

# Identification of microsatellite loci, gene ontology and functional gene annotations in Indian salmon (*Eleutheronema tetradactylum*) through next-generation sequencing technology using illumina platform

Shihab Ismail\*, N. Vineesh, Reynold Peter, P. Vijayagopal, A. Gopalakrishnan

ICAR Central Marine Fisheries Research Institute, Post Box No: 1603, Ernakulam North P.O., Kochi, Kerala, 682018, India

## ARTICLE INFO

### Keywords:

Microsatellites  
Indian salmon  
Population genetics  
Genetic stocks  
Gene ontology

## ABSTRACT

Whole genome sequencing was performed on three samples of four finger threadfin *Eleutheronema tetradactylum* (KET<sub>25</sub>, KET<sub>29</sub> and KET<sub>30</sub>) using illumina NextSeq500 platform using 2 × 150 bp chemistry. 8,390,317, 7,085,775 and 8,461,589 high quality reads were obtained after trimming low quality reads and adapter sequence. These high quality reads obtained were used for de novo assembly and obtained a number of scaffolds. From these scaffolds of vast sequenced data, we were able to identify 60246, 46107 and 60907 Simple Sequence Repeats (SSR) markers in KET<sub>25</sub>, KET<sub>29</sub> and KET<sub>30</sub> respectively, which will be useful in population genetic analysis and other diversity studies in Indian salmon. The gene prediction on assembled scaffolds predicted 31,943 genes for KET<sub>25</sub>; 26,487 genes for KET<sub>29</sub> and 31,654 genes for KET<sub>30</sub> with average gene size of 458bp, 424bp and 459bp respectively. A total of 30,209, 25,107 and 29,943 genes were annotated against the NCBI Nr database for the samples respectively. *E. tetradactylum* is a commercially important fish species for many countries. This is the first report on the identification of genomic SSR markers in *E. tetradactylum* using NGS technology. This study provides an insight of baseline knowledge of the genome sequence of Indian salmon for future studies.

## 1. Introduction

NGS technologies are making a notable footprint on many areas of biology, including the genetic diversity in populations [1]. NGS is not only a simple genome sequencing method, but have greatly beneficial to the fields of biology, epidemiology, evolutionary biology, phylogenetics, comparative genomics, microbial diversity, DNA marker discovery and studies of gene function and expression [2]. NGS technologies have radically changed the way genetic sequence data are generated and have accelerated a revolution in biological research [3]. Recent developments in NGS technologies have sophisticated the rapid and economical discovery of molecular markers from non-model organisms [4]. Recent developments in sequencing technology, short read sequencers (90–400 bp), such as Illumina and Ion Torrent, are starting to be more frequently used for the generation of large NGS data sets, with affordable price range [5]. NGS has enabled the rapid and cost effective genetic marker discovery, including microsatellites and SNPs and has empowered the massive increase in the number of sequence attained per sequencing effort and also led to the development of high-output genotyping by sequencing [6]; Davey et al., 2011 and [7].

The vast and huge sequencing data generated through NGS technology and whole genome sequencing enables the development of molecular markers for population genetics, molecular systematics, evolutionary developmental biology and gene mapping studies. Population genetics rely essentially on two types of genetic markers; microsatellite markers or short tandem repeats (STRs) and single-nucleotide polymorphisms (SNPs). Microsatellites have been used widely since the late eighties for applications such as parentage analysis, population genetic structure and conservation genetics because of their high level of polymorphism (allelic richness), higher mutation rate, relatively small size and higher statistical power per locus (rapid analysis protocol) [8,9]. Microsatellites, also known as SSRs, are tandem repeated motifs of 1–6 bases and serve as the most important molecular markers in population and conservation genetics, molecular epidemiology and pathology and gene mapping. Majority of population genetic studies in marine fisheries have employed microsatellites, due to its high mutation rate which may results in extremely high level of variation in marine fish [10]. The development of suitable molecular markers would favour studies of wild population structure that will finally result in improved broodstock selection [11].

\* Corresponding author.

E-mail address: [shihabismail51@gmail.com](mailto:shihabismail51@gmail.com) (S. Ismail).

<https://doi.org/10.1016/j.egg.2019.100038>

Received 21 May 2018; Received in revised form 17 December 2018; Accepted 1 April 2019

Available online 02 April 2019

2405-9854/ © 2019 Elsevier Inc. All rights reserved.

Microsatellites generally require species-specific marker development that can be expensive and laborious and limited by the difficulties of de novo development in species without any prior genomic information [12,13]. Other problems analogous with microsatellites include poor level of inter-laboratory calibration with genotype based on fragment size, fragment size-homoplasmy and laborious genotyping [14]. Many of these issues associated with microsatellite based population studies could be eliminated using NGS based microsatellite approach leading to faster and cheaper genotyping in large-scale population genetic studies [15]. The use of microsatellites will only outweigh the use of SNPs if microsatellites can be generated from NGS platforms and used with programs such as MEGASAT [16] for genotype calling.

*Eleutheronema tetradactylum* or four finger threadfin commonly known as blue threadfin also called Indian salmon, is a marine protandrous hermaphrodite species [17] belong to the family polynemidae, and they are distributed in tropical and subtropical waters throughout the world. They are generally found in coastal marine waters, estuaries or rivers in the tropics [18]. They mainly feed on small fishes, prawns, shrimps and mysids and adult fish prey on other fishes [18]. In Western Australia they are considered fully or over-exploited [19,20]. Stock structure analysis of blue threadfin in Australia using parasite and tag-recapture data [21] and by comparing the life history parameters [22] suggests the possibility of separate sub-stocks. There is limited information about the genetic stock structure of this commercially important fish species from Indian waters. As *E. tetradactylum* is commercially important species of the region, it is very essential to utilize the fishery in a sustainable manner.

Here we present the microsatellite identification, gene ontology and functional gene annotations for a hermaphroditic species Indian salmon through NGS technology using Illumina platform. The microsatellite markers identified through this NGS technology herein offer important genetic resources for the assessment, understanding and conservation of this hermaphroditic species and facilitate future work on the other related species.

## 2. Materials and methods

### 2.1. Sampling and DNA isolation

Fresh samples of 6 Indian salmon were collected from commercial fishing operations from Indian waters. Each fish was measured for Total Length (TL) to the 1.0 mm below and total body weight to the nearest 0.01 g. The body cavity was dissected out to assess sex and maturity. A piece of tissue was excised from caudal peduncle of each specimen and stored at 4 °C in absolute ethanol. Total genomic DNA was extracted from 10 mg of tissue from each fish using a chloroform/isoamyl alcohol protocol. Extracted DNA was quantified on a NanoDrop ND-1000 spectrophotometer and was stored at -20 °C until further use. Isolated genomic DNA was sent to the Eurofins Genomics India Private limited (Bangalore, India) for library preparation and sequencing.

### 2.2. Qualitative and quantitative analysis of isolated gDNA

Quality of the fish muscle gDNA was checked on NanoDrop and 2 µl of DNA was resolved on 0.8% Agarose gel at 120 V for approximately 60 min or until the samples reached 3/4th of the gel. 1 µl of sample was loaded in NanoDrop for determining A260/280 ratio and also quantified by Qubit 3.0. All the DNA samples passed the QC and of these three samples (KET<sub>25</sub>; KET<sub>29</sub> and KET<sub>30</sub>) were randomly selected for paired-end sequencing library preparation.

### 2.3. Preparation of 2 × 150 NextSeq500 libraries

The paired-end sequencing libraries were prepared from the QC passed gDNA samples (KET<sub>25</sub>; KET<sub>29</sub> and KET<sub>30</sub>) using TruSeq Nano DNA library Prep Kit (Illumina, San Diego, CA). 200 ng of Genomic

**Table 1**  
Summary of gene Assembly in *E. tetradactylum*.

Description	KET25	KET29	KET30
Number of scaffolds	289,461	280,260	289,650
Total size of assembly	301,704,044	263,802,638	306,253,351
Average size of scaffolds	1042	941	1057
Scaffold N50	1157	1001	1180
Maximum size of scaffold	16,681	17,051	16,622
Minimum size of scaffold	500	500	500

**Table 2**  
Gene statistics in three samples of *E. tetradactylum*.

Description	KET25	KET29	KET30
Number of genes	31,943	26,487	31,654
Average gene length	458	424	459
Maximum gene length	5070	5493	4848
Minimum gene length	201	201	201

DNA from the three samples was fragmented by Covaris M220 (Covaris, Woburn, MA) to generate a mean fragment distribution of 400bp. The focussed ultrasonic shearing using Covaris generates dsDNA fragments with 3' to 5' overhangs. The fragments were then subjected to end-repair. This process converts the overhangs resulting from fragmentation in to blunt ends using End Repair Mix. The 3' to 5' exonuclease activity of this mix removes the 3' overhangs and the 5' to 3' polymerase activity fills in the 5' overhangs followed by adapter ligation to the fragments. This strategy ensures a low rate of chimera formation. The ligated products were size selected using AMPure XP beads (NEB, Ipswich, MA). The size selected product range between 478bp to 492bp was PCR amplified with the index primer. Indexing adapters were ligated to the ends of the DNA fragments, preparing them for hybridization on to a flow cell.

### 2.4. Quality check (QC) of library, cluster generation and sequencing

The PCR amplified libraries were analysed on Tape Station 4200 (Agilent Technologies) using High sensitivity D1000 Screen Tape assay kit as per manufacturer instructions. After obtaining the Qubit concentration for the libraries and the mean peak size from Agilent Tape Station profile, the PE illumina libraries were loaded onto NextSeq 500 for cluster generation and sequencing. Paired-End sequencing allows the template fragments to be sequenced in both the forward and reverse directions on NextSeq 500. The kit reagents were used for binding of samples to complementary adapter oligos on paired-end flow cell. The adapters were designed to allow selective cleavage of the forward strands after re-synthesis of the reverse strand during sequencing. The copied reverse strand will then use to sequence from the opposite end of the fragments.

### 2.5. Statistical analysis

The sequenced raw data was processed to obtain high quality clean reads using Trimmomatic v0.35 to remove adapter sequences, ambiguous reads (reads with unknown nucleotides "N" larger than 50%), and low-quality sequences (reads with more than 10% quality threshold (QV) < 20 phred score). A minimum length of 75 nucleotide after trimming was applied. After removing the adapter and low quality sequences from the raw data, 8,390,317 (2 × 150bp), 7,085,775 (2 × 150bp) and 8,461,589 (2 × 150bp) high quality reads were retained for KET<sub>25</sub>, KET<sub>29</sub> and KET<sub>30</sub> samples respectively. This high quality (QV > 20), paired-end reads were used for de novo assembly of all the samples.



**Table 3**  
Gene Annotation Distributions in three samples of *E. tetradactylum*.

Sr. No	Sample Name	No. of genes	No. of genes with Blast Hit	No. of genes without Blast Hit
1	KET25	31,943	30,209	1734
2	KET29	26,487	25,107	1380
3	KET30	31,654	29,943	1711

**Table 4**  
Summary of Gene Ontology annotations in three samples of *E. tetradactylum*.

Sr. No	Sample Name	Biological Processes	Molecular Functions	Cellular Component
1	KET25	2980	3223	2459
2	KET29	2633	2866	2264
3	KET30	3054	3355	2647

## 2.7. Gene prediction

AUGUSTUS (V.3.2.2) was used to predict genes from the assembled scaffolds with default parameters and *Danio rerio* was used as the model species.

## 2.8. Functional annotation of genes and gene ontology

The predicted genes of samples KET<sub>25</sub>, KET<sub>29</sub> and KET<sub>30</sub> were searched against NCBI non redundant protein (Nr) database using basic local alignment search tool (BlastX). Gene ontology (GO) annotations of the genes were determined by the Blast2GO programs. GO mapping was carried out in order to retrieve GO terms for all the BlastX functionally annotated genes. BlastX result accession IDs are used to retrieve gene names or symbols, identified gene name or symbols are then searched in the species specific entries of the gene-product tables of GO database.

## 2.9. Simple sequence repeat (SSR) and polymorphic SSR identification

The potential SSRs from three individual assemblies (KET<sub>25</sub>; KET<sub>29</sub> and KET<sub>30</sub>) were identified as ranging from dinucleotide motifs with a minimum of ten repeats, tri nucleotide motif with minimum of three repeats, tetra, penta and hexa nucleotide motifs with a minimum of five

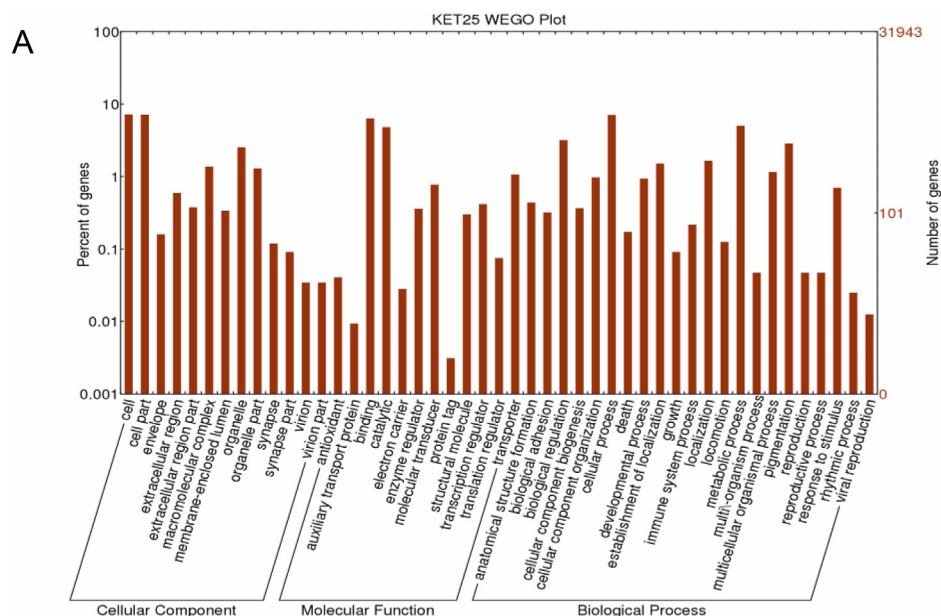
repeats. A maximum distance of 200 nucleotides was allowed between two SSRs.

The scaffolds from three individual assemblies (KET25\_Scaffold.fa, KET29\_Scaffold.fa and KET30\_Scaffold.fa) were clustered using CD-HIT v4 to generate a comprehensive reference. SSR prediction was done from the clustered reference assembly using MISA v1.0. Then we used PSR\_read\_retrieval script for identification of Polymorphic SSR Count based on reads mapping to the SSR region. The consensus sequence was called for each individual samples using high quality reads of KET<sub>25</sub>; KET<sub>29</sub> and KET<sub>30</sub> mapped on reference assembly.

## 3. Results and discussion

Whole genome sequencing was performed on three samples of Indian salmon (*E. tetradactylum*) (KET<sub>25</sub>, KET<sub>29</sub> and KET<sub>30</sub>) using illumina NextSeq500 platform using 2 × 150 bp chemistry. 8,390,317, 7,085,775 and 8,461,589 high quality reads were obtained after trimming low quality reads and adapter sequence. These high quality reads obtained after the trimming were used for de novo assembly and obtained 289,461; 280,260 and 289,650 scaffolds for KET<sub>25</sub>, KET<sub>29</sub> and KET<sub>30</sub> samples respectively. The detailed assembly statistics are provided in Table 1. The gene prediction on assembled scaffolds predicted 31,943 genes for KET<sub>25</sub>; 26,487 genes for KET<sub>29</sub> and 31,654 genes for KET<sub>30</sub> with average gene size of 458bp, 424bp and 459bp respectively. Predicted gene statistics are provided in Table 2.

A total of 30,209, 25,107 and 29,943 genes were annotated against the NCBI Nr database for the samples respectively. The majority of sequence similarity hits were found to be against the *Larimichthys crocea* (large yellow croaker) followed by *Stegastes partitus* (bicolor damselfish) for KET<sub>25</sub>, KET<sub>29</sub> and KET<sub>30</sub> samples (Fig. 1a, b & c). Gene annotation distribution statistics are provided in Table 3. From the gene ontology (GO) analysis, 4,068, 3616 and 4182 genes were annotated for the given samples respectively. Summary of gene ontology annotations are provided in Table 4. GO assignments were used to classify the functions of the predicted genes. The GO mapping also provides ontology of defined terms representing gene product properties which are grouped in to three main domains such as biological process, molecular function and cellular component (Fig. 2a, b & c). The GO category with the highest number of terms assigned biological processes, followed by cellular component while molecular functions had the least contigs



**Fig. 2a.** WEGO plot showing the Gene Ontology distribution in *E. tetradactylum* (sample - KET<sub>25</sub>).



**Table 5**  
SSR Prediction statistics in three samples of *E. tetradactylum*.

	KET25	KET29	KET30
Total number of scaffolds examined	289,461	280,260	289,650
Total size of examined sequences (bp)	301,704,044	263,802,638	306,253,351
Total number of identified SSRs	60,246	46,107	60,904
Number of SSR containing sequences	48,664	38,543	49,191
Number of sequences containing more than 1 SSR	8902	5950	9004
Number of SSRs present in compound formation	5565	3918	5628
Di-nucleotide repeat	15,430	9529	15,750
Tri-nucleotide repeat	40,425	33,132	40,652
Tetra-nucleotide repeat	3565	2846	3667
Penta-nucleotide repeat	668	494	666
Hexa-nucleotide repeat	158	106	169

SSR prediction statistics are represented in Table 5 In KET<sub>30</sub> sample, out of 289,650 scaffolds examined, 60,904 SSRs were identified. 23,983, 17,693 and 24,033 validated SSRs were obtained based on flanking sequence for KET<sub>25</sub>, KET<sub>29</sub> and KET<sub>30</sub> samples respectively. Based on common *in silico* validated SSR we have identified 1695 polymorphic SSR.

NGS methods can be used to address new and long-standing questions previously hindered by technological and financial limitations [2]. Microsatellite loci remain one of the most popular options for population genetic studies. The accessibility and throughput of NGS technologies has entitled the rapid and efficient microsatellite discovery by providing a greater amount of DNA sequencing reads at lower costs compared to other techniques (Irias et al., 2013). Massive and vast amounts of sequence data for a single or multiple individuals in a single run, low sequencing cost per base, reduction of the role of cloning and polymerase chain reaction (PCR) and, thus, reduced bias in resulting sequences; and the ability to identify rare variant sequences rather than a single sequence are the benefits of NGS technology when compared to the conventional capillary based sequencing [2]. NGS tools are also valuable for the discovery, validation and assessment of genetic markers in populations (Davey et al., 2011).

Accurate de novo assembly is critical for NGS projects in non-model organisms. The NGS and mining of the *E. tetradactylum* genome helped in identification of thousands of SSR markers. This vast sequenced data obtained by the de novo assembly and scaffolding can be used for developing polymorphic microsatellite markers and which will be useful in population genetic analysis and genotyping and conservation strategies in Indian salmon. The functional gene annotations and gene ontology results are useful for further gene expression studies in *E. tetradactylum*.

Here, we provide a gateway for the fishery biologist by providing NGS data to make such methods more broadly applicable and potential applications in several subfields of fishery biology. The vast data generated here in this study through NGS technology will be useful for the microsatellite marker development, gene expression studies, gene mining for novel genes and other related studies in this low volume high value species.

#### Conflicts of interest

The authors declare that they have no competing interests.

#### Data availability

The raw sequencing reads were deposited to the NCBI Sequence

Read Archive as part of Bioproject PRJNA450893. The Illumina NextSeq500 reads are available through the NCBI SRA accession number SRP141102.

#### Acknowledgments

The authors deeply thank Indian Council of Agricultural Research for providing the fund and deeply acknowledge the support and help of Dr. Prathibha Rohit, Mohammed Koya, Ms Muktha Menon, Dr. Akhilesh KV, Dr. MP Paulton, Dr Jeena NS. We thank Eurofins Genomics India Private limited (Bangalore, India) for the NGS library preparation and microsatellite identification.

#### References

- [1] J.W. Davey, M.L. Blaxter, RADSeq: next-generation population genetics, *Briefings Funct. Genomics* 9 (2010) 416–423.
- [2] H.R. Lerner, R.C. Fleischer, Prospects for the use of next-generation sequencing methods in ornithology, *Auk* 127 (2010) 4–15.
- [3] E.R. Mardis, Next-generation sequencing platforms, *Annu. Rev. Anal. Chem.* 6 (2013) 287–303.
- [4] T. Malausa, A. Gilles, E. Meglécz, H. Blanquart, S. Duthoy, C. Costedoat, N. Feau, et al., High-throughput microsatellite isolation through 454 GS-FLX Titanium pyrosequencing of enriched DNA libraries, *Mol. Ecol. Res.* 11 (2011) 638–644.
- [5] L. Liu, Y. Li, S. Li, N. Hu, Y. He, R. Pong, D. Lin, L. Lu, M. Law, Comparison of next-generation sequencing systems, *BioMed Res. Int.* 5 (2012) 2012.
- [6] T.N. Jennings, B.J. Knaus, T.D. Mullins, S.M. Haig, R.C. Cronn, Multiplexed microsatellite recovery using massively parallel sequencing, *Mol. Ecol. Resour.* 11 (2011) 1060–1067.
- [7] S.R. Narum, C.A. Buerkle, J.W. Davey, M.R. Miller, P.A. Hohenlohe, Genotyping by sequencing in ecological and conservation genomics, *Mol. Ecol.* 22 (2013) 2841–2847.
- [8] J. Wang, X. Yu, K. Zhao, Y. Zhang, J. Tong, Z. Peng, Microsatellite development for an endangered bream *Megalobrama pellegrini* (Teleostei, Cyprinidae) using 454 sequencing, *Int. J. Mol. Sci.* (2012) 3009–3021 2012 Mar 6.
- [9] R.J. Haasl, B.A. Payseur, Multi-locus inference of population structure: a comparison between single nucleotide polymorphisms and microsatellites, *Heredity* 106 (2010) 158–171.
- [10] L. Hauser, J.E. Seeb, Advances in molecular technology and their impact on fishery genetics, *Fish. Fish.* 9 (2008) 473–486.
- [11] D.A. Chistiakov, B. Hellemans, F.A. Volckaert, Microsatellites and Their Genomic Distribution, Evolution, Function and Applications: a Review with Special Reference to Fish Genetics, *Aquaculture*, 2006 2551–29.
- [12] L. Zane, L. Bargelloni, T. Patarnello, Strategies for microsatellite isolation: a review, *Mol. Ecol.* 11 (2002) 1–16.
- [13] T.C. Glenn, N.A. Schable, Isolating microsatellite DNA loci, *Methods Enzymol.* 395 (2005) 202–222.
- [14] A.I. Putman, I. Carbone, Challenges in analysis and interpretation of microsatellite data for population genetic studies, *Ecol. Evol.* 4 (2014) 4399–4428.
- [15] B.J. Darby, S.F. Erickson, S.D. Hervey, S.N. Ellis-Felege, Digital fragment analysis of short tandem repeats by high-throughput amplicon sequencing, *Ecol. Evol.* 6 (2016) 4502–4512.
- [16] L. Zhan, I.G. Paterson, B.A. Fraser, B. Watson, I.R. Bradbury, P. Nadukkalam Ravindran, ... P. Bentzen, MEGASAT: automated inference of microsatellite genotypes from sequence data, *Mol. Ecol. Resour.* 17 (2) (2017) 247–256.
- [17] I. Shihab, A. Gopalakrishnan, N. Vineesh, M. Muktha, K.V. Akhilesh, P. Vijayagopal, Histological profiling of gonads depicting protandrous hermaphroditism in *Eleutheronema tetradactylum*, *J. Fish Biol.* 90 (6) (2017) 2402–2411.
- [18] H. Motomura, Threadfins of the World. An Annotated and Illustrated Catalogue of Polynemid Species Known to Date. Family Polynemidae. FAO Species Catalogue for Fishery Purpose No. 3, Rome, (2004), p. 117.
- [19] M.B. Pember, Characteristics of Fish Communities in Coastal Waters of North-Western Australia, Including the Biology of the Threadfin Species *Eleutheronema tetradactylum* and *Polydactylus macrochir*, Doctoral dissertation, Murdoch University, 2006.
- [20] S.J. Newman, M.B. Pember, B.M. Rome, G.E. Mitsopoulos, C.L. Skepper, Q. Allsop, T. Saunders, A.C. Ballagh, L. Van Herwerden, R.N. Garrett, N.A. Gribble, Stock structure of blue threadfin *Eleutheronema tetradactylum* across northern Australia as inferred from stable isotopes in sagittal otolith carbonate, *Fish. Manag. Ecol.* 18 (2011) 246–257.
- [21] M.T. Zischke, T.H. Cribb, D. Welch, W. Sawynok, R.J. Lester, Stock structure of blue threadfin *Eleutheronema tetradactylum* on the Queensland east coast, as determined by parasites and conventional tagging, *J. Fish Biol.* 75 (2009) 156–171.
- [22] D.J. Welch, A. Ballagh, S.J. Newman, R.J. Lester, B. Moore, L. Van Herwerden, J. Horne, Q. Allsop, T. Saunders, J. Stapley, N.A. Gribble, Defining the Stock Structure of Northern Australia's Threadfin Salmon Species, (2010).