 40, 1–29.