

**WHOLE GENOME SEQUENCING FOR
ESTABLISHMENT OF PHYLOGENETIC
POSITION AND DISCOVERY OF
MINIATURIZATION AND LIPID METABOLISM
GENES IN *PAEDOCYPRIS***

SAM KA KEI

UNIVERSITI SAINS MALAYSIA

2022

**WHOLE GENOME SEQUENCING FOR
ESTABLISHMENT OF PHYLOGENETIC
POSITION AND DISCOVERY OF
MINIATURIZATION AND LIPID METABOLISM
GENES IN *PAEDOCYPRIS***

by

SAM KA KEI

**Thesis submitted in fulfilment of the requirements
for the degree of
Doctor of Philosophy**

May 2022

ACKNOWLEDGEMENT

I would like to express my deepest gratitude to my supervisor, Prof. Alexander Chong Shu Chien, for his invaluable support, suggestions and guidance throughout the study. He has been giving me constructive advice and support with his expertise in the academic field. He is serious all the time but undeniable, he is a good philosopher, a good leader, and a great teacher to follow. "I cannot teach anybody anything, but I can only make them think" his strong mind, hard work, critical thinking and suggestion have given me effective training before stepping out into my career after graduation. I sincerely appreciate my best friend or senior in Chinese we called "shifu", Dr. Lau Nyok Sean, for her constant support, guidance, patience, and encouragement throughout my PhD study. She is helping me all the time in terms of research and also giving me unlimited support.

I would like to thank the assistance offered by all members in the labs: Dr. Yasmin, Lishen, Melissa, MC Teoh, Wong She Cheng, Seng Yeat, Jana, Diyana, Wong, Yeap, Shima, Liana, Diana, Jess, Fiza, Fatanah, Aliza, Fadzli, Luqman, and all the staff from the Centre for Chemical Biology. Not forget to thanks to my best friends, Leong Kar Xin and Sumaliny. All of them are very supportive and helped me complete my study successfully. I also thank Universiti Sains Malaysia for providing financial assistance during my PhD study.

I am grateful for the unconditional love and support my parents gave, without which the completion of this study would be impossible. The encouragement and trust served as a source of motivation to overcome the challenges I have encountered during my PhD study.

TABLE OF CONTENTS

ACKNOWLEDGEMENT	ii
TABLE OF CONTENTS	iii
LIST OF TABLES	viii
LIST OF FIGURES	xi
LIST OF ABBREVIATIONS AND SYMBOLS	xvi
LIST OF APPENDICES	xvii
ABSTRAK	xviii
ABSTRACT	xx
CHAPTER 1 INTRODUCTION	1
1.1 Overview	1
1.2 Problem statement	5
1.3 Research objectives	6
CHAPTER 2 LITERATURE REVIEW	7
2.1 Miniaturization	7
2.2 Peat swamp forest.....	9
2.3 <i>Paedocypris</i> species	10
2.3.1 <i>Paedocypris</i> serve as an important model fish to study	11
2.3.2 Morphological classification and phylogenetic position of <i>Paedocypris</i>	12
2.4 Genome assembly	13
2.4.1 Fish genome study.....	14
2.5 DNA sequencing	19
2.5.1 Next generation sequencing	19
2.5.2 Chromosome conformation capture techniques	22
2.6 Fatty acids	24

2.6.1	Long-chain polyunsaturated fatty acids (LC-PUFAs)	25
2.6.2	LC-PUFA biosynthetic pathway	26
2.6.3	Adaptive plasticity of LC-PUFA biosynthesis mechanisms in fish.....	28
CHAPTER 3 METHODOLOGY.....		29
3.1	Sample collection and morphological identification.....	29
3.2	Genomic DNA isolation.....	31
3.3	Experimental workflow	32
CHAPTER 4 EVOLUTION OF MINIATURIZATION AS REVEALED BY GENOME OF THE WORLD'S SMALLEST FISH PAEDOCYPRIS.....		33
4.1	Introduction	33
4.2	Materials and methods	35
4.2.1	Sample collection and genomic DNA isolation	35
4.2.2	Genome sequencing and assembly.....	35
4.2.3	Genome annotation	39
4.2.4	GC content analysis.....	40
4.2.5	Repeat analysis	41
4.2.6	<i>Hox</i> gene identification	41
4.2.7	Phylogenetic analysis	42
4.2.8	Macrosynteny analysis	43
4.2.9	Gene family analysis	43
4.2.10	Putative gene loss analysis	44
4.2.11	Data availability	44
4.3	Results	45
4.3.1	Genome sequencing, assembly and annotation.....	45
4.3.2	Comparisons with other fish genomes	61
4.3.2(a)	Synteny analyses.....	61
4.3.2(b)	Gene characteristics	65

4.3.3	Gene family analysis	67
4.3.3(a)	Discovery of gene loss in <i>Paedocypris</i>	68
4.3.3(b)	Homeobox genes and clusters	70
4.3.3(c)	SCPP gene cluster.....	72
4.4	Discussion	74
4.4.1	Genome streamlining in the <i>Paedocypris</i>	74
4.4.2	Evolution of <i>Hox</i> genes and clusters	75
4.4.3	SCPP gene repertoires	77
4.4.4	Loss of <i>bglap</i>	80
4.4.5	Loss of <i>spp2, fgf8b, disc1, tbx22</i>	81
4.4.6	Loss of <i>mtmr14, ccdc78, marcksa</i>	84
4.5	Summary	87
CHAPTER 5 COMPLETE MITOCHONDRIAL GENOME OF PAEDOCYPRIS MICROMEGETHES AND P. CARBUNCULUS REVEAL CONSERVED GENE ORDER AND PHYLOGENETIC RELATIONSHIPS OF MINIATURIZED CYPRINIDS.....		88
5.1	Introduction	88
5.2	Material and Methods.....	90
5.2.1	Sample collection, DNA isolation and mitochondrial genome sequencing.....	90
5.2.2	Mitochondrial genome assembly, PCR amplification and sequencing.....	91
5.2.3	Mitochondrial genome annotation and analyses	93
5.2.4	Phylogenetic analysis	95
5.3	Results	96
5.3.1	PCR amplification of the mitochondrial genome.....	96
5.3.2	Mitochondrial genome structure and composition.....	97
5.3.3	Base composition bias of the mitochondrial genome.....	103
5.3.4	Protein-coding gene features and codon usage	106

5.3.5	Analysis of PCGs evolutionary rate	112
5.3.6	Transfer RNAs and ribosomal RNAs.....	113
5.3.7	Control region and O _L	116
5.3.8	Gene arrangement	118
5.3.9	Phylogenetic analysis	120
5.4	Discussion	124
5.5	Summary	131
CHAPTER 6 DE NOVO TRANSCRIPTOME ASSEMBLY AND FUNCTIONAL CHARACTERIZATION OF THE LONG-CHAIN POLYUNSATURATED FATTY ACID BIOSYNTHETIC PATHWAYS IN PAEDOCYPRIS MICROMEGETHES.....		132
6.1	Introduction	132
6.2	Material and Methods.....	135
6.2.1	Sample collection, total RNA extraction, and sequencing.....	135
6.2.2	<i>De novo</i> assembly, filtering, and annotation	135
6.2.3	Validation, identification and phylogenetic analysis of the fatty acid genes	137
6.2.4	cDNA synthesis and molecular cloning of the <i>fads2</i> , <i>elovl5</i> , <i>elovl2</i> , <i>elovl4a</i> , and <i>elovl4b</i>	138
6.2.5	Functional characterization of the <i>fads2</i> , <i>elovl5</i> , <i>elovl2</i> , <i>elovl4a</i> , and <i>elovl4b</i>	140
6.2.6	Lipid extraction, FAME preparation and gas chromatographic...	141
6.3	Results	143
6.3.1	<i>De novo</i> transcriptome assembly and annotation.....	143
6.3.2	Sequence and phylogenetic analyses of the Fads and Elovl enzymes in <i>P. micromegethes</i>	149
6.3.3	<i>In-vitro</i> functional characterization of Fads2 desaturase and Elovl elongases from <i>P. micromegethes</i>	158
6.4	Discussion	177
6.5	Summmary	185

CHAPTER 7	SUMMARY OF FINDINGS, CONCLUSIONS AND FUTURE STUDIES.....	186
7.1	Summary of findings and conclusions	186
7.2	Future study.....	188
REFERENCES		189
APPENDICES		
LIST OF PUBLICATIONS		

LIST OF TABLES

		Page
Table 2.1	Summarize of the fish genome study.....	18
Table 2.2	Sequencing approach utilized in the fish mitochondrial genome assembly.....	20
Table 2.3	Sequencing approach applied for fish chromosome-level genome assembly.....	23
Table 4.1	Summary of genomic sequencing data.	46
Table 4.2	Summary of <i>P. micromegethes</i> and <i>P. carbunculus</i> genome assemblies.	46
Table 4.3	Mapped reads to the genome assembly.....	50
Table 4.4	Evaluation of gene coverage using RNA-seq assembled transcripts.	50
Table 4.5	Assessment of assembled genomes by BUSCO and CEGMA.	51
Table 4.6	Repeat composition in <i>P. micromegethes</i> , <i>P. carbunculus</i> and other teleost genomes.....	53
Table 4.7	Number of gene models with homology or functional classification.....	55
Table 4.8	Summary of assembly statistics for the genomes of <i>P. micromegethes</i> and <i>P. carbunculus</i>	56
Table 4.9	Non-coding RNAs in <i>P. micromegethes</i> and <i>P. carbunculus</i> genomes.....	60
Table 4.12	Gene characteristics of <i>P. micromegethes</i> , <i>P. carbunculus</i> and the comparison with other teleost	66
Table 4.13	Gene family size evolution analysis using CAFÉ (Han et al., 2013)	67

Table 5.1	Primers used to verify the accuracy of the assembled mitochondrial genome sequence	92
Table 5.2	Summary on the annotation of the <i>P. micromegethes</i> and <i>P. carbunculus</i> mitochondrial genomes	100
Table 5.3	Base composition and skewness of <i>P. micromegethes</i> and <i>P. carbunculus</i> mitochondrial genomes	104
Table 5.4	Whole mitogenome composition and skewness in Danionidae species	105
Table 5.5	Summary for the base composition and skewness of the mitochondrial genomes in Danionidae species	108
Table 6.1	Primers used for cloning of full-length open reading frame (ORF) of miniature fish <i>P. micromegethes fads2, elovl5, elovl2, elovl4a,</i> and <i>elovl4b.</i>	139
Table 6.2	General statistic for the <i>de novo</i> assembly Trinity of miniature fish <i>Paedocypris micromegethes</i>	145
Table 6.3	List of KEGG annotated transcripts in the biosynthesis of unsaturated fatty acids for miniature fish <i>P. micromegethes</i>	147
Table 6.4	Functional characterization of miniature fish <i>P. micromegethes</i> Fads2 through heterologous expression of the respective transcript ORF in yeast (<i>Saccharomyces cerevisiae</i>), followed by incubation with PUFA substrate to assay the desaturation rate leading to PUFA product.	160
Table 6.5	Functional characterization of miniature fish <i>P. micromegethes</i> Elov15 and Elov12 through heterologous expression of the respective transcript ORF in yeast (<i>Saccharomyces cerevisiae</i>), followed by incubation with PUFA substrate to assay the elongation rate leading to PUFA product.....	161
Table 6.6	Functional characterization of miniature fish <i>P. micromegethes</i> Elov14a and Elov14b through heterologous expression of the respective transcript ORF in yeast (<i>Saccharomyces cerevisiae</i>),	

followed by incubation with SFA/PUFA substrate to assay the
elongation rate leading to SFA/PUFA product 162

LIST OF FIGURES

		Page
Figure 2.1	Combination of long inserts from mate pair sequencing with that from short-insert paired-end sequencing for <i>de novo</i> assembly.....	21
Figure 2.3	Overview of polyunsaturated fatty acids synthesis pathways in fish, demonstrating the involvement of Fads2 desaturase and elongases (Elov15, Elov12, and Elov14) in endogenous LC/VLC-PUFAs production.	27
Figure 3.1	Sample collection and identification of <i>Paedocypris</i> spp. for this current research.	30
Figure 3.2	Experimental workflow.....	32
Figure 4.1	The workflows applied in <i>P. micromegethes</i> and <i>P. carbunculus</i> genome assemblies in this work.....	38
Figure 4.2	Distribution of GC contents in <i>P. micromegethes</i> , <i>P. carbunculus</i> , <i>D. rerio</i> and <i>C. carpio</i> genomes. Sliding windows of 10 Kb were used.	48
Figure 4.3	Genome size estimation based on <i>k</i> -mer analysis.....	49
Figure 4.4	Gene ontology level two classification of <i>P. micromegethes</i> genes.	58
Figure 4.5	Gene ontology level two classification of <i>P. carbunculus</i> genes.....	59
Figure 4.6	Dot-plot analyses between <i>P. micromegethes</i> (pm) and (a) <i>D. rerio</i> (dr), (b) <i>C. carpio</i> (cc), (c) <i>C. auratus</i> (ca) and (d) <i>L. waleckii</i> (lw) chromosomes.	62
Figure 4.7	Macro-synteny between <i>P. micromegethes</i> (pm) and (a) <i>D. rerio</i> (dr), (b) <i>C. carpio</i> (cc), (c) <i>C. auratus</i> (ca) and (d) <i>L. waleckii</i> (lw)	63
Figure 4.8	Dot-plot analyses between (a) <i>P. micromegethes</i> (pm) chromosomes and (b) <i>P. carbunculus</i> (pc) scaffolds.	64

Figure 4.9	Relationships between genome size, average intron length and repeat content for selected teleost species.....	66
Figure 4.10	The summary of the important gene loss identified in the <i>Paedocypris</i>	69
Figure 4.11	Maximum likelihood phylogeny inferred from 188 single-copy orthologs.....	71
Figure 4.12	SCPP gene clusters in the <i>P. micromegethes</i> and other teleost.	73
Figure 5.1	Analysis of PCR amplification products from the mitochondrial genome of <i>P. micromegethes</i> and <i>P. carbunculus</i>	96
Figure 5.2	Circular genetic map and organization for the <i>P. micromegethes</i> mitochondrial genome.....	98
Figure 5.3	Circular genetic map and organization for the <i>P. carbunculus</i> mitochondrial genome.....	99
Figure 5.4	Amino acid composition in the mitogenomes of <i>P. micromegethes</i> and <i>P. carbunculus</i>	110
Figure 5.5	Relative synonymous codon usage of protein-coding genes in <i>P. micromegethes</i> and <i>P. carbunculus</i> mitogenomes.	111
Figure 5.6	The ratio of nonsynonymous and synonymous substitution (Ka/Ks) estimated in all 13 mitochondrial protein-coding genes of <i>Paedocypris</i> species	112
Figure 5.7	Predicted secondary structure for 22 tRNA genes in the <i>P. micromegethes</i> (Pm) and <i>P. carbunculus</i> (Pc) mitochondrial genomes as predicted by Mitos and visualized in Forna web server	115
Figure 5.8	Control region features of the (A) <i>P. micromegethes</i> and (B) <i>P. carbunculus</i> mitogenomes. Arrows below the nucleotide sequences indicate repeat region.....	117
Figure 5.9	Comparison of the mitochondrial gene arrangement in representative miniature fish and cyprinid species.	119

Figure 5.10	Phylogenetic tree inferred from concatenated nucleotide sequences of 13 protein-coding genes and two rRNA genes of the mitogenomes using maximum likelihood method.	122
Figure 5.11	Phylogenetic tree inferred from concatenated nucleotide sequences of 13 protein-coding genes and two rRNA genes of the mitogenomes using Bayesian method.	123
Figure 6.1	Analyses performed on the transcriptome assembly of <i>P. micromegethes</i>	145
Figure 6.2	Level two classification of the gene ontology functional annotation for <i>P. micromegethes</i> transcripts.	146
Figure 6.3	Comparison of the deduced amino acid (aa) sequence of the newly cloned fatty acyl desaturase 2 (Fads2) of the miniature fish <i>P. micromegethes</i> (MZ352750) and orthologues from <i>Homo sapiens</i> (AAD20018), <i>Clarias gariepinus</i> (AMR43366), <i>Barbonymus gonionotus</i> (AXF92413), <i>Danio rerio</i> (AAG25710), <i>Tinca tinca</i> (QIA97820), and <i>Tor tambroides</i> (AZL94116) through alignment using Clustal omega.	153
Figure 6.4	Maximum likelihood phylogenetic analysis inferred from deduced amino acid sequence of the miniature fish <i>P. micromegethes</i> Fads2 with Fads1 and Fads2 orthologs from other vertebrate species rooted with invertebrate group.	154
Figure 6.5	Comparison of the deduced amino acid (aa) sequence of the newly cloned elongases (Elovl5) of the miniature fish <i>P. micromegethes</i> (Elovl5- MZ352751, Elovl2- MZ352752, Elovl4a- MZ352753, and Elovl4b- MZ352754) and orthologues from <i>Homo sapiens</i> (Elovl5-NP_068586, Elovl2-NP_060240, and Elovl4-NP_073563), <i>Clarias gariepinus</i> (Elovl4a-ASY01350 and Elovl4b-ASY01351), <i>Danio rerio</i> (Elovl5-AAN77156, Elovl2-NP_001035452, Elovl4a-NP_957090, and Elovl4-NP_956266), <i>Barbonymus gonionotus</i> (Elovl5-AXG50646, Elovl2-AZN23179, Elovl4-QEE04393), <i>Tor tambroides</i> (Elovl5-QJU12162 and	

	Elov12-QJU12161), and <i>Misgurnus anguillicaudatus</i> (Elov14a-QFR04622) through alignment using Clustal Omega.....	156
Figure 6.6	Phylogenetic tree inferred from deduced amino acid sequence of the miniature fish <i>P. micromegethes</i> Elov15, Elov12, Elov14a, and Elov14b with orthologs from other vertebrate species using maximum likelihood analysis.	157
Figure 6.7	Functional characterization of the <i>P. micromegethes</i> Fads2 in transgenic yeast (<i>Saccharomyces cerevisiae</i>) cultured in one of the exogenously added fatty acid (FA) substrates: $\Delta 6$: 18:3n-3 (B), 18:2n-6 (D), 24:5n-3 (R) and 24:4n-6 (T); $\Delta 8$: 20:3n-3 (F) and 20:2n-6 (H); $\Delta 5$: 20:4n-3 (J) and 20:3n-6 (L); and $\Delta 4$: 22:5n-3 (N) and 22:4n-6 (P).....	166
Figure 6.8	Functional characterization of the <i>P. micromegethes</i> Elov15 in transgenic yeast (<i>Saccharomyces cerevisiae</i>) cultured in one of the exogenously added fatty acid (FA) substrates: 18:3n-3 (B), 18:2n-6 (D), 18:4n-3 (F), 18:3n-6 (H), 20:5n-3 (J), 20:4n-6 (L), 22:5n-3 (N) and 22:4n-6 (P).	168
Figure 6.9	Functional characterization of the <i>P. micromegethes</i> Elov12 in transgenic yeast (<i>Saccharomyces cerevisiae</i>) cultured in one of the exogenously added fatty acid (FA) substrates: 18:3n-3 (B), 18:2n-6 (D), 18:4n-3 (F), 18:3n-6 (H), 20:5n-3 (J), 20:4n-6 (L), 22:5n-3 (N) and 22:4n-6 (P).	170
Figure 6.10	Functional characterization of the <i>P. micromegethes</i> Elov14a in transgenic yeast (<i>Saccharomyces cerevisiae</i>) cultured in one of the exogenously added fatty acid (FA) substrates: 24:0 (B), 18:3n-3 (D), 18:2n-6 (F), 18:4n-3 (H), 18:3n-6 (J), 20:5n-3 (L), 20:4n-6 (N), 22:5n-3 (P), 22:4n-6 (R), and 22:6n-3 (T).....	173
Figure 6.11	Functional characterization of the <i>P. micromegethes</i> Elov14b in transgenic yeast (<i>Saccharomyces cerevisiae</i>) cultured in one of the exogenously added fatty acid (FA) substrates: 24:0 (B), 18:3n-3 (D), 18:2n-6 (F), 18:4n-3 (H), 18:3n-6 (J), 20:5n-3 (L), 20:4n-6 (N), 22:5n-3 (P), 22:4n-6 (R), and 22:6n-3 (T).....	176

Figure 6.12 Proposed biosynthetic pathways of LC-PUFA (C₂₀₋₂₄) and VLC-PUFA (>C₂₄) from linolenic (18:3n3) and linoleic (18:2n6) acids in miniature fish *P. micromegethes* by *in-vitro* characterization of Fads2 desaturase (“Δx” coloured in orange) and Elovl elongases (“Elox” coloured in green) capacities using yeast. 184

LIST OF ABBREVIATIONS AND SYMBOLS

%	Percentage
°C	Degree Celsius
Δ	Delta
α	Alpha
β	Beta
g	gram
mL	Millilitre
μg	Microgram
v/v	Volume per volume
C	Carbon
cDNA	Complementary deoxyribonucleic acid
DEPC	Diethylpyrocarbonate
dH ₂ O	Distilled water
DHA	Docosahexaenoic acid
gDNA	Genomic deoxyribonucleic acid
dNTP	Deoxyribonucleoside triphosphate
ALA	α – linolenic acid
ARA	Arachidonic acid
ATP	Adenosine triphosphate
BF ₃	Boron trifluoride
EFA	Essential fatty acid
Elovl	Eicosapentaenoic acid
FA	Fatty acid
Fads	Fatty acid desaturase
FAME	Fatty acid methyl ester
GC-MS	Gas chromatography–mass spectrometry
LA	Linoleic acid
LB	Luria-bertani
LC-PUFA	Long-chain polyunsaturated fatty acid
LDL	Low density lipoprotein
mRNA	Messenger RNA
MUFA	Monounsaturated fatty acid
n-3	Omega-3
n-6	Omega-6
OD	Optical density
ORF	Open reading frame
PCR	Polymerase chain reactions
RNA	RNA
SCMM-U	<i>S. cerevisiae</i> minimal medium plates without uracil
VLC-PUFA	Very long-chain polyunsaturated fatty acid
VLC-SFA	Very long-chain saturated fatty acid
YPD	Yeast extract peptone dextrose

LIST OF APPENDICES

Appendix A	gDNA sample quality assess using gel electrophoresis, clear bands with no smearing were obtained before sending for library preparation and sequencing
Appendix B	Hi-C dataset for <i>P. micromegethes</i> chromosome-scale assembly
Appendix C	Hi-C interaction heat maps of <i>P. micromegethes</i>
Appendix D	Summary of the <i>P. micromegethes</i> Hi-C assisted assembly
Appendix E	RNA ladder and sample result using an Agilent 2100 Bioanalyzer (Agilent Technologies, USA)
Appendix F	Transformation of plasmid DNA into competent yeast
Appendix G	Total number of KEGG annotated pathways with transcripts for the miniature fish <i>P. micromegethes</i>
Appendix H	Biosynthesis of unsaturated fatty acids pathway

**PENJUJUKAN KESELURUHAN GENOM UNTUK PENETAPAN
KEDUDUKAN FILOGENETIK DAN PENEMUAN GEN PENGKERDILAN
DAN METABOLISME LIPID DALAM *PAEDOCYPRIS***

ABSTRAK

Paedocypris adalah merupakan antara vertebrata terkecil di dunia dengan jumlah panjang ≤ 20 mm yang terhad kepada persekitaran hutan paya gambut yang berasid. Anatomi *Paedocypris* telah banyak diubahsuai dengan ciri perkembangan yang terpenggal. Fenotip kerdil yang unik menjadikan *Paedocypris* sebagai model yang sesuai untuk memahami perkembangan dan evolusi. Tiga *Paedocypris* spesies telah dihuraikan, termasuk *P. micromegethes* dari Sarawak dan dua spesies dari Indonesia, *P. progenetica* dan *P. carbunculus*. Setakat ini, *Paedocypris* tidak mempunyai genom skala kromosom, dan hanya draf genom dengan perancah N50 ~59-61 kb. Kajian ini telah menghasilkan genom bertahap kromosom berkualiti tinggi untuk *P. micromegethes* dengan perancah N50 saiz 29.8 Mb menggunakan gabungan teknologi penjujukan Illumina, PacBio dan Hi-C. Sementara itu, draf genom untuk *P. carbunculus* telah ditambah baik dengan panjang perancah N50 sebanyak 125 Kb. Kajian ini mendedahkan bahawa saiz genom yang kecil dan padat *Paedocypris* berciri kandungan ulangan yang rendah, intron pendek, dan kekurangan duplikasi genom keseluruhan. Sebagai tambahan, kajian ini juga menemui kehilangan gen daripada genom *Paedocypris*, seperti *fgf8b*, *mtmr14*, *scpp3a*, *scpp3b*, *scpp6*, *scpp7*, *scpp8*, *gsp37*, *spp2*, dan *tbx22*. Analisis filogenetik untuk menyelesaikan pertikaian kedudukan taksonomi ikan kerdil cyprinid dengan menggunakan 13 gen pengekodan protein (PCG) dan dua gen RNA ribosom (rRNA) daripada mitokondria menunjukkan nilai sokongan tinggi untuk hubungan kumpulan beradik bagi *Paedocypris* dan

Sundadanio kepada kelompok keluarga Danionidae. Perbezaan kedudukan dalam topologi pokok untuk ikan kerdil mencadangkan bahawa mereka tidak membentuk klad monofiletik. Kebolehan biosintesis asid lemak rantai panjang politaktepu dalam *P. micromegethes* juga telah dikaji. Walaupun saiz genom *Paedocypris* telah mengalami penyusutan, bilangan salinan gen desaturase dan elongases yang dijumpa dalam perhimpunan data genom dan transkrip adalah sama dengan *Danio rerio*. Tambahan pula, keputusan fungsi *in-vitro* menunjukkan bahawa *P. micromegethes* mempunyai keupayaan untuk biokonversi C18:2n-6 dan C18:3n-3 kepada asid arakidonik (ARA, 20:4n-6), asid eicosapentaenoic (EPA, 20:5n-3) dan asid docosahexaenoic (DHA, 22:6n-3), masing masing. Pengekalan laluan biosintesis LC-PUFA yang lengkap dalam *Paedocypris* adalah untuk memastikan penjajahan dan adaptasi dalam persekitaran akuatik yang berasid.

**WHOLE GENOME SEQUENCING FOR ESTABLISHMENT OF
PHYLOGENETIC POSITION AND DISCOVERY OF MINIATURIZATION
AND LIPID METABOLISM GENES IN *PAEDOCYPRIS***

ABSTRACT

Paedocypris is the world's smallest vertebrate with a total length of ≤ 20 mm restricted to acidic freshwater peat swamp environments. The anatomy of *Paedocypris* is highly modified with developmental truncated characters. The unique miniature phenotypes make the *Paedocypris* a good model fish to study development and evolution. Three *Paedocypris* species have been described, including *P. micromegethes* from Sarawak and another two species, *P. progenetica* and *P. carbunculus* from Indonesia. Thus far, *Paedocypris* has no chromosome-scale genome, with only a draft genome with scaffold N50 ~59-61 kb. In this study, a high-quality chromosome level genome assembly for the *P. micromegethes* was generated with improved scaffold N50 of 29.8 Mb using a combination of Illumina, PacBio and Hi-C sequencing technologies. Meanwhile, the draft genome for the *P. carbunculus* was improved with N50 scaffold length of 125 Kb. This study revealed that the small and compact genome sizes of miniature *Paedocypris* are characterized by a low repeat content, short introns, and a lack of recent genome duplications. Additionally, this study discovered gene loss from the improved genome assemblies, such as *fgf8b*, *mtmr14*, *scpp3a*, *scpp3b*, *scpp6*, *scpp7*, *scpp8*, *gsp37*, *spp2*, and *tbx22*. Phylogenetic analysis to resolve the contentious taxonomy position of miniaturized cyprinids using mitochondrial 13 protein-coding genes (PCGs) and two ribosomal RNA genes (rRNAs) showed improved supportive values for the sister-group relationship for *Paedocypris* and *Sundadanio* to the Danionidae family. A different

position for the miniature fish within the tree topology suggested that they do not form monophyletic clade. The capacity for long-chain polyunsaturated fatty acid (LC-PUFAs) biosynthesis in *P. micromegethes* was also studied. Identical copy numbers of desaturase and elongases genes with *Danio rerio* were identified in the *Paedocypris* genome and transcriptome assemblies, despite its genome size reduction. *In-vitro* functional study results revealed that the *P. micromegethes* has the ability for bioconversion of C18:2n-6 and C18:3n-3 to arachidonic acid (ARA, 20:4n-6), eicosapentaenoic acid (EPA, 20:5n-3), and docosahexaenoic acid (DHA, 22:6n-3), respectively. The retaining of a complete LC-PUFA biosynthesis pathway in this species most likely ensures successful colonization and adaptation in the acidic aquatic environment.

CHAPTER 1

INTRODUCTION

1.1 Overview

The evolution of miniaturization is observed in all vertebrate lineages. However, among the miniature organisms, *Paedocypris* comprised the world's smallest vertebrate. Therefore, to understand the evolution of miniaturization, the two best candidates, *P. micromegethes* and *P. carbunculus*, were selected for this study. *Paedocypris* species are restricted to the highly acidic (pH 3~5) dark-coloured freshwater peat swamp forest distributed in Southeast Asia. The acidic freshwater ecosystem in the peat swamp is depleted in nutrients but high in tannins (Yule and Gomez, 2009;Page and Rieley, 2016). Interestingly, adaptive evolution has emerged specialized life forms to adapt to this harsh aquatic environment, and most of them are restricted to this habitat (Leete, 2006;Posa et al., 2011). Acidic peat swamp forest consists of many miniature creatures, including fish (Kottelat et al., 2006). *Paedocypris* represents a phenotypic extreme in adult body size reduction among vertebrates and is classified as the smallest vertebrate known in the world. They have a sexually matured body size for an adult fish that can reach only 7.9-11.5 mm standard length (Kottelat et al., 2006;Britz and Kottelat, 2008). Despite body size reduction, *Paedocypris* possesses developmentally truncated characters, resulting in massive bones loss and reduction (Britz and Conway, 2009). Owing to its simplified morphology and remarkable biological characteristics, including its well-adaptation in

such a hostile environment, *Paedocypris* serve as an excellent model for the studies of evolution, development, and acidic freshwater adaptation.

Sequencing the genome is an important task that can help provide a clue to scientists on the gene functions and how genes work together to direct an organism's growth, development and maintenance through analyzing the sequenced genome. The advent of next-generation sequencing technology has enabled scientists to discover the organism's genetic information by sequencing its whole genome and now can even sequence up to a chromosome level. Nevertheless, the paucity of genome information of this tiny teleost fish has hindered understanding of its ecology and evolution. Furthermore, over-exploitation of peat swamp forests has threatened the survival of the fish population, urging the need for a genomic study to conserve the genomic information before the fish is going to be extinct. To date, a draft whole-genome sequence of *Paedocypris* (*P. micromegethes* and *P. carbunculus*) was released. Nevertheless, the employment of Illumina short-read sequencing technology has led to the creation of highly fragmented genome assembly with scaffold N50 of ~59-61 Kb (Malmstrøm et al., 2018). Without a high-quality reference genome and the lack of completeness and continuity of the current *Paedocypris*'s draft genome may become one of the most critical factors limiting its application for in-depth downstream analysis. Therefore, it is essential to improve the genome assembly and create a model reference genome for *Paedocypris* species to get more detailed information on their genome architectures and gene contents to understand the genotype-phenotype relationships towards the evolution of miniaturization.

Mitochondria DNA is the small circular chromosome found in all eukaryotes and serves as an essential component responsible for supplying energy for all eukaryotes (Saraste, 1999). This membrane-bound organelle is an effective genetic

marker in phylogeny, taxonomy and biological evolution analyses, owing to its maternal inheritance, simple genetic structure, easy detection, low recombination frequencies and rapid evolutionary rate (Miya et al., 2003; Gissi et al., 2008). Mitochondria genome (mitogenome) is a popular DNA molecule that has been successfully employed in studying the phylogenetic taxonomy among fishes (Miya and Nishida, 2015). However, the mitogenome of the miniature fish is still less well-understood, with only a few *Paedocypris* mitogenomes have been deposited in the GenBank database. The inconsistencies of the *Paedocypris* phylogenetic position arise from contradicting morphological identification and molecular analyses. No study has been done using a complete mitochondria gene set to examine the phylogenetic position of this miniaturized fish species among cyprinids. However, this method has been successfully applied to improve the controversial phylogenetic position of fish (Sun et al., 2021). Therefore, it is essential to assemble and study the mitogenome of *P. micromegethes* and *P. carbunculus* to provide a good resolution into its features, evolution and phylogenetic relationships among miniature fish and other cyprinids.

Long-chain polyunsaturated fatty acids (LC-PUFAs; \geq C20), consisting of eicosapentaenoic acid (EPA; 20:5n-3), docosahexaenoic acid (DHA; 22:6n-3), and arachidonic acid (ARA; 20:4n-6), are important compounds for normal cell functions in vertebrates, including humans. These essential LC-PUFAs play vital roles in manipulating the cell's membrane phospholipid composition, modulating gene expression, regulating eicosanoid synthesis, and as cell-signaling molecules (Jump, 2002; Schmitz and Ecker, 2008). Due to the absence of D12 and D15 desaturases, fish, like all vertebrates, are incapable of *de novo* synthesizing linolenic acid (LA) and α -linolenic acid (ALA) and are solely dependent on dietary intake. Limited or inability of marine fish species in LC-PUFA biosynthesis has been hypothesized to adapt to the

DHA-rich marine environment. In contrast, enzymatic functions required for C18 PUFAs to LC-PUFAs conversion are conserved in most fish species living in DHA-deficient freshwater environments (Tocher, 2010). Studies on fish suggested that evolutionary history, environmental factors (habitats), trophic levels, and ecological backgrounds can influence the fatty acid biosynthesis capability (Kabeya et al., 2017; Xie et al., 2020). The miniature *Paedocypris* species mainly depend on planktonic rotifers and cladocerans for food (Kottelat et al., 2006). *Paedocypris* species are adapted well to this acidic freshwater environment. They might possess a functional LC-PUFA biosynthesis capacity to cope with the stress environment with limited food resources for LC-PUFA production. Thus, it is interesting to investigate the *Paedocypris* LC-PUFA biosynthesis pathways from the genome, and transcriptomic assemblies and then functional characterize the isolated fatty acid genes to provide an insight into how they colonize and adapt this acidic aquatic environment.

1.2 Problem statement

Understanding the genome features and genotype-phenotype relationship is vital for answering the molecular basis of evolutionary history and diversification of life; however, the current *Paedocypris* draft genomes are fragmented, and part of the genome sequences are missing or consisting of gaps. Thus, improving the contiguity of the genome assembly is essential to capture complete genetic information for deeper insights into the characteristics of the genome, the evolution of miniaturization and acidic adaptation, which can tell us much about the organisms in question. In this study, Malaysia's species *P. micromegethes* was selected to sequence up to a chromosomal-scale assembly using Illumina, PacBio and Hi-C sequencing techniques. Meanwhile, the contiguity of the *P. carbunculus* draft genome was improved to examine the genotype difference between the two miniature *Paedocypris* species.

Besides that, the mitogenome of the miniature fish is still less well-understood. Limited mitogenome resources are reported in the GenBank database, encouraging the assembly of the complete mitogenome from different *Paedocypris* species. *Paedocypris* has reduced genome size; thus, examining the miniaturization effects toward the mitogenome features is interesting. Mitogenome is a good marker for studying the phylogenetic relationships among organisms. Therefore, a complete set of mitochondrial genes was applied in this study to examine the phylogeny topology of miniature *Paedocypris* to explore the evolutionary relationships between different fish species.

Furthermore, the capability of LC-PUFA biosynthesis in *Paedocypris* is still unknown. This is due to the non-discovery of *fads* and *elovl* genes in this miniature species. Therefore, it is important to investigate the capacity of LC-PUFA in *Paedocypris* species based on the assembled genome and transcriptome to understand its LC-PUFA biosynthesis capability. Acquiring this knowledge could provide fundamental insight into how fatty acid biosynthesis capacity contributes to acidic freshwater adaptation in organisms.

1.3 Research objectives

The objectives of this study are stated below:

- To improve genome assemblies for *P. micromegethes* and *P. carbunculus*.
- To identify potential factors contributing to the evolution of miniaturization.
- To generate *Paedocypris* complete mitogenomes and determine their phylogenetic position.
- To functionally characterize putative *fads2* and *elovls* genes of *Paedocypris* in adaptation to the acidic peat swamp.

CHAPTER 2

LITERATURE REVIEW

2.1 Miniaturization

Miniaturization, the evolutionary reduction of adult body size, occurs among vertebrate lineages such as fish, amphibians, reptiles, mammals (Hanken and Wake, 1993) and within invertebrates, including insects (Polilov, 2015). Miniaturization is commonly associated with novel ecology, physiology and morphology (Hanken and Wake, 1993). The miniaturized phenotypes are unique because of the complex derived traits with reduction and structural simplification not developed in their larger relatives or ancestors. The consequences of miniaturization for the organism in physical (obvious) features are reduction and structural simplification, morphological novelty and increased morphological variability, whereas biologically features (nonobvious) are genome and cell size. Miniaturization has emerged as an adaptive evolution for living organisms to survive under critical environments. The evolution of miniaturization has influenced organismal physiology, ecology, life history, and behaviour (Hanken and Wake, 1993). Dwarfism can help organisms avoid predators, obtain new food sources, reduce intra-and interspecific competitions, invade new ecological niches or habitats and promote sex maturity at an early age (Hanken and Wake, 1993;Blanckenhorn, 2000). Extreme body size reduction is reported in teleosts, where species maturing at sizes less than 20 mm are defined as miniature fishes (Weitzman and Vari, 1988). Miniaturization has independently evolved several times within the bony fish from family Cyprinidae, where abundant miniature cyprinids are under subfamily Danioninae, the zebrafish relatives, such

as *Boraras*, *Danionella*, *Horadandia*, *Microboraras*, *Sundadanio*, and *Paedocypris* (Rüber et al., 2007). In addition, there are some miniature cyprinids from the subfamily Cyprininae, including African genera “*Barbus*” and *Barboides* and South Asian Sawbwa (Conway et al., 2017). Cyprinidae (Teleostei: Ostariophysi: Cypriniformes) is the most diverse vertebrate family consisting mainly of freshwater fishes with more than 3,000 species distributed in Africa, Eurasia, and North America (Nelson et al., 2016;Fricke et al., 2018). Two main phenotypic forms of dwarfism were observed and distinguished in danionine cyprinids: (1) proportionate dwarfism, tiny at the adult stage but almost identical copies to their larger relatives; and (2) developmental truncation, the acceleration of maturation, leading to dwarf sexually matured adults (progenetic miniatures), resembling the juvenile/larvae stage (an early developmental stage) of their larger relatives. These two miniature groups are similar in size range, but the proportioned dwarfs exhibit few reductions compared to their large ancestors, whereas developmentally truncated miniatures possess numerous losses (Rüber et al., 2007;Britz and Conway, 2009;Britz et al., 2014). Miniaturized cyprinids with remarkably developmental truncation features are *Barboides*, *Danionella*, *Sundadanio*, and *Paedocypris* (Rüber et al., 2007;Conway et al., 2017). Generally, reduction and structural simplifications are common among miniature cyprinids. However, morphological novelties are only exhibited in miniature species with a high degree of developmental truncation, and the majority are sexually dimorphic (Rüber et al., 2007;Britz, 2009;Britz and Conway, 2009). Thus, there is a possible relationship between developmental truncation and the novelty in morphological evolution (Britz, 2009;Britz and Conway, 2009;Britz et al., 2009).

2.2 Peat swamp forest

Peat swamp forest (PSF) is formed by accumulating partially decomposed dead leaves and wood over time under waterlogged conditions. This forest is one of the unique and harsh ecosystems in the tropical rainforest biome distributed extensively in Southeast Asia. About 60% of the world's peat swamp forests are distributed in Southeast Asia, and Malaysia is the second-largest country where the peat swamp forests are localized (Posa et al., 2011). PSF is also the most endangered ecosystem globally, under legal and illegal logging, fire and land conversion. Nevertheless, the biodiversity in this forest is still poorly understood. Peat swamp, also called "black water", is highly acidic (pH 3 to 5). The water is almost dark in colour due to the high concentration of humic acids and other phenolic compounds like tannins in this waterlogged environment. Peat swamp is a poor nutrient environment since leaf litters are the main build-up for the peat, and the water was attributed to the rainfall with no inflow of nutrients from either stream or river (Yule and Gomez, 2009). Highly acidic, low nutrients and poor light conditions (primarily due to the dark brown colour of water and shared with the forest canopy) have restricted algae or plankton growth in this harsh environment (Wehr and Sheath, 2003; Yule, 2008; Yule and Gomez, 2009). Surprisingly, this unique acidic ecosystem supports many flora and fauna despite unfavourable conditions like low pH, anaerobic, and poor nutrients (Yule, 2008). PSF endemism is likely due to the inhabitant of organisms in rigorous microhabitats provided by the peat. Adaptive evolution enables them to uniquely adapted to this stressed environment (Posa et al., 2011). Most of the fish species found in PSF are mainly under the Cyprinidae family, followed by Osphronemidae, Bagridae and Siluridae (Sule et al., 2016). Meanwhile, a high number of miniature fish species found in the aquatic ecosystem of PSF are under the Cyprinidae family (Kottelat et al., 2006),

such as *Boraras* (Conway and Kottelat, 2011), *Sundadanio* (Conway et al., 2011), *Fangfangia* (Britz et al., 2011), and the smallest recorded fish belong to *Paedocypris* species (Kottelat et al., 2006). In PSF, small size is advantageous for fish, as they can survive droughts in shallow pools with low water levels, animal burrows, or soil (Kottelat et al., 2006).

2.3 *Paedocypris* species

Paedocypris (Teleostei; Order: Cypriniformes) is the genus of the family Cyprinidae (Rüber et al., 2007; Britz and Conway, 2009) and *P. progenetica*, *P. micromegethes* and *P. carbunculus* are the three described species to date (Kottelat et al., 2006; Rüber et al., 2007). *Paedocypris* is derived from the Greek words "*Paideios*" (children) and "*Cypris*" (Venus), a common suffix for cyprinid genera. *Paedocypris* is the world's smallest vertebrate and was first discovered in Southeast Asia's peat swamp forests (Kottelat et al., 2006). These miniature fish species are restricted and stereotopic to acidic peat swamp forests, for example, *P. progenetica* at Sumatra and Bintan island, Indonesia, *P. micromegethes* at Matang, Sarawak, and *P. carbunculus* at Kalimantan, Indonesia (Kottelat et al., 2006; Rüber et al., 2007). Besides that, there are some undescribed *Paedocypris* species, such as *Paedocypris* sp "North Selangor", "Pahang", and "Pondok Tanjung" were distributed in the peat swamp forests of Peninsular Malaysia (Sule et al., 2016; Ng et al., 2019), and also from Indonesia, namely *Paedocypris* sp "Pulau Singkep" and "Banka" (Rüber et al., 2007).

Paedocypris is the smallest fish and vertebrate known in the world, displaying a remarkably larval-like appearance with a sexually matured body size for adult fish that can reach only 7.9-11.5 mm standard length. (Kottelat et al., 2006; Britz and

Kottelat, 2008). The anatomy of *Paedocypris* is simplified, combined with developmentally truncated features, including loss of scales, reduction in the number of fin rays, absence of some bones, poorly ossified skeleton, presence of post-anal-fin fold along the ventral edge of caudal peduncle that resembles the fish larvae, and a marked reduction in size (Kottelat et al., 2006; Britz and Kottelat, 2008; Britz and Conway, 2009). In addition, *Paedocypris* is sexually dimorphic. This unique morphological novelty was observed in the modified pelvic fins in males associated with the greatly enlarged and supported keratinized pads of skin, termed "flange and hook" (Kottelat et al., 2006).

2.3.1 *Paedocypris* serve as an important model fish to study

Model organisms play a vital role to provide insight into the general principles in biological research related to genetics, development and evolution. They usually possess unique or interesting characteristics helped to decipher specific processes. In the past decades, model organisms often possess favourable traits such as short generation times, easy handling and culturing, and facilitating experimental laboratory research. However, due to the advancement of sequencing technologies, many whole-genome assemblies, "genomic" models (from organisms with unique/favourable traits) are available, and most of the research studies are based on the analyzed data obtained from the assembled genome (Hedges, 2002). *Paedocypris* species is miniature fish, truncated in development with unusual reproductive mode and is adapted well in an acidic environment, making this fish a potential model organism for genetic study of vertebrate development, the evolution of miniaturization, acidic adaptation, and reproductive behaviour (Liu et al., 2012).

2.3.2 Morphological classification and phylogenetic position of *Paedocypris*

There was considerable interest to resolve the phylogenetic position of *Paedocypris* using morphological characteristics (Britz and Conway, 2009; Mayden and Chen, 2010; Britz et al., 2014), a combination of nuclear and/- or mitochondrial markers (Fang et al., 2009; Mayden and Chen, 2010; Tang et al., 2010; Tang et al., 2013; Stout et al., 2016; Hirt et al., 2017) or single marker gene alone (Rüber et al., 2007). As a result, it was suggested that *Paedocypris* was a member of the Danioninae subfamily (Rüber et al., 2007; Fang et al., 2009; Tang et al., 2010; Tang et al., 2013), within Cyprinidae family but not within Danioninae subfamily (Yang et al., 2015), or as a lineage sister to all Cypriniformes (Mayden and Chen, 2010; Stout et al., 2016; Hirt et al., 2017). The differences in *Paedocypris* phylogenetic position observed from several molecular studies could be due to the high evolutionary rate and base compositional heterogeneity in miniature fishes. It is known that base compositional heterogeneity and long branches (associated with high molecular evolution rates) can influence phylogenetic reconstruction (Hirt et al., 2017).

2.4 Genome assembly

Each species adapts and lives in a specific environmental condition, an ecological niche, that shapes its unique genome sequence and expression (Stange et al., 2021). Thus, it is crucial to sequence and assemble the organisms' genome to provide insights into their genome diversity. Genome assembly constructs a representative genome sequence of that particular organism through a computational process of taking a large number of short or long DNA fragments and aligning them together to form the original genome sequence (Baker, 2012). Genome assembly is a hierarchical process with continuous improvement of sequence length; for example, the shortest assembly at the beginning are contigs derived from the assembly of sequence reads. Hereafter, the contigs are assembled into longer scaffolds, and then the scaffold lengths are further improved to chromosome-scale (Choudhuri, 2014). There are *de novo* and reference-based of genome assemblies, *de novo* refers to sequencing a novel genome without using a reference sequence, whereas reference-based assembly uses a reference-guided assembly strategy (with the availability of a related genome) during mapping/alignment (Lischer and Shimizu, 2017).

2.4.1 Fish genome study

Many fish genomes were sequenced and assembled to investigate or identify genetic factors underpinning evolution that contributed to the diverse traits among fish species. Tine et al., 2014 generated a high-quality chromosome-scale genome of European sea bass *Dicentrarchus labrax*, a temperate-zone euryhaline teleost. Gene families related to ion and water regulation were identified expanded in the genome, providing insights into euryhaline adaptation. Mudskippers are amphibious teleost fishes that are uniquely evolved to live on mudflats. You et al., 2014 sequenced and assembled the genome of four mudskippers (*Boleophthalmus pectinirostris*, *Scartelaos histophorus*, *Periophthalmodon schlosseri* and *Periophthalmus magnuspinnatus*) to examine genetic changes underlying adaptation towards terrestrial life. Comparative genomic analyses, including gene family and positive selection, were carried out. Expansion of genes involved in the innate immune system was identified from gene family analysis, suggesting their roles in defence against terrestrial pathogens. For positive selection, genes involved in the ammonia excretion pathway were positively selected that could help mudskippers tolerate environmental ammonia. Genomic analysis also revealed genetic changes in mudskipper genomes associated with aerial vision, such as the absence or mutation of certain vision-related genes. The genome from another economically important marine fish species, the large yellow croaker *Larimichthys crocea* was sequenced and assembled (Ao et al., 2015). This fish species displays peculiar behavioral and physiological features, particularly sensitive to environmental stresses, a good model to study. Analysis of the genome revealed expansions of several gene families, specifically associated with vision-related crystallins, olfactory receptors, and the auditory sense related genes, providing insights into the genetic mechanisms of environmental stress adaptation. Whole-

genome sequencing of three endemic cavefish *Sinocyclocheilus* species: the surface-dwelling *S. grahami*; the semi-cave-dwelling *S. rhinoceros*; and the cave-restricted *S. anshuiensis* was performed by Yang et al., 2016 to investigate and understand how genomic change in adaptation to the isolated cave life. Comparative analyses of the *Sinocyclocheilus* genomes showed many genetic changes, including loss of gene, pseudogenes, mutations, and down-regulation that contributed to the regressive characters, such as degeneration of eye, albinism, rudimentary scales, and low fecundity. They also found the expansion of genes associated with the sense of taste in the cave-restricted species, providing a better understanding of cavefish biology. Besides that, the genome of the largest bony fish, ocean sunfish *Mola mola* was sequenced and assembled by Pan et al., 2016. This world's largest fish can grow up to 2.7 m and weigh 2.3 tons. It possessed notable characteristics, such as fast growth rate and largely cartilaginous skeleton. Genomic and comparative analyses showed that growth hormone and insulin-like growth factor 1 (GH/IGF1) were positively selected, pointing to its fast growth rate and large body size. Meanwhile, several genes associated with the extracellular matrix involved in the regulation of bone and cartilage development were under positive selection, which could be the factors contributing to the cartilaginous bone. In Braasch et al., 2016, spotted gar *Lepisosteus oculatus*, whose lineage diverged from teleost before teleost genome duplication, was chosen to be sequenced to facilitate in the connection between human and teleost biomedical models. Analyses of the slowly evolving gar genome showed that many entire gar chromosomes conserved with some tetrapods, including human and some undetectable human-teleost conserved noncoding elements become apparent when using the gar genome. Thus, the gar genome could assist in identifying potential regulatory gene elements shared by teleost and humans. The genomes of the world's smallest

vertebrate, *Paedocypris micromegethes* and *P. carbunculus* were sequenced and assembled by Malmström et al., 2018. The unique features of this miniature species are developmental truncation and remaining larval phenotype with an adult body size of < 8mm. Analyses of the two miniature *Paedocypris* genomes showed evolutionary simplification through extensive loss of *Hox* genes and genome reduction associated with their developmental truncation features. In addition, loss of genes related to the development of muscle, nervous system and skeletal were identified from the genomes, contributing to their progenetic phenotype. Genome size reduction of the *Paedocypris* was associated with short intron and low repeat content. Amazonian freshwater fish *Arapaima gigas* is the largest freshwater fish in the world. Du et al., 2019 have sequenced and assembled the *Arapaima* genome to understand this giant fish's biology. Gene family and positive selection analyses identified genes related to its unique characters, especially its large size and fast growth. The application of genome sequence from both sexes (male and female) and RAD-tag analyses enable the identification of male- (Y-) specific scaffolds, which is useful in developing the XY sex-determination system in *Arapaima*. The Antarctic environment is the coldest place on Earth where blackfin icefish *Chaenocephalus aceratus* live. Icefishes are “white blooded” because they are the only vertebrate without functional haemoglobin genes and red blood cells. Whole-genome sequencing on blackfin icefish was performed by Kim et al., 2019. Gene family analysis on the icefish genome identified genes involved in protection from ice damage and control cellular redox state are highly expanded, as an evolutionary adaptation to this cold and high oxygen concentration Antarctic environment, respectively. Loss of genes involved in regulating circadian homeostasis was detected in the blackfin icefish genome. The availability of the blackfin icefish genome enables the understanding of how they adapt

to an extreme environment. Anglerfish *Lophius litulon* is a widely consumed fish with unique features, including a scaleless body, diverse feeding habits and large liver. A high-quality genome assembly of anglerfish *L. litulon* was generated (Lv et al. 2020a). Comparative genomic analyses, including gene family and positive selection, were performed. Gene family analysis revealed expanded and contracted gene families are correlated with adaptation to the benthic environment and its scaleless phenotype; for positive selection, genes involved in metabolic processes are positively selected, which could help diversify the prey for anglerfish. Hadal zone is an extreme environment with high hydrostatic pressure, low temperatures, lack of food supply, and limitation of light. Snailfish is the fish inhabiting the hadal environment with a depth below 6,000 m, and its genome has been sequenced and assembled by Mu et al., 2021. Comparative genomic analyses showed that genes involved with DNA repair were positively selected and expanded in the snailfish genome, showing the importance of maintaining DNA integrity under high hydrostatic pressure. Gene families associated with taste receptors, olfactory receptors, and vision-related genes were significantly changed in the snailfish genome, suggesting these changes might help snailfish adapt to nutrient-poor and dark hadal environments. All findings from fish genome studies mentioned above were summarized in Table 2.1.

Table 2.1 Summary of the fish genome study.

Fish species	Findings from genome
European sea bass (<i>Dicentrarchus labrax</i>) (Tine et al., 2014)	Expansion of gene families contribute to euryhaline adaptation
Mudskipper <i>Boleophthalmus pectinirostris</i> , <i>Scartelaos histophorus</i> , <i>Periophthalmodon schlosseri</i> and <i>Periophthalmus magnuspinnatuto</i> (You et al., 2014)	Genome analyses revealed expansion of genes involved in the innate immune system, positive selection on genes involved in the ammonia excretion pathway, and absence or mutation of certain vision-related genes associated with terrestrial life adaptation in mudskipper.
Large yellow croaker <i>Larimichthys crocea</i> (Ao et al., 2015)	Expansion of gene families associated to the adaptation and response to environmental stress
<i>Sinocyclocheilus</i> cavefish (Yang et al., 2016)	Genetic changes, including loss of gene, pseudogenes, mutations, down-regulation, and gene expansion is associated with the evolved cavefish features for adaptation
Ocean sunfish <i>Mola mola</i> (Pan et al., 2016)	Identify positively selected genes corresponding to the fast growth rate and largely cartilaginous skeleton
Spotted gar <i>Lepisosteus oculatus</i> (Braasch et al., 2016)	Many entire gar chromosomes are conserved from bony vertebrate ancestors, facilitating the comparison of human and teleost medical models.
Miniature fish <i>Paedocypris</i> (Malmstrøm et al., 2018)	Developmental truncation of <i>Paedocypris</i> or evolutionary simplification through the loss of <i>Hox</i> genes or genes related to muscle, nervous system, and skeleton development, whereas short intron and low repeat content contribute to their small genome size. Findings from the genomic analyses have provided insights into miniaturization.
Amazonian freshwater fish <i>Arapaima gigas</i> (Du et al., 2019)	Gene family and positive selection analyses identified genes related to its unique characters, especially its large size (gigantism) and fast growth. Isolation of the male- (Y-) specific scaffolds and supports the XY sex determination system.
Antarctic blackfin icefish <i>Chaenocephalus aceratus</i> (Kim et al., 2019)	Discovered the expanded genes involved in protection from ice damage and control cellular redox state and absence of genes involved in regulating circadian homeostasis in the blackfin icefish genome, providing insights on Antarctic adaptation.
Anglerfishes <i>Lophius litulon</i> (Lv et al. 2020a)	Comparative genomic analyses, including gene family and positive selection, discovered the genes that contribute to the scaleless phenotype, adaptation to the benthic environment and diverse feeding habits.
Yap hadal snailfish (YHS) (Mu et al., 2021)	Comparative genomic analyses showed that genes involved with DNA repair were positively selected and expanded in the snailfish genome; the gene families associated with taste receptors, olfactory receptors, and vision-related genes were significantly changed in the snailfish genome providing insight into adaptation in hadal environments.

2.5 DNA sequencing

DNA sequencing is the process to determine nucleotide content (Adenosine, Guanine, Cytosine, and Thymine) in a DNA fragment. The earliest sequencing (first generation) strategy, known as Sanger sequencing, was developed in 1977 by F. Sanger and colleagues based on a chain-termination approach (Sanger et al., 1977). The Sanger sequencing was commonly applied in single sequencing reactions using a specific DNA primer, especially in verifying PCR products or plasmid constructs. Since then, the sequencing technologies have been constantly improved, and “next-generation” sequencing has begun to emerge.

2.5.1 Next generation sequencing

Next-generation sequencing methods enable the examination of billions of DNA and RNA templates via a high-throughput sequencing technique. Short-read-based second-generation sequencing provides low cost and higher accurate data (Goodwin et al., 2016; Levy and Myers, 2016; Slatko et al., 2018). Short-reads paired-end is being used for genome survey, genomic base correction and assembly in fish, such as *Oxygymnocypris stewartii* (Liu et al., 2019), golden pompano *Trachinotus ovatus* (Zhang et al., 2019) and *Labeo rohita* (Das et al., 2020). Meanwhile, this low-cost and high-throughput Illumina sequencing approach has been successfully used to sequence and assemble the small organelle, the chloroplast and mitochondria (Nunez and Oleksiak, 2016; Cheng et al., 2020). The sequencing approach has been applied to sequence and assemble many fish mitochondria DNA (mtDNA), and the studies are provided in Table 2.2.

Table 2.2 Sequencing approach utilized in the fish mitochondrial genome assembly.

Fish species	Sequencing approach	Genome assembly	Application
Tonguefish, <i>Cynoglossus trigrammus</i> (Mu et al., 2015)	Illumina 2×500 bp paired-end reads	Complete	Unusual mitochondrial genome structure and phylogenetic analysis
Endangered Japanese Swellshark <i>Cephaloscyllium umbratile</i> (Zhu et al., 2017)	Illumina 2×101 bp paired-end reads	Complete	mtDNA features and phylogenetic analysis
Banded cusk-eel <i>Raneya brasiliensis</i> (Fromm et al., 2019)	Illumina 2×150 bp paired-end reads	Complete	Gene rearrangement and phylogenetic reconstructions
Medicinal fish <i>Cyprinion semiplotum</i> (Sharma et al., 2020)	Illumina 2×150 bp paired-end reads	Complete	mtDNA features and phylogenetic implications

Besides that, in short-insert paired-end sequencing, Illumina also provides long-insert mate-pair sequencing, which can sequence more than 1kb DNA fragments. A combination of short-insert paired-end and long-insert mate-pair sequences can improve genome assembly quality. Long insert reads (mate pairs) facilitate mapping across a greater distance that can better cover highly repetitive regions. In contrast, short insert reads (paired-end) can fill in gaps missed by the long inserts (van Heesch et al., 2013). This combination approach (Figure 2.1) is widely adopted in many whole-genome sequencing projects, including vertebrate and invertebrate organisms. For vertebrates, especially in fish, both Illumina paired-end and mate-pair were applied in the generation of the draft genome in whale shark (Weber et al., 2020), forming a hybrid assembly approach with Nanopore in *Danionella translucida* to achieve high assembly contiguity (Kadobianskyi et al., 2019), and to perform initial genome assembly prior to chromosome construction in Anglerfish genome (Lv et al., 2020a). In invertebrates, both Illumina paired-end and mate-pair were applied to generate draft genome with PacBio in jellyfish (Gold et al., 2019) and cricket (Ylla et al., 2021) for

evolutionary study. Thousands of eukaryotic genomes were assembled based on Illumina reads; however, the mate-pair libraries are still unable to span all the repetitive regions, resulting in highly fragmented assemblies. Therefore, long-read sequencing approaches such as Pacific Bioscience (PacBio) and Oxford Nanopore Technology (ONT) and chromosome conformation capture methods, including Hi-C and Dovetail Genomics Chicago libraries were, emerged to resolve and span the repetitive regions in order to improve the continuity of genome assemblies (Elbers et al., 2019).

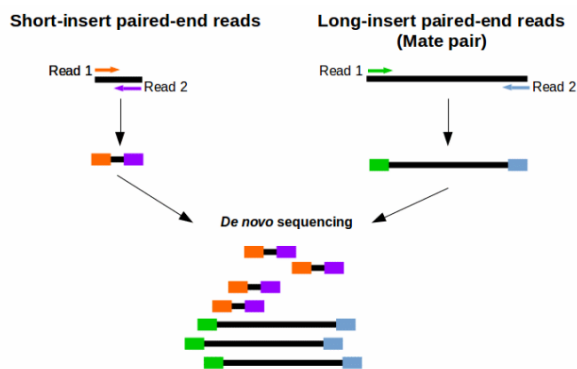


Figure 2.1 Combination of long inserts from mate pair sequencing with that from short-insert paired-end sequencing for *de novo* assembly. This figure was adopted from website: <https://www.ecseq.com/support/ngs/what-is-mate-pair-sequencing-useful-for>

With the rapid evolution of sequencing technologies, third-generation sequencing (also called long-read sequencing) technology focuses on generating long reads DNA (or can be RNA) longer than second-generation sequencing was developed. The technology adopted is called Single-Molecule Sequencing in Real-Time (SMRT), which enables the sequencing of very long fragments up to 30 to 50 kb or longer (average 10-25 kb). The long-continuous reads can help produce large scaffolds to overcome limitations encountered when using second-generation sequencing, specifically to assemble the highly repetitive genomic regions that mainly caused numerous short breaks in short-reads assemblies. Long-read PacBio or ONT can generate long sequence reads for prokaryotic genome assemblies. However, for eukaryotic genomes, this technique can improve their genome assemblies to scaffold-level, but insufficient decipher to the chromosome-scale level (Belser et al., 2018).

2.5.2 Chromosome conformation capture techniques

New technologies utilize long-range linking information, including linked reads (10x Genomics) (Zheng et al., 2016;Yeo et al., 2017), optical mapping (BioNano) (Schwartz et al., 1993;Shelton et al., 2015;Vij et al., 2016) and chromosome conformation capture such as Hi-C and Dovetail Genomics Chicago libraries (Dovetail) (Simonis et al., 2006;Burton et al., 2013;Elbers et al., 2019) were being adopted to validate orientation and further improve draft genome assemblies by spanning the previously unscaffolded contigs. Hi-C data can provide long-range linkage information across a variety of length scales, spanning tens of megabases. Meanwhile, Hi-C data can be mapped to draft genome assemblies to scaffold the contigs, generating super/ultra-long-range-scaffolds in the size range of chromosomes (Elbers et al., 2019;Ghurye et al., 2019). Therefore, Hi-C data is being adopted in many

genome projects to improve the draft genome assemblies up to a chromosome scale to obtain a high-quality reference genome (Burton et al., 2013;Marie-Nelly et al., 2014;Dudchenko et al., 2017;Zhang et al., 2019). Illumina, PacBio, and Hi-C high-throughput sequencing technologies have been used in *de novo* fish genome studies, and chromosome-level genome assemblies were reported, as summarized in Table 2.3.

Table 2.3 Sequencing approach applied for fish chromosome-level genome assembly.

Fish species	Sequencing approach	Genome assembly
<i>Takifugu favidus</i> (Zhou et al., 2019)	Illumina short reads*, PacBio Sequel**, Hi-C***	Chromosome-level
<i>Ancherythroculter nigrocauda</i> (Sun et al., 2020b)	Illumina short reads*, PacBio Sequel**, Hi-C***	Chromosome-level
<i>Argyrosomus japonicus</i> (Zhao et al., 2021)	Illumina short reads*, PacBio Sequel**, Hi-C***	Chromosome-level

*Genome survey and genomic base correction; **Genome assembly; ***Chromosome construction

2.6 Fatty acids

Fatty acids are building blocks of lipids, made up of an aliphatic chain with a carboxyl group (-COOH) at one end and methylene group at the other end that can be classified into four categories: saturated, monounsaturated, polyunsaturated, and *trans* fats (Rustan and Drevon, 2005; White, 2009). Saturated fatty acids (SFAs) have no double bonds of linear hydrocarbon chains (CH₂), with one carboxyl group attached at the terminal end. A common source of SFAs can be derived from animals and plants (i.e. palm oil and coconut oil). High consumption of SFAs is not suitable for human health that can cause health problems such as raised cholesterol levels, cardiovascular disease, obesity, inflammation, and increased diabetes risk (Iggman and Ris érus, 2011). However, a study using a tropical fish, barramundi, which is living in a habitat with high temperature, showed no effect on the growth of fish with a dietary supplement of SFAs, suggesting the SFAs might be able to utilize and digest for energy in some fish species (Salini et al., 2017).

Unsaturated fatty acids have one or more double bonds (C=C) in their carbon chains. Monounsaturated fatty acids (MUFAs) consist of only one double bond, existing in either *cis* (Z-configuration) or *trans* (E-configuration) configuration. In *cis*-configuration, the two adjacent hydrogens lie on the same side of the double-bonded carbons, whereas in *trans* configuration, the two adjacent hydrogens lie on opposite sides double-bonded carbons (Rustan and Drevon, 2005; Christie and Han, 2012). In nature, *cis*-MUFAs are the predominant form compared with *trans* and oleic acid (C18:1n-9), palmitoleic acid (C16:1n-7), and vaccenic acid (C18:1n-7) are common *cis*-MUFAs found in the dietary supplements. For daily nutrition, more MUFAs can obtain from vegetable oils, including olive oil, canola oil, hazelnut oil, mustard oil, almond oil and avocado oil (Schwingshackl and Hoffmann, 2014). Dietary