



RESEARCH ARTICLE

Whole mitogenome analysis and phylogeny of freshwater fish red-finned catopra (*Pristolepis rubripinnis*) endemic to Kerala, India

S. CHANDHINI¹, YUSUKE YAMANOUE², SNEHA VARGHESE¹, P. H. ANVAR ALI³, V. M. ARJUNAN⁴ and V. J. REJISH KUMAR^{1*}

¹Department of Aquaculture, Kerala University of Fisheries and Ocean Studies, Panangad, Kochi 682 506, India

²The University Museum, University of Tokyo, 7-3-1 Hongo Bunkyo-ku, Tokyo 113-0033, Japan

³Department of Fisheries Resource Management, Kerala University of Fisheries and Ocean Studies, Panangad, Kochi 682 506, India

⁴School of Ocean Studies, Kerala University of Fisheries and Ocean Studies, Panangad, Kochi 682 506, India

*For correspondence. E-mail: rejishkumar@gmail.com.

Received 10 September 2020; revised 27 January 2021; accepted 8 February 2021

Abstract. The freshwater leaf fish *Pristolepis rubripinnis* belongs to the family Pristolepididae, restricted to Pamba and Chalakudy rivers of Kerala, India. In the present study, we sequenced the complete mitogenome of *P. rubripinnis* and analysed its phylogeny in the order Anabantiformes. The 16622-bp long genome comprised of 13 protein-coding genes, two rRNA genes, 22 transfer RNAs (tRNAs) genes and had a noncoding control region. All the protein-coding genes, tRNA and rRNA were located on the heavy strand, except nad6 and eight tRNAs (glutamine, alanine, asparagine, cysteine, tyrosine, serine, glutamic acid and proline) transcribed from L strand. The genome exhibited an overlapping between atp8 and atp6 (2 bp), nad4 and nad4l (2 bp), tRNA^{Ile} and tRNA^{Gln} (1 bp), tRNA^{Thr} and tRNA^{Pro} (1 bp). Around 157 bp, an intergenic spacer was identified. The overall GC-skews and AT-skews of the H-strand mitogenome were -0.35 and 0.079, respectively, revealing that the H-strand consisted of equal amounts of A and T and that the overall nucleotide composition was C skewed. All tRNA genes exhibited cloverleaf secondary structures, while the secondary structure of tRNA^{Ser} lacked a discernible dihydrouridine stem. The phylogenetic analysis of available mitogenomes of Anabantiformes revealed a sister group relationship between Pristolepididae and Channidae. The whole mitogenome of *Pristolepis rubripinnis* will form a molecular resource for further taxonomic and conservation studies on this endemic freshwater fish.

Keywords. Anabantiformes; mitogenome; red-finned catopra; *Pristolepis*; *Pristolepis rubripinnis*.

Introduction

Pristolepis, the only genus in the family Pristolepididae are native to freshwater streams of Western Ghats and Southeast Asia. This genus belongs to the clade Percomorphs, but is unique among Nandidae and Polycentridae in position of parasphenoid tooth patch (Nelson *et al.* 2016). The family Nandidae, Polycentridae and Pristolepididae are commonly known as leaf fishes (Collin *et al.* 2015). The *Pristolepis* genus consist of eight species: *P. fasciata*, *P. grootii*, *P. malabarica*, *P. marginata*, *P. pauciradiata*, *P. pentacantha*, *P. procerus* and *P. rubripinnis* (Froese and Pauly 2019).

Pristolepis species are valued as food and ornamental fish and contributes to the riverine fishery of central Kerala, India (Renjithkumar *et al.* 2011; Renjithkumar *et al.* 2016).

Pristolepis rubripinnis, also known as red-finned Catopra is distributed in the Pamba and Chalakudy rivers of Kerala, India. The colour pattern of *P. rubripinnis* is as follows, it has an orange red soft dorsal fin, soft anal and caudal fins, and a yellow to orange pelvic fin which greatly differ from other *Pristolepis* species. The 4–5 scales above and 10 scales below the lateral line distinguish it from *P. grootii* and *P. marginata*, and the absence of prominent bars on the body helps to identify it from *P. fasciata* (Britz *et al.* 2012). *P. rubripinnis* has not

Supplementary Information: The online version contains supplementary material available at <https://doi.org/10.1007/s12041-021-01292-4>.

Published online: 15 May 2021

Table 1. The organization and characterization of the complete mitochondrial genome of *P. rubripinnis*.

Gene	Position		Size		Codon		Intergenic nucleotide	Anti-codon	Strand
	From	To	Nucleotide	Amino-acid	Start	Stop			
tRNA Phe (tRNA F-TTC)	1	68	68				0	GAA	H
12S rRNA	69	1015	947				0		H
tRNA Val (tRNA V-GTA)	1016	1087	72				1	TAC	H
16S rRNA	1089	2780	1692				0		H
tRNA Leu (tRNA L-TTA)	2781	2854	74				15	TAA	H
nad1	2870	3823	954	318	ATC	T-	10		H
tRNA Ile (tRNA I-ATC)	3834	3903	70				-1	GAT	H
tRNA Gln (tRNA Q-CAA)	3903	3973	71				0	TTG	L
tRNA Met (tRNA M-ATG)	3974	4042	69				0	CAT	H
nad2	4043	5080	1038	346	ATG	TG-	8		H
tRNA Trp (tRNA W-TGA)	5089	5160	72				1	TCA	H
tRNA Ala (tRNA A-GCA)	5162	5230	69				1	TGC	L
tRNA Asn (tRNA N-AAC)	5232	5304	73				36	GTT	L
tRNA Cys (tRNA C-TGC)	5341	5406	66				0	GCA	L
tRNA Tyr (tRNA Y-TAC)	5407	5476	70				7	GTA	L
cox1	5484	7019	1536	512	ATC	T-	9		H
tRNA Ser (tRNA S-TCA)	7029	7099	71				2	TGA	L
tRNA Asp (tRNA D-GAC)	7103	7174	72				6	GTC	H
cox2	7181	7864	684	228	ATG	T-	7		H
tRNA Lys (tRNA K-AAA)	7872	7946	75				1	TTT	H
atp8	7948	8109	162	54	ATG	TGA	-2		H
atp6	8106	8786	681	227	ATG	TA-	2		H
cox3	8789	9571	783	261	ATG	T-	2		H
tRNA Gly (tRNA G-GGA)	9574	9644	71				0	TCC	H
nad3	9645	9992	348	116	ATG	TGA	1		H
tRNA Arg (tRNA R-CGA)	9994	10062	69				0	TCG	H
nad4l	10063	10356	294	98	ATG	TG-	-2		H
nad4	10353	11726	1374	458	ATG	TGA	7		H
tRNA His (tRNA H-CAC)	11734	11802	69				0	GTG	H
tRNA Ser (tRNA S-AGC)	11803	11869	67				4	GCT	H
tRNA Leu (tRNA L-CTA)	11874	11946	73				18	TAG	H
nad5	11965	13776	1812	604	ATA	TA-	8		H
nad6	13785	14303	519	173	ATG	T-	0		L
tRNA Glu (tRNA E-GAA)	14304	14372	69				4	TTC	L
cob	14378	15511	1134	378	ATG	TAA	7		H
tRNA Thr (tRNA T-ACA)	15519	15590	72				-1	TGT	H
tRNA Pro (tRNA P-CCA)	15590	15658	69				0	TGG	L
D loop	15659	16622	963				0		

been assessed for its conservation status based on the IUCN Redlist of threatened species (Dahanukar and Raghavan 2013). Currently no complete mitogenome is available for the family Pristolepididae. Hence the present study aims to characterize the whole mitogenome of the *P. rubripinnis*, endemic to Kerala, India and to analyse its phylogenetic position (based on the mitogenome) among the order Anabantoformes.

Materials and methods

Isolation, qualitative and quantitative analysis of mitochondrial DNA

The *P. rubripinnis* collected from Prayikara, Achankovil River (9°19'0"N 76°28'0"E), Kerala, India on May 2019

was used in this experiment. The mitochondria were isolated using organelle isolation protocol described by Frezza *et al.* (2007). Mitochondrial DNA was isolated from fin tissue using DNeasy 96 Blood & Tissue Kit by Qiagen. It was validated using mtDNA specific (cox primers F: 5'-TCAACCAACCACAAA-GACATTGGCAC-3' R: 5'-TAGACTTCTGGGTGGC-CAAAGAACATCA-3' & cobF: 5'-GGCTGATTTCG GAATATGCAYGCNAAYGG-3' R: 5'-GGGAATGGATCGTAGAATTGCRTANGCRAA-3') genes for qualitative analysis. Amplification was checked on 1.5% agarose gel (loaded with 3 µL) for the single intact band of cox and cob genes. One microlitre of the DNA sample was used for determining the concentration using Nanodrop spectrophotometer.

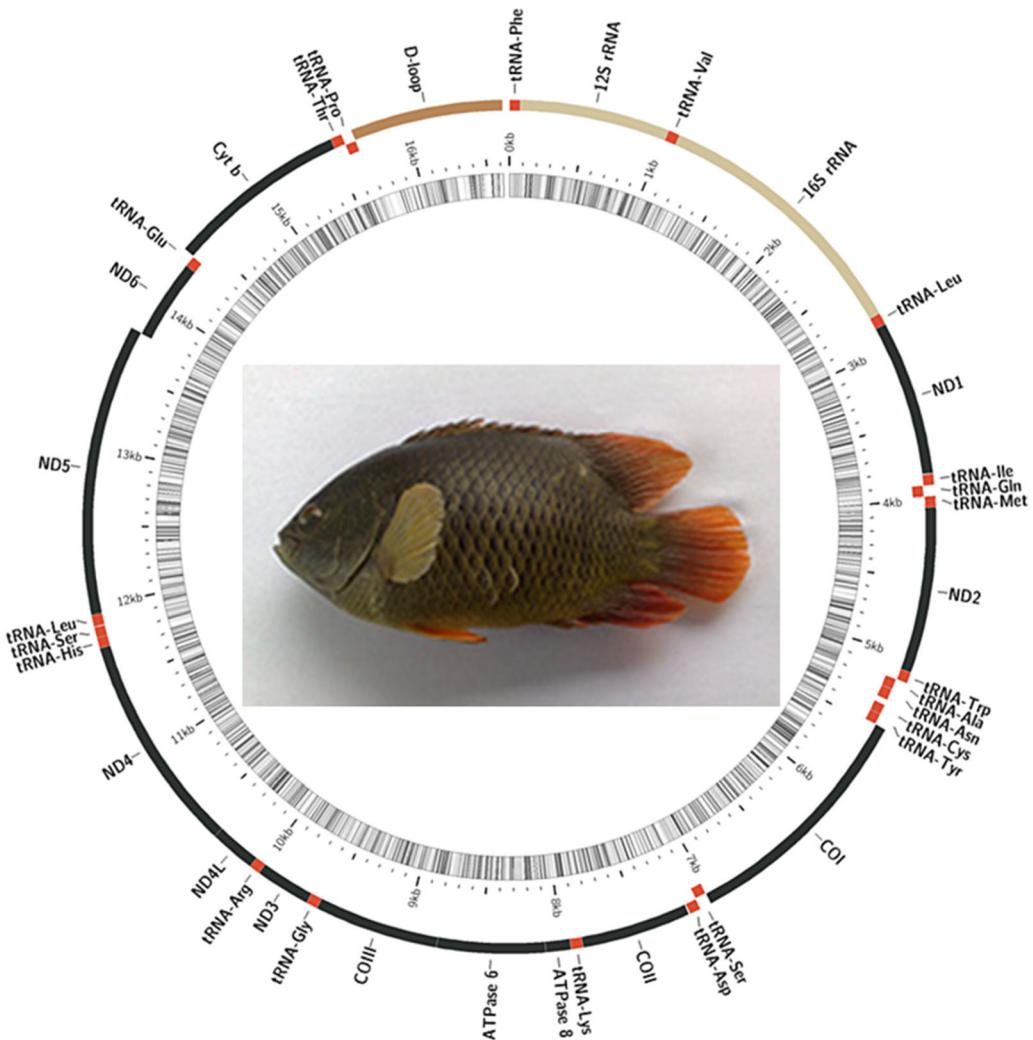


Figure 1. The complete mitochondrial genome organization of *P. rubripinnis*. tRNAs are labelled by their corresponding amino acids and are shown in red colour. Protein-coding genes are marked in black and D-loop in brown. 12S rRNA and 16S rRNA, respectively are shown in beige colour.

Library preparation

The paired-end sequencing library was prepared using QIAseq FX DNA Library kit. The library preparation process was initiated with 100 ng DNA. The DNA was enzymatically sheared into smaller fragments followed by continuous step of end-repair where an 'A' is added to the 3' ends making the DNA fragments ready for adapter ligation. Following this step, platform-specific adapters are ligated to both ends of the DNA fragments. These adapters contain sequences essential for binding dual-barcoded libraries to a flow cell for sequencing, allowing PCR amplification of adapter-ligated fragments, and binding standard Illumina sequencing primers. To ensure maximum yields from limited amounts of starting material, a high-fidelity amplification step was performed using HiFi PCR Master Mix.

Sequencing and analysis

After obtaining the Qubit concentration for the library and the mean peak size from Bioanalyser profile, library was loaded to Illumina HiSeq ((2 × 150 bp chemistry) able to generate ~2.5 gb data) for cluster generation and sequencing. The library molecules were bound to complementary adapter oligos on paired end flow cell. The adapters were designed to allow selective cleavage of the forward strands after resynthesis of the reverse strand during sequencing. The copied reverse strand was then used to sequence from the opposite end of the fragment. The mitogenome was annotated using MITOS (Bernt *et al.* 2013). The transfer RNAs (tRNAs) genes and secondary structures were identified by tRNA Scan SE using default parameters with a cut-off score of 2. Base composition, amino acid composition, skewness, codon usage and other

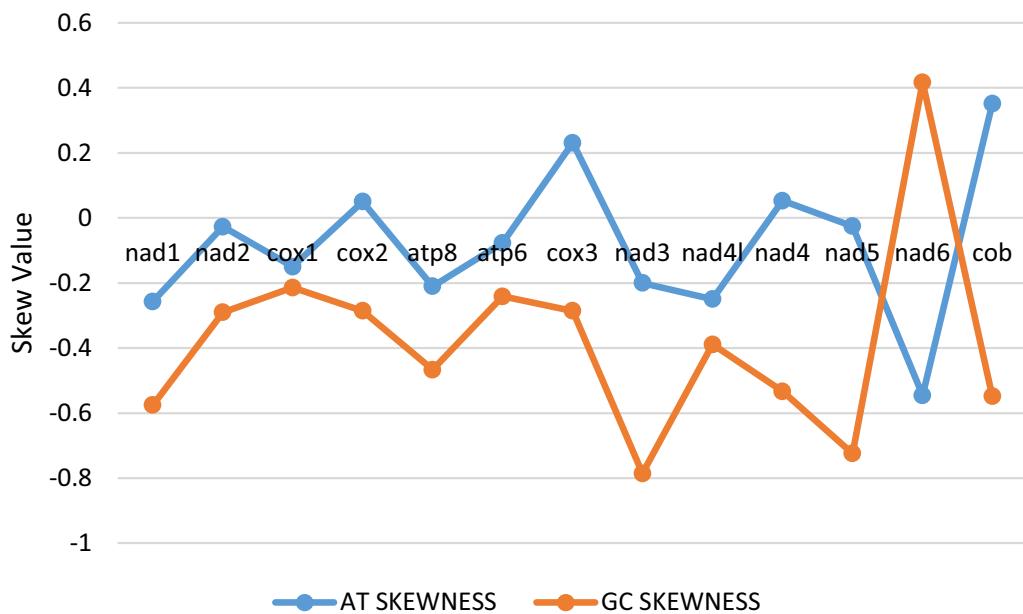


Figure 2. AT-skew and GC-skew in all the 13 protein-coding genes of coding (sense) strands of *P. rubripinnis* mitogenome.

statistics were done by MEGA 7 (Kumar *et al.* 2016). AT skew $[(A-T)/(A+T)]$ and GC skew $[(G-C)/(G+C)]$ were estimated to find the base bias by analysing normalized excess of A over T and C over G, respectively. The A+T content was calculated to estimate strand asymmetry, which is high in many fishes. Total codon distribution (count) and relative synonymous codon usage (RSCU) of mitogenome was analysed with the help of MEGA7. RSCU is the ratio of the number of times a particular codon was observed to the expected frequency of all synonymous codons of same amino acid. All synonymous codons were identified to find RSCU values. A codon that is used less frequently than expected will have an RSCU value of less than 1.00 and vice versa for a codon that is used more frequently than expected. The complete mitochondrial genomic DNA sequence was deposited in the GenBank database (GenBank accession number MK482630).

Phylogenetic analysis

The complete mitochondrial sequence of *P. rubripinnis* along with that of the 27 available species with complete mitogenome of Anabantiforms plus two outgroups, *Mastacembelus favus* and *Indostomus pacificus* were considered for phylogenetic analysis (table 1 in electronic supplementary material at <http://www.ias.ac.in/jgenet/>). The dataset was aligned using MAFFT v.7 (Katoh and Standley 2013) and the aligned sequences were automatically trimmed using trimAL ('-automated1' option; Capella-Gutiérrez *et al.* 2009). Based on the patterns of sequence variations, the dataset with the third codon positions converted by RY-

coding ($12_n3_rT_n$) was expected to effectively remove the likely noise from quickly saturated transitional changes in the third codon positions and avoids a lack of signal by retaining all available positions in the dataset (Phillips and Penny 2003). Accordingly, the three datasets of 123_nRT_n , $12_n3_rT_n$, and 12_nRT_n were prepared and subjected to the phylogenetic analyses. Phylogenetic relationships were inferred with the maximum likelihood method (ML) and Bayesian methods by RAxML v.8 (Stamatakis 2014) and MrBayes v.3.2 (Ronquist *et al.* 2012), respectively. To calculate the robustness of each branch of the resultant tree, a bootstrap analysis with 1000 replications was performed in ML analyses and posterior probabilities were determined in Bayesian analyses based on 10,000 trees, which resulted from every 100 trees sampled from two independent runs of 500,000 replications, after their likelihood scores reached a plateau.

Result and discussion

Genome organization and composition

The *P. rubripinnis* mitochondrial genome size was 16,622 bp with 13 protein-coding genes in order of occurrence: *nad1*, *nad2*, *cox1*, *cox2*, *atp8*, *atp6*, *cox3*, *nad3*, *nad4l*, *nad4*, *nad5*, *nad6*, *cob*, 22 interspersed transfer RNA genes, two ribosomal RNA genes (12S and 16S rRNA) and the noncoding control region (also termed displacement loop region or D-loop) (table 1). Most of the genes were encoded on heavy strand, whereas only *nad6* and eight tRNA (glutamine, alanine, asparagine, cysteine, tyrosine, serine, glutamic acid

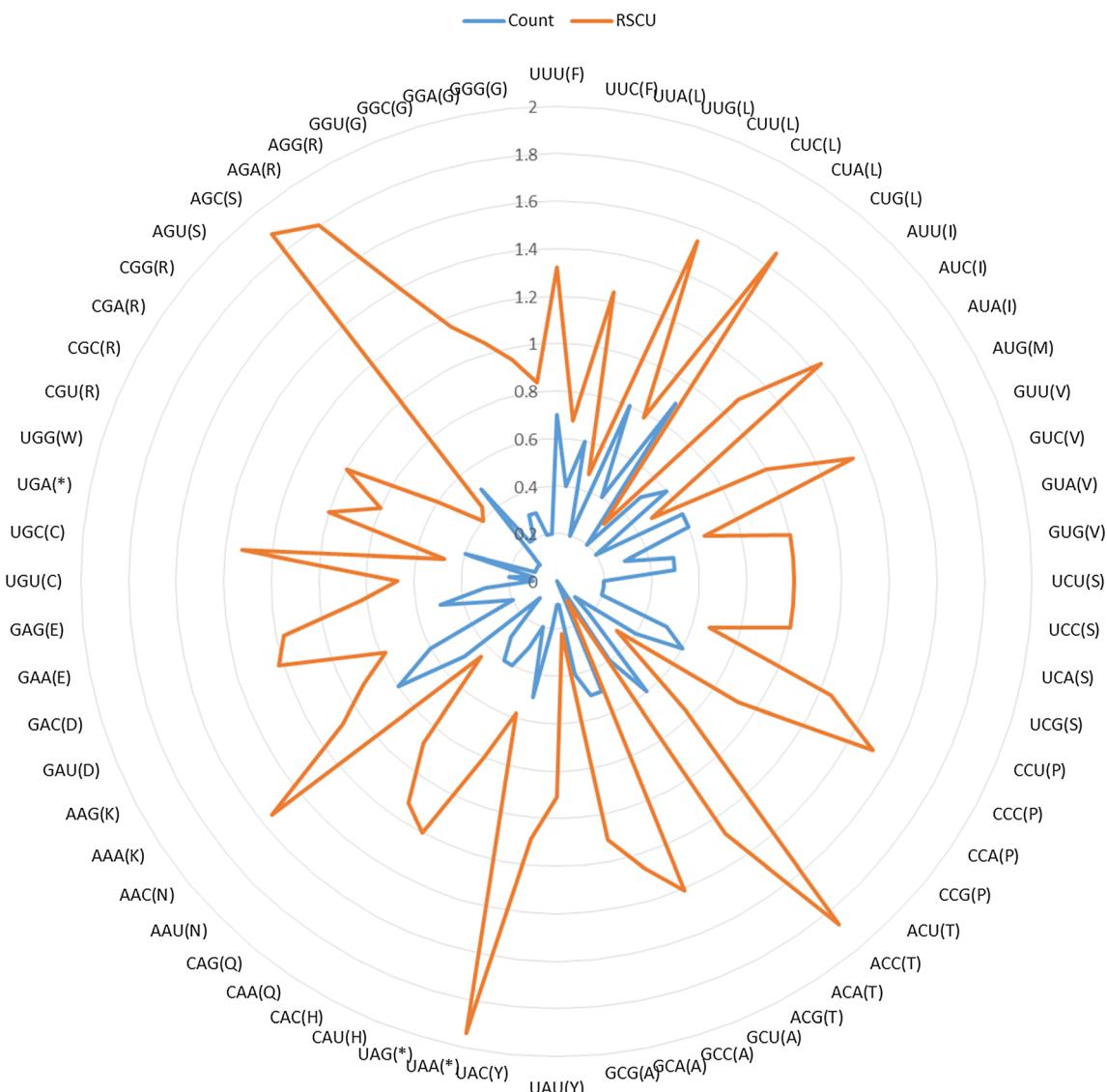


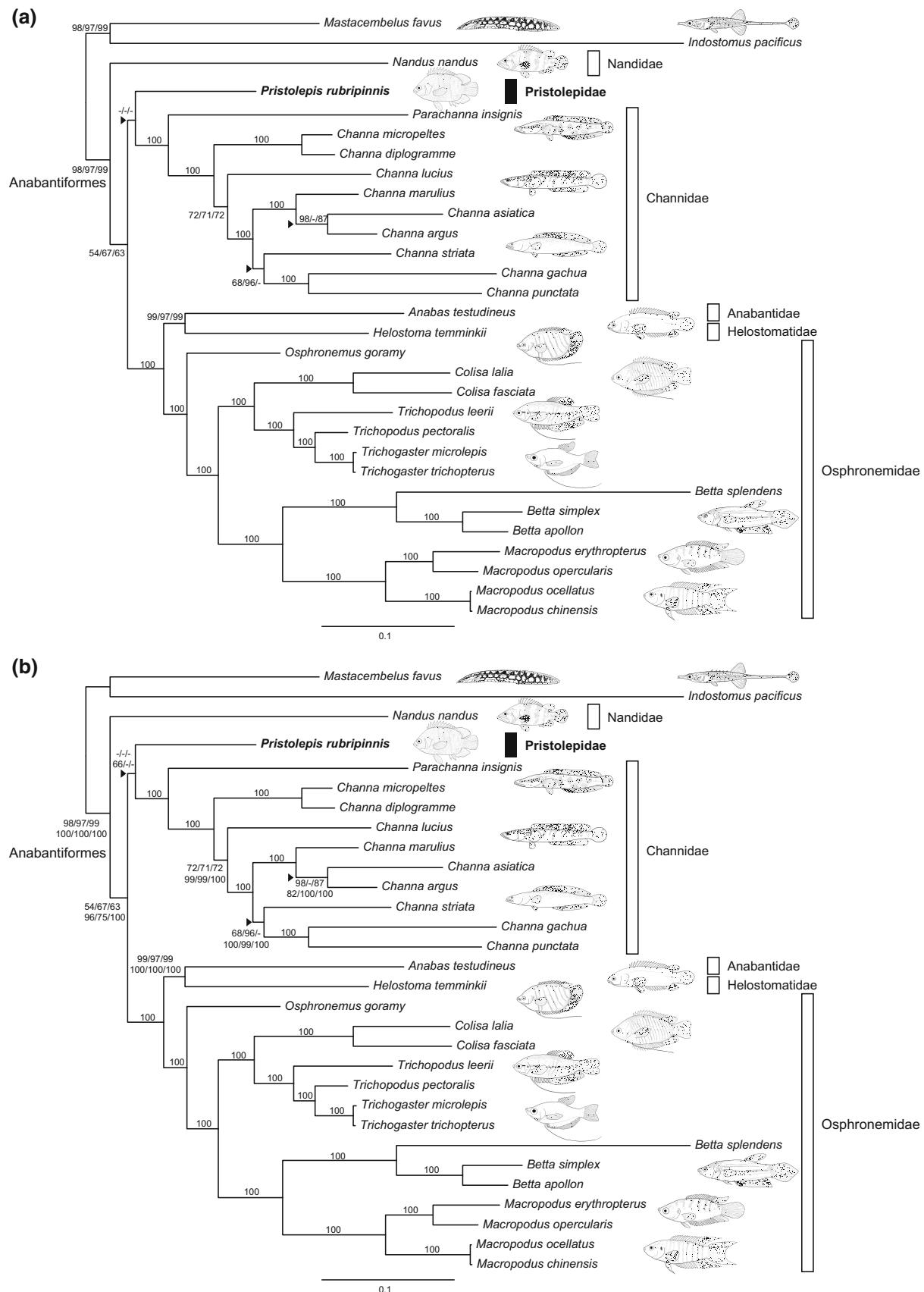
Figure 3. Count and RSCU of *P. rubripinnis*. Blue colour indicates codon count and brown colour indicates RSCU value. For instance, six codons (UUA(L), UUG(L), CUU(L), CUC(L), CUA(L) and CUG(L)) coded for leucine with preference for CUA. RSCU values for these six codons were 1.24, 0.47, 1.55, 0.78, 1.66 and 0.31, respectively in which CUA and UUA were frequently used codons and CUG and UUG least used.

and proline) genes were encoded on light strand (figure 1). Gene overlapping was observed between *atp8* and *atp6* (2 bp), *nad4* and *nad4l* (2 bp), tRNA^{Ile} and tRNA^{Gln} (1 bp), tRNA^{Thr} and tRNA^{Pro} (1 bp). Intergenic spacer was about 157 bp in length. The overall base composition was found to be A: 28.2, T: 24, G: 15.5 and C: 32.4 with increased A+T (52.2) contents as the pattern observed in all other vertebrate. The overall AT skew and GC skew were 0.079 and -0.35, respectively, which confirms base bias towards C over G nucleotide and slight excess of A over T nucleotide. The AT-skew and GC-skew values of the protein-coding genes are shown in figure 2. AT skew and GC skew of protein-coding genes were 0.004 and -0.14, respectively, but only *cox2*, *nad4* and *nad6*, and *cob* showed positive AT skew. All the

protein-coding genes showed negative GC skew, except L strand gene *nad6*.

Protein-coding genes

All the 13 protein-coding genes consisted 11,319 bp and 3773 amino acid, i.e. 68% of the total nucleotides in the mitogenome. The *nad5* gene, transcribed in heavy strand with 1812 bp and 604 amino acids, makes the longest gene in mitogenome for dehydrogenation reaction. *Nad 5* codes for respiratory protein NADH dehydrogenase subunit 5. Other protein-coding genes such as *nad1*, *nad2*, *nad3*, *nad4*, *nad4l*, *nad6*, *cox1*, *cox2*, *cox3*, *atp8*, *atp6* and *cob* encodes six subunit of NADH



► **Figure 4.** Maximum likelihood (ML) (a) and Bayesian (b) tree using the 12n3rRTn data set. Arrowheads indicate the difference from the tree from the 123nRTn and 12nRTn datasets. The numbers near internal branches indicate ML bootstrap probabilities (a) and Bayesian posterior probabilities (b) of 12n3rRTn (left), 123nRTn (middle), and 12nRTn datasets, respectively (values less than 50% are not shown). Single numbers indicate that ML and Bayesian analyses for the three datasets resulted in identical values.

dehydrogenase, three subunits of cytochrome C oxidase, two subunits of ATP synthase and cytochrome b respectively. Shortest gene on whole mitogenome was atp8 (162 bp). Gene rearrangement of vertebrate mitochondrial DNA is relatively conserved and in this mitogenome, gene order was nad1, nad2, cox1, cox2, atp8, atp6, cox3, nad4, nad4l, nad5, nad6 and cox3 as in all other vertebrate, in which nad6 transcribed in light strand to the remaining oriented in heavy strand. The average base composition of protein-coding genes was A: 27.4%, T: 27.1%, G: 19.4% and C: 25.9%. The nad6 contained highest A+T content (64.7%) and cox3 contained high G+C content (61.76%) compared to other protein-coding genes. Only nad6 gene had positive GC skew (0.416), but showed negative AT skew (-0.54). All the 64 codons coded for 20 essential amino acids in which leucine was most frequently used amino acid (14.8) and cysteine (1.14) was the least used.

Codon usage

Among the 64 codons, six codons (UUA(L), UUG(L), CUU(L), CUC(L), CUA(L) and CUG(L)) coded for leucine amino acid with preference for CUA. RSCU values for these six codons were 1.24, 0.47, 1.55, 0.78, 1.66 and 0.31, respectively. CUA and UUA were the frequently used codons and CUG and UUG were least used. Codon usages for all the codons are given in figure 3. A in third position was high (35.6%) to that of second and first positions. First position of all the codons were occupied by A: 28.1, T: 21, G: 24.8 and C: 26%. Second position consisted of low purine content (60.5%) and high pyrimidine content (39.4%). All the protein-coding genes had ATG as start codon except nad1 (ATA), cox1 (ATC), and nad5 (ATA). Only atp8, nad3 and nad4 had TGA as stop codons. Incomplete stop codons such as TA- (atp6, nad5), TG- (nad2, nad4l) and T- (nad1, cox1, cox2, cox3 and nad6) were observed. Cob had TAA as stop codon. Generally, some of the mitochondrial protein-coding genes lack a complete stop codon at their 3' ends, and addition of the polyA tail to the 3' end of processed mRNAs generates a stop codon, such as TAA.

Transfer RNA and ribosomal RNA genes

Twenty-two tRNA genes, which are interspersed between the rRNA genes and protein-coding genes, and two rRNA genes were identified. Base composition of tRNA genes were A: 27.4%, T: 27.2%, G: 19.5%, C: 25.9%. There were

two forms of tRNA^{Ser} (TCA and AGC) among 22 tRNAs, so is tRNA^{Leu} (TTA and CTA). tRNA^{Ser} (AGN) is different from other tRNAs, which lacks the DHC stem. Mismatches were identified in the stem of 22 tRNAs, in the TC stems, in the anticodon stems, in the DHU stems and in the amino acid acceptor stems. Secondary structures of tRNAs are given in figure 1 in electronic supplementary material). The 12S and 16S rRNA genes were 954 and 1675 nucleotide long, respectively. Like all the other mitogenome, 22 tRNA and two rRNA were observed. As in other vertebrates, they are located between tRNA^{Phe} (TTC) and tRNA^{Leu} (TTA) genes, and are separated by the tRNA^{Val} (GTA). They consisted A: 35%, T: 18%, G: 10% and C: 37%.

Phylogenetic analysis

Molecular phylogeny of Anabantiformes based on the available mitogenomes placed *P. rubripinnis* as a sister group of Channidae (figure 4, a&b). Both ML and Baysesian analysis exhibited similar topology. This is the first attempt on phylogenetic analysis based on complete mitochondrial genome to demonstrate the phylogeny of the family Pristolepididae among the order Anabantiformes. Order Anabantiformes consists of two suborders (Anabantoidei and Channoidei) and four families of freshwater fishes (Channidae, Anabantidae, Helostomatidae and Osphronemidae) (Nelson *et al.* 2016). An alternative classification based on osteological and molecular phylogenetic analysis of five nuclear genes expanded the order incorporating the suborder Nandoidei, which includes the families Nandidae, Badidae and Pristolepididae (Collins *et al.* 2015). According to the recent phylogenetic analysis of the bony fishes, order Anabantiformes consists of seven families (Nandidae, Badidae, Pristolepididae, Channidae, Anabantidae, Helostomatidae and Osphronemidae) under three suborders (Nandoidei, Channoidei and Anabantoidei) (Bentancur *et al.* 2017). Previous phylogenetic analysis of leaf fishes included *Pristolepis* in Nandoidei, however the authors were very much ambiguous about this inclusion and pointed out further requirement of molecular analysis to corroborate this phylogeny. In contradiction, our study on mitogenome-based phylogeny has placed *Pristolepis* as a sister group of Channidae. However, the limited availability of the whole mitogenome sequences of the species of Anabantiformes makes it premature to state the exact phylogenetic position of *Pristolepis*. With the advent of next-generation sequencing technologies, mitogenome database of fishes are expanding at high pace and in future, more mitogenomes will make a robust analysis on the phylogeny of Anabantiformes.

Acknowledgement

This work was supported by Centre of Excellence in Sustainable Aquaculture and Aquatic Health Management (CAAHM), Plan fund 2019-20, Kerala University of Fisheries and Ocean Studies.

References

- Bernt M., Donath A., Juhling F., Externbrink F., Florentz C., Fritzsch G. et al. 2013 MITOS: improved de novo metazoan mitochondrial genome annotation. *Mol. Phylogenet. Evol.* **69**, 313–319.
- Betancur R. R., Wiley E. O., Arratia G., Acero A., Bailly N., Miya M., Lecointre G. and Ortí G. 2017 Phylogenetic classification of bony fishes. *BMC Evol. Biol.* **17**, 162.
- Britz R., Kumar K. and Baby F. 2012 *Pristolepis rubripinnis*, a new species of fish from southern India (Teleostei: Percomorpha: Pristolepididae). *Zootaxa* **3345**, 59–68.
- Capella-Gutiérrez S., Silla-Martínez J. M. and Gabaldón T. 2009 trimAl: a tool for automated alignment trimming in large-scale phylogenetic analyses. *Bioinformatics* **25**, 1972–1973.
- Collins R. A., Britz R. and Rüber L. 2015 Phylogenetic systematics of leaf fishes (Teleostei: Polycentridae, Nandidae). *J. Zool. Syst. Evol. Res.* **53**, 259–272.
- Dahanukar N. and Raghavan R. 2013 Freshwater Fishes of Western Ghats - Checklist V.01. August 2013. MIN-Newsletter of IUCN SSC/WI Freshwater Fish Specialist Group-South Asia and the Freshwater Fish Conservation Network of South Asia (FFCNSA). Vol. 1. 6–16.
- Frezza C., Cipolat S. and Scorrano L. 2007 Organelle isolation: functional mitochondria from mouse liver, muscle and cultured fibroblasts. *Nat. Protoc.* **2**, 287–295.
- Froese R. and Pauly D. 2019 *FishBase* (www.fishbase.org).
- Katoh K. and Standley D. M. 2013 MAFFT multiple sequence alignment software version 7: improvements in performance and usability. *Mol. Biol. Evol.* **30**, 772–780.
- Kumar S., Stecher G. and Tamura K. 2016 MEGA7: molecular evolutionary genetics analysis version 7.0 for bigger datasets. *Mol. Biol. Evol.* **33**, 1870–1874.
- Nelson J. S., Grande T. C., Wilson M. V. 2016 *Fishes of the World*. New Jersey: Wiley.
- Phillips M. J. and Penny D. 2003 The root of the mammalian tree inferred from whole mitochondrial genomes. *Mol. Phylogenet. Evol.* **28**, 171–185.
- Renjithkumar C. R., Harikrishnan M. and Kurup B. M. 2011 Exploited fishery resources of Pampa River, Kerala, India. *Indian J. Fish.* **58**, 13–22.
- Renjithkumar C. R., Roshni K. and Kurup B. M. 2016 Exploited fishery resources of Muvattupuzha River, Kerala, India. *Fish. Technol.* **53**, 177–182.
- Ronquist F., Teslenko M., van der Mark P., Ayres D. L., Darling A., Höhna S. et al. 2012 MrBayes 3.2: efficient Bayesian phylogenetic inference and model choice across a large model space. *Syst. Biol.* **61**, 539–542.
- Stamatakis A. 2014 RAxML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies. *Bioinformatics* **30**, 1312–1313.

Corresponding editor: PUNYASLOKE BHADURY