

1 **SI Appendix for:**

2 **The subgenomes show asymmetric expression of alleles in hybrid lineages of *Megalobrama***
3 ***amblycephala* × *Culter alburnus***

5 **Short title: Asymmetric expression in hybrid fish lineages**

7 Li Ren,^{1,b} Wuhui Li,^{1,b} Qinbo Qin,^{1,b} He Dai,^{2,b} Fengming Han,² Jun Xiao,¹ Xin Gao,¹ Jialin
8 Cui,¹ Chang Wu,¹ Xiaojing Yan,¹ Guoliang Wang,³ Guiming Liu,³ Jia Liu,¹ Jiaming Li,¹ Zhong
9 Wan,⁴ Conghui Yang,¹ Chun Zhang,¹ Min Tao,¹ Jing Wang,¹ Kaikun Luo,¹ Shi Wang,¹ Fangzhou
10 Hu,¹ Rurong Zhao,¹ Xuming Li,² Min Liu,² Hongkun Zheng,² Rong Zhou,¹ Yuqin Shu,¹ Yude
11 Wang,¹ Qinfeng Liu,¹ Chenchen Tang,¹ Wei Duan,¹ and Shaojun Liu^{1,a}

13 ¹State Key Laboratory of Developmental Biology of Freshwater Fish, College of Life Sciences,
14 Hunan Normal University, Changsha, China

15 ²Biomarker Technologies Corporation, Beijing 101300, China

16 ³Beijing Agro-Biotechnology Research Center, Beijing Academy of Agriculture and Forestry
17 Sciences, 100097, Beijing, China

18 ⁴School of Mathematics and Statistics, Central South University, Changsha, 410083, Hunan,
19 P.R. China

21 **^aCorresponding author**

22 Professor Shaojun Liu

23 Email: lsj@hunnu.edu.cn

24 State Key Laboratory of Developmental Biology of Freshwater Fish, Hunan Normal University,
25 Changsha 410081, China

26 Tel./Fax: +86-073188873074

28 **^bThese authors contributed equally to this work.**

30 **Keywords:** Hybridization, Whole-genome sequencing, Allelic exchange, Expressed
31 recombinant transcripts, *Cis*- and *trans*-regulatory expression

1 **1 SI methods and results of genome analysis**

2 **1.1 Sample preparation for genome and transcriptome sequencing**

3 A wild, water-captured male adult of topmouth culter (*Culter alburnus*, TC) and a
4 gynogenetic individual of blunt snout bream (*Megalobrama amblycephala*, BSB) were
5 obtained for whole-genome sequencing. To obtain the gynogenetic BSB individual, BSB eggs
6 were stimulated by UV-inactivated sperm of *Cyprinus carpio haematopterus* (Gong et al. 2019).
7 One minute after stimulation, the eggs were placed in a water bath at 4-8°C for cold shock
8 treatment for 30-40 min. The purpose of this treatment was to suppress the discharge of the
9 second polar body, resulting in diploid gynogenetic individuals. Then, we collected 26 fish from
10 the Engineering Center of Polyploid Fish Breeding of the National Education Ministry at Hunan
11 Normal University, China, for RNA-seq. The samples included the following fishes: three 2-
12 year-old TC males whose testes were mature; three 2-year-old BSB females whose ovaries were
13 mature; nine 2-year-old 2nF₁-2nF₃ female hybrids of BSB (♀) × TC (♂) (abbreviated as BT);
14 and nine 2-year-old 2nF₁-2nF₃ female hybrids of TC (♀) × BSB (♂) (abbreviated as TB) (Fig.
15 1 A-H). Among these hybrid offspring, only “real hybrid” with two inbred parental
16 chromosomes were identified with 45S rDNA and used in our study (Xiao et al. 2016). All
17 tissue samples were excised carefully and subsequently stored at -80°C. The Administration of
18 Affairs Concerning Animal Experimentation guidelines state that approval from the Science
19 and Technology Bureau of China and the Department of Wildlife Administration is not
20 necessary when the fish in question is not rare or near extinction, as is the case in this paper.

21
22 **1.2 Sequencing and assembly**

23 High-quality DNA was extracted from the muscle of one adult male of TC and one adult
24 female of BSB using the DNeasy Blood and Tissue Kit (Qiagen). We constructed 13 paired-
25 end libraries (220 bp to 17 kb) of TC and 10 paired-end libraries (220 bp-17 kb) of BSB to
26 produce the raw data according to the Illumina standard operating procedure. A total of ~171
27 Gb (~194 X) and ~181 Gb (~161 X) of reads was obtained for every sequenced locus for TC
28 and BSB, respectively (Supplemental Tables S1, S2). Before the beginning of assembly,
29 sequencing adaptors of raw reads were removed. Then, contaminated reads (mitochondrial,
30 bacterial and viral sequences, etc.) were screened by alignment to the NCBI-NR database using
31 BWA (v 0.7.13) (Li and Durbin 2010) with default parameters. FastUniq (v1.1) (Xu et al. 2012)
32 was used to remove the duplicated read pairs, while the low-quality Illumina reads were
33 removed based on the following criteria: reads with ≥ 10% unidentified nucleotides (N); reads
34 with ≥ 10 nt aligned to the adaptor, allowing ≤ 10% mismatches; and reads with ≥ 50% of bases
35 having a Phred quality < 5.

36 Illumina reads (44.27 Gb and 51.96 Gb of data for TC and BSB, respectively) from the
37 220-bp library were selected to estimate genome size. This analysis was performed using

1 “kmer_freq_stat” software (developed by Biomarker Technologies) and showed a major peak
2 at 48X and 38X with a 21-kmer frequency distribution for TC and BSB, respectively
3 (Supplemental Fig. S1). Based on the formula Genome size = kmer_Number/Peak_Depth, the
4 genome sizes of TC and BSB were estimated to be ~946.91 Mb and ~1.12 Gb, respectively. At
5 the same time, we determined the genome size of TC and BSB with flow cytometry, using
6 zebrafish (1.37 Gb, GRCz10 downloaded from the Ensembl database) as a reference standard
7 and propidium iodide as the stain. Our flow cytometry analysis showed that the genome sizes
8 of TC and BSB were ~ 1.12 Gb and 1.15 Gb, respectively.

9 Genome assembly was performed with ALLpaths-LG (Gnerre et al. 2011), joining the
10 contigs into scaffolds with SSPACE (v 2.3) (Boetzer et al. 2011) and filling the gaps within
11 scaffolds with GapCloser (v 1.12) in the SOAPdenovo package (Luo et al. 2012), resulting in
12 two preassemblies: 1) a total length of 993.5 Mb and a contig and scaffold N50 size of 19.6 kb
13 and 3.4 Mb, respectively, for TC and 2) a total length of 1.04 Gb and a contig and scaffold N50
14 size of 33.7 kb and 1.33 Mb, respectively, for BSB. We also sequenced the genomes of TC and
15 BSB with ~6X (mean subread length: 8,523 bp, subread N50: 11,677 bp) and ~10X (mean
16 subread length: 14,092 bp, subread N50: 19,206 bp) coverage, respectively, using long PacBio
17 reads and used these reads with the PBJelly algorithm to close gaps in the ALLpaths-LG
18 assembly above. The addition of long PacBio reads strikingly increased the length of the contig
19 N50 from 19.6 kb to 72.2 kb for TC and from 33.7 kb to 142.7 kb for BSB (Table 1).

20

21 **1.3 Assessment of assembly quality**

22 To evaluate the integrity of the assembly, we assembled 45,055 and 13,871 unigenes (\geq
23 500 bp) with RNA-seq reads from multiple tissues using Trinity and aligned these unigenes to
24 the assemblies of TC and BSB using BLAT (Kent 2002), yielding an average of more than 99%
25 (44,787 for TC and 13,837 for BSB) coverage of these transcripts based on sequence similarity.
26 The two assemblies were also analyzed using CEGMA and BUSCO (v. 2.0.1) with default
27 parameters to evaluate their completeness. Our results suggested that more than 99% (456 for
28 TC and 458 for BSB) conserved core eukaryotic genes (CGEs) and more than 98% (244 for TC
29 and 246 for BSB) highly conserved CGEs were found in the TC and BSB assemblies. In
30 addition, more than 90% of metazoan and vertebrate BUSCOs were covered by the assembled
31 genomes (Supplemental Table S3). These evaluations suggested a high-quality assembly for
32 TC and BSB.

33

34 **1.4 Gene prediction and annotation**

35 Protein-coding genes were predicted using *de novo*, protein homology-based and
36 transcript-based approaches. To elaborate, GENSCAN (v 1.0) (Burge and Karlin 1997),
37 Augustus (v 2.5.5) (Stanke and Waack 2003), GlimmerHMM (v 3.0.1) (Majoros et al. 2004),

1 GeneID v1.3 (Parra et al. 2000) and SNAP (Johnson et al. 2008) with the default parameters
2 were used for *de novo* gene prediction, and all of these software packages were tested using the
3 zebrafish (GCF_000002035.6) and grass carp (<http://www.ncbi.nlm.nih.gov/grasscarp/>) gene models
4 before gene prediction. Homologous peptides from the zebrafish, grass carp, and common carp
5 (GCF_000951615.1) genomes were aligned to our assembly to identify homologous genes with
6 GeMoMa (v 1.4.2) (Keilwagen et al. 2016). The RNA-seq reads of multiple organs of TC and
7 BSB, including muscle, blood, liver, testis, and ovary, were assembled into transcripts using
8 Trinity (v 2.4.4) (Grabherr et al. 2011) and aligned to the corresponding repeat-masked TC and
9 BSB assemblies with BLAT, and the gene structures were modeled using the Program for
10 Automated Sequential Assignment. Consensus gene models were generated by integrating the
11 *de novo* predictions and protein and transcript alignments using EVidenceModeler (Haas et al.
12 2008). Functional annotation was performed based on comparisons with the SwissProt,
13 TrEMBL, InterPro, COG and KEGG protein databases. The gene ontology (GO) for each gene
14 was assigned by Blast2GO (Conesa et al. 2005) based on NCBI databases (Supplemental Table
15 S7). Overall, 30,835 (99.40%) and 28,930 (93.72%) of the protein-coding genes could be
16 annotated, for which more than 70% of the predicted exonic regions were covered by aligning
17 RNA-seq reads and assembled transcripts to the genome assembly using TopHat v2.1.2 and
18 BLAT, respectively. Finally, we obtained 30,443 and 29,994 protein-coding genes
19 (Supplemental Table S8). Then, comparisons of coding sequence (CDS) number, CDS length,
20 intron length and gene length were performed in four Cyprinid fish (*Danio rerio*,
21 *Ctenopharyngodon idellus*, BSB and TC) (Supplemental Fig. S3).

22

23 **1.5 Noncoding RNA prediction**

24 Noncoding RNAs play an important role in life processes, such as the rRNAs and tRNAs
25 involved in mRNA translation. Therefore, the noncoding RNAs in our TC and BSB assemblies
26 were predicted. The rRNA fragments were identified by aligning the rRNA template sequences
27 (Pfam database v 22.0) using BLAST with an *e*-value of 1e⁻¹⁰ and identity cutoff of 95% or
28 more. The tRNAscan-SE v 2.0 algorithms (Lowe and Chan 2016) with default parameters were
29 applied to predict tRNA genes. The miRNA, snRNA and snoRNA genes were identified by
30 mapping the genome sequences to the Rfam v11.0 database using INFERNAL (v 1.1)
31 (Nawrocki and Eddy 2013). Finally, we identified 335 miRNAs, 628 rRNAs, 891 tRNAs, 106
32 snRNAs and 203 snoRNAs in TC, while 302 miRNAs, 474 rRNAs, 470 tRNAs, 92 snRNAs
33 and 198 snoRNAs were detected in BSB (Supplemental Table S9).

34

35 **1.6 Evolutionary analysis**

36 Protein sequences of *D. rerio* (zebrafish), *C. idellus* (grass carp), *C. carpio* (common carp),
37 *Xiphophorus maculatus* (platy fish), *Oryzias latipes* (medaka), *Fugu rubripes* (fugu),

1 *Larimichthys crocea* (large yellow croaker) and *Cynoglossus semilaevis* (flatfish) were
2 downloaded from the Ensembl or GenBank database. The proteomes of the above eight species
3 and our proteomes of TC and BSB, comprising 26,003 and 26,566 protein sequences,
4 respectively, were clustered into 27,498 orthologous groups using OrthoMCL (Li et al. 2003)
5 based on an all-to-all BLASTP strategy with an *e*-value of 1e⁻⁵ and a Markov chain clustering
6 (MCL) default inflation parameter of 1.5. Based on the clustering results, TC- and BSB-specific
7 gene clusters were obtained and annotated. To infer phylogenetic relationships, we extracted
8 796 single-copy clusters from all ten species, and multiple sequence alignment of proteins for
9 each cluster was performed by MUSCLE (Edgar 2004) (Fig. 2A; Supplemental Table S10); all
10 the alignments were combined into one supergene to construct a phylogenetic tree using
11 RAxML (v 7.0.0) (Stamatakis 2014) with 1,000 rapid bootstrap analyses followed by a search
12 of the best-scoring maximum-likelihood (ML) tree in a single run (Fig. 2A). Finally, divergence
13 time was estimated using MCMCTree from the PAML package (Yang 2007), together with the
14 molecular clock model. Two reference calibrated time points provided by the TimeTree
15 database (<http://timetree.org/>) were used to assess the divergence times of nodes of interest.
16 Expansion and contraction of OrthoMCL-derived homologous clusters were determined by
17 CAFÉ (v 2.1) calculation based on changes in gene cluster size, utilizing the phylogeny and the
18 species divergence time.

19

20 **1.7 Diversifying selection analysis**

21 Orthologous gene pairs between TC and BSB were determined by best reciprocal BLAST
22 hits with an *e*-value of 1e⁻⁵. The *K_s* value between the orthologous pairs was calculated by the
23 yn00 program in the PAML package. The all-to-all BLASTP method was used to detect the
24 orthologous genes in the TC-*C. idellus*, TC-BSB, BSB-*C. idellus*, *D. rerio*-TC, *D. rerio*-*C.*
25 *idellus* and *D. rerio*-BSB pairs with an *e*-value threshold of 1e⁻⁵. Homologous blocks were
26 detected using MCScanX (Wang et al. 2012), and the 4DTV (transversions at fourfold
27 degenerate sites) values of the blocks in the CDS alignments were calculated using the HKY
28 model. The distribution of *K_s* values was used to determine the events of species divergence
29 (Fig. 2B).

30

31 **1.8 Genetic map construction and scaffold anchoring**

32 High-quality genomic DNA extracted from 106 individuals of TC (including one male
33 parent, one female parent, and 104 F₁ progenies) was used to construct Illumina sequencing
34 libraries following the manufacturer's protocol. The genomic libraries were sequenced on an
35 Illumina HiSeq 2500 system, and a total of 133.1 M reads were obtained. The average
36 sequencing depths of sequenced loci per parent and per progeny were ~28.0X and ~6.5X,
37 respectively. Genotyping and evaluation of the quality of genetic markers were performed as

1 described in Sun et al. After filtering, a total of 6,515 high-quality SNP markers were identified
2 and used to construct the genetic map with the double pseudo-testcross strategy using HighMap
3 software (Liu et al. 2014). Using the nearest-neighbor method, a total of 6,377 markers were
4 clustered into 24 linkage groups with a total genetic length of 3,867.4 cM. The SNP markers
5 were then aligned to the TC assembled scaffolds, and only uniquely aligned markers were used
6 to anchor and orient the scaffolds onto the 24 TC pseudochromosomes.

7

8 **1.9 Hi-C assembly of the BSB genome**

9 According to the Hi-C procedure, nuclear DNA from the blood of the BSB individual was
10 cross-linked and then cut with a restriction enzyme, leaving pairs of distally located but
11 physically interacting DNA molecules attached to one another. The sticky ends of these digested
12 fragments were biotinylated and then ligated to each other to form chimeric circles. Biotinylated
13 circles, which are chimeras physically associated with DNA molecules from the original
14 crosslinking, were enriched, sheared and sequenced. In our study, a total of 202.3 million clean
15 Hi-C read pairs (60.58 Gb) with approximately 55.64-fold coverage of the BSB genome were
16 obtained. Of these reads, 87.32% were mapped to the BSB genome, and 58.53% showed unique
17 alignment to the genome. After filtering, we obtained 88.53 M valid interaction pairs for the
18 chromosome-level assembly. Subsequently, the scaffolds within the BSB assembly were broken
19 into 50-kb fragments and clustered by LACHESIS software (Burton et al. 2013) using valid
20 interaction read pairs. Finally, 2,867 scaffolds with a total length of 1,048 Mb were anchored
21 to the BSB chromosomes, of which 1,901 scaffolds (983.74 Mb in length) were assigned,
22 ordered and oriented to 24 chromosome-level groups (Supplemental Fig. S2). The syntenic
23 relationships among the TC, BSB and *D. rerio* chromosomes were constructed (Fig. 2C;
24 Supplemental Fig. S4).

25

26 **2 SI methods and results of allelic recombinant analysis**

27 **2.1 RNA isolation and transcriptome sequencing**

28 To sequence the transcriptomes of reciprocal cross hybrids and their inbred parents, total
29 RNA was isolated and purified from the liver (24 samples), muscle (6 samples) and gonad
30 (ovary) (6 samples) by a TRIzol extraction method. The RNA concentration was measured
31 using NanoDrop technology. Total RNA samples were treated with DNase I (Invitrogen) to
32 remove any contaminating genomic DNA. The purified RNA was quantified using a 2100
33 Bioanalyzer system (Agilent, Santa Clara, CA, USA). We fragmented 1 μ g of isolated mRNA
34 with fragmentation buffer. The resulting short fragments were reverse transcribed and amplified
35 to produce cDNA. An Illumina RNA-seq library was prepared according to a standard high-
36 throughput method (Dillies et al. 2013). The cDNA library concentration and quality were
37 assessed by the Agilent Bioanalyzer 2100 system, after which the library was sequenced with

1 a paired-end setting using the Illumina HiSeq 2000/4000/X Ten platform. An RNA-seq
2 experiment was conducted with three biological replicates. Then, the raw reads containing
3 adapters and poly-N tails and of low quality were removed using a custom computational script
4 (see Supplemental Script 1: filtering raw reads). The high-quality reads were used in the next
5 analysis.

6

7 **2.2 Ortholog identification and allelic recombinant sequence detection by Illumina**
8 **sequencing**

9 The orthologous gene alignments between BSB and TC were obtained from all-against-
10 all reciprocal BLASTP (v 2.2.26) comparisons with the parameters "-e 10-5 -F F -v 1 -m 8"
11 based on the protein sequences. The alignments were used in MCScanX to determine syntenic
12 blocks, which were then displayed as a schematic diagram created with Circos (v 0.69-6)
13 (Krzywinski et al. 2009) (Fig. 2C; Supplemental Fig. S4). Transcripts lacking annotated CDSs
14 and those with lengths < 100 bp were discarded. A total of 12,322 orthologous gene pairs were
15 selected based on the best match scores and filtering.

16 Allelic recombinant genes in hybrids originate from the recombination of inbred parental
17 genome sequences. The expression of these proteins can be detected by Illumina sequencing
18 based on a few methods (Francesc et al. 2009; Sabio and Davis 2010; Liu et al. 2016). The
19 high-quality Illumina reads of hybrids in the above analysis were mapped to the mixed genome
20 of two parents using STAR (v 2.4.0) with the parameter “--outFilterMismatchNoverReadLmax
21 0.01/0.02/0.03 --chimSegmentMin 50/70/90/105/120” (Dobin et al. 2013). Focusing on
22 mapping files, pairwise alignments of orthologous gene pairs were used to assess the
23 differences between hybrids. The “--outFilterMismatchNoverReadLmax 0.02” parameter was
24 selected based on low mutation rates and sequencing errors. The other parameter, “—
25 chimSegmentMin”, was related to the standard for recombinant reads, which required at a
26 minimum that fragments be identified in two paired-end reads. Regarding the two styles of
27 paired-end (PE) libraries used in our analysis, the parameter “--chimSegmentMin 70” with 100
28 PE reads*2 samples and “--chimSegmentMin 105” with 150 PE reads*2 samples was selected
29 to detect the recombinant reads in hybrids. Then, only the reads distributed in orthologous gene
30 pairs between BSB and TC were considered recombinant reads. After removing the repeated
31 recombinant reads in one orthologous gene pair, the distribution of recombinant events in
32 reciprocal cross offspring was enriched in the corresponding genes (Branchetti et al. 2013).

33

34 **2.3 PCR validation of allelic recombinant genes**

35 From the short fragments obtained by Illumina sequencing, 19 genes were selected based
36 on the above prediction results. The total DNA of BSB, TC, TBF₁, BTF₁, TBF₃ and BTF₃ was
37 extracted in triplicate. The set of primers used in PCRs and clone numbers are provided in

1 Supplemental Table S14. Amplified products were evaluated using an ABI 3730 DNA Analyzer
2 (Applied Biosystems, Carlsbad, CA, USA). To determine sequence homology and variation
3 among the fragments, sequences were aligned using BioEdit (Huang et al. 2012) and
4 CLUSTALW2 (Huang and Parmacek 2012). The *tgfb1b* sequences of grass carp and zebrafish
5 were downloaded from the NCBI database (accession nos. EU099588.1 and XM_687246.8).

6

7 **2.4 Library construction and PacBio sequencing**

8 To verify the accuracy of the allelic recombinant genes, the unassembled long reads from
9 PacBio sequencing were considered useful for studying allelic recombinant events. Thus, the
10 long transcripts of reciprocal cross F₃ hybrids were obtained with PacBio SMRT sequencing.
11 The total RNA of TBF₃ and BTF₃ from five tissues, including liver, muscle, gonad, heart and
12 hypophysis, was obtained and mixed in equal amounts. The RNA was reverse transcribed using
13 the SMARTer PCR cDNA Synthesis Kit, and PCR amplification was performed using KAPA
14 HiFi PCR Kits. The PCR product (size = 0.5-6 K and > 6 K) was selected based on the agarose
15 gel electrophoresis method. Then, libraries were constructed from these cDNA products using
16 the SMRTbell Template Prep Kit 1.0. After library preparation, the library template and enzyme
17 mixture were used in the PacBio Sequel™ system for sequencing.

18

19 **2.5 Allelic recombinant gene detection by genome-wide long-read alignments**

20 After obtaining the sequencing data, low-quality data (adapter sequences, length of
21 subreads < 50 bp, accuracy rate < 0.75) were deleted from the raw data. Sequence reads from
22 the PacBio RS SMRT chip were processed through PacBio's SMRT-Portal analysis suite to
23 generate circular consensus sequences (CCSs). To obtain more accurate reads, the reads
24 (number of cycles of CCS > 1 and accuracy > 0.8) were used to obtain full-length reads (lengths >
25 300 bp, poly (A) tails, 5' primers and 3' primers) based on the SMRT Iso-Seq analysis pipeline
26 (<http://www.pacb.com/products-and-services/analytical-software/smrt-analysis/>).

27 We detected allelic recombinant transcripts in TBF₃ and BTF₃ by the following steps: a
28 recombinant full-length read was split into two fragments that were mapped to orthologous
29 gene pairs between BSB and TC. To find the paired fragment alignment of a recombinant long
30 read, we aligned reads to the reference genome (mixed genome of two parents) by GMAP with
31 the parameters identity > 90% and coverage \geq 99% (Parmacek and Epstein 2013). In addition,
32 we examined all mutant pairs of fragment alignments and defined them as candidate expressed
33 recombinant transcripts if the following conditions were met: (i) both aligned fragments were
34 of sufficient length (> 200 bp); (ii) two fragments were mapped to an orthologous gene pair;
35 and (iii) exon regions were distributed in both aligned fragments. Long repetitive genomic
36 elements and TGS long reads with sequencing errors were eliminated as described by Jason et
37 al. (Nguyen et al. 2013).

1 The advantages of sequencing technology, including the high throughput of Illumina
2 sequencing (Li et al. 2012; Shen et al. 2011), the long reads in PacBio sequencing (Rhoads and
3 Au 2015), and the high accuracy of Sanger sequencing (Huang et al. 2009; Parmacek and
4 Epstein 2009), and the disadvantages, including short reads in Illumina sequencing (Rong et al.
5 2012), high error rate in PacBio sequencing (Rong et al. 2012), and the low sequencing data
6 content and short fragments in Sanger sequencing (Huang et al. 2009; Parmacek and Epstein
7 2009), compelled us to combine the three methods for predicting recombinant events in hybrids.
8 As a comparison of Illumina and PacBio sequencing, the differences in the number and length
9 of sequencing reads between the two methods may be the main causes leading to a large number
10 of allelic recombinant events in Illumina sequencing and few such events in PacBio sequencing.
11 However, the shared allelic recombinant events detected between the methods proved the
12 accuracy of these methods in predicting allelic recombinant events.
13

14 **3 SI methods and results of expression analysis**

15 **3.1 Mapping of RNA-seq data**

16 All Illumina reads of BSB and TC were aligned to the BSB and TC genomes using STAR
17 (v 2.4.0) with the default parameters, respectively (Dobin et al. 2013). In addition, the other
18 RNA-seq reads of reciprocal cross hybrids were aligned to the two reference genomes of BSB
19 and TC. The number of mapped reads in each gene was calculated with some in-house Perl
20 scripts (see Supplemental Scripts 2: calculating mapped reads). Therefore, two mapping results
21 based on alignment to the two reference genomes were obtained in hybrid offspring. The silent
22 genes were screened in comparisons of both parents with each hybrid, with a threshold of no
23 mapped reads detected in hybrids and mapped read counts ≥ 5 in both parents, with three
24 biological replicates. To avoid the negative effect of expression noise, the expression analysis
25 was performed on only filtered genes with mapped read counts ≥ 5 in all three biological
26 replicates in each comparison. To further avoid the interference of expression noise, the
27 expression values of three biological replicates were screened with a mean ± 2 standard
28 deviation (SD) threshold in each gene (Quackenbush 2002). Moreover, the total expression
29 value was normalized based on the ratio of the number of mapped reads for each gene to the
30 total number of mapped reads for the entire genome. For the transcriptome data of the two
31 parents, 68.10 million clean reads (84.68%) from the liver of BSB and 64.34 million clean reads
32 (84.91%) from the liver of TC were obtained from mapping to the respective reference genome.
33 Additionally, 626.33 million clean reads and 627.92 million clean reads in the hybrid mapped
34 to the two parental genomes, respectively.
35

36 **3.2 Differential expression analysis**

37 To investigate the expression levels of the two parents and hybrids in one comparison, we

1 set the “average parental expression” (APE) value (H_{APE}) related to orthologous gene pairs
2 between BSB (B) and TC (T). The value was equal to the average expression level of
3 orthologous gene pairs in the parents ($H_{APE} = (B + T)/2$). Moreover, the real expression level of
4 the hybrid (H_{real}) was equal to the average of the two expression values aligned to the two
5 parental genomes ($H_{real} = (H_B + H_T)/2$). To compare H_{APE} with H_{real} , the clear changes in
6 expression after hybridization were described in one comparison. The expression values of each
7 gene with an FDR ≤ 0.05 and a \log_2 -transformed ratio of 1.25 were considered differentially
8 expressed genes. The DE analysis was performed with the DESeq2 package in the R program
9 (v 1.10.0) (R Foundation for Statistical Computing, Vienna, Austria) (Varet et al. 2016).

10

11 **3.3 Analysis of expression dominance**

12 To compare the expression levels of hybrids with those of both inbred parents, the mode
13 of the expression level relative to the total expression level of orthologous genes between BSB
14 and TC was used. The comparison for each gene of the hybrid was performed in three steps: 1) The
15 BSB and hybrid expression values were obtained by alignment with the BSB genome, and
16 significant DE was detected. 2) The TC and hybrid expression values were obtained by
17 alignment with the TC genome, and significant DE was detected. 3) Information on the
18 significant DE between the two parents based on orthologous gene pairs between BSB and TC
19 was obtained. The analysis of significant DE was performed with Fisher’s exact tests (FDR)
20 (Robinson et al. 2010), and the distribution of P -values was controlled for with an FDR by the
21 BH method at $\alpha = 0.05$ (Francesc et al. 2009). Genes in the results of the above three
22 comparisons were classified as additive, BSB/TC dominant, underdominant, and overdominant
23 based on the magnitude of the expression difference, as described by Gibson et al. (Gibson et
24 al. 2004).

25

26 **3.4 Species-specific SNP identification**

27 To investigate the expression levels of parent-of-origin genes (allelic expression), the
28 LASTZ pairwise alignment tool (v 1.02.00) (Harris 2007) with default parameters was used to
29 obtain the corresponding loci from the orthologous gene pairs between BSB and TC. The InDel
30 loci were discarded, and only the aligned loci with the best match scores were used for the next
31 analysis. Then, the SNPs were collected using the SNP Calling pipeline with GATK (v 3.8)
32 based on the results of parental transcriptome mapping to the respective genome as described
33 above (McKenna et al. 2010). According to the comparison of SNPs and other loci in
34 orthologous gene pairs between BSB and TC, the differential loci, including heterozygous and
35 homozygous loci, were considered species-specific SNPs, as in Schaeck et al. (McManus et
36 al. 2010; Schaeck et al. 2013). To prevent the negative effect of sequencing and mapping, the
37 screening of species-specific SNPs must be checked for consistency with three biological

1 replicates, and loci possessing read counts ≥ 1 in all accessions and biological replicates in each
2 comparison are retained. After screening 9,753 gene pairs, 103,190 SNPs, including 268
3 heterozygous loci in TC and 41 heterozygous loci in BSB, were obtained from 60,909,895 and
4 52,151,477 clean reads mapped to the BSB and TC reference genomes, respectively
5 (Supplemental Fig. S27).

6

7 **3.5 Detection of allelic expression levels in hybrids**

8 To describe the allelic expression in hybrids, the hybrid transcriptome mapping results
9 (bam files) described above were used in our next analysis. The map files of each hybrid were
10 divided into two categories based on the two different parental reference genomes. The BSB/TC
11 allelic reads in the hybrid were calculated using some in-house Perl scripts based on
12 corresponding BSB/TC species-specific SNPs in the corresponding map files (see
13 Supplemental Scripts 3: calculating allelic reads with SNPs). In addition, the expression levels
14 of parents in the allelic analysis were also detected based on the respective reference genome
15 related to species-specific SNPs. In the above analysis, to remove the negative effect of
16 mutation sites in the hybrid in biological replicates, if the mapped reads of a species-specific
17 SNP did not comply with the mean ± 2 SD threshold in the three biological replicates, the
18 abnormal value was discarded when estimating the BSB and TC allelic expression levels. Then,
19 the total numbers of BSB and TC alleles in each gene were normalized based on the ratio of the
20 number of mapped reads for each gene to the total number of mapped reads for the entire
21 genome (Quackenbush 2002). The sum of BSB/TC allelic reads for all species-specific SNPs
22 of each gene was used to assess BSB/TC allelic expression.

23

24 **3.6 Allelic expression silencing and bias**

25 After obtaining BSB/TC allelic expression in hybrids, we detected allelic expression
26 silencing in different tissues and generations of reciprocal cross hybrids. This screening
27 complied with the thresholds in which the total number of BSB and TC allelic reads for each
28 gene in both parents was larger than four (≥ 5), either the number of BSB allelic reads or the
29 number of TC allelic reads in the hybrid was zero in the three biological replicates, and the total
30 number of other BSB or TC allelic reads was greater than four (≥ 5).

31 Equal expression of both alleles and unbalanced expression of the two alleles were used
32 for simple classification of hybrids. For this analysis, expression levels were considered “equal”
33 if the \log_2 -transformed ratio of allelic expression of TC and BSB ($|\log_2 (TC/BSB)|$) in hybrid
34 samples was less than 1. Genes with $|\log_2 (TC/BSB)|$ values in the hybrid sample greater than
35 or equal to 1 were considered “unbalanced”, and TC and BSB allelic expression bias was
36 classified based on the plus-minus of these values. Additionally, we chose another APE value
37 for allelic analysis based on the \log_2 -transformed ratio of allelic expression of TC and BSB in

1 the two parental samples, which was used as a reference for the beginning of changes in allelic
2 expression.

3

4 **3.7 Cis- and trans-regulatory differences underlying expression divergence between BSB**
5 **and TC**

6 To further investigate the mechanism of expression divergence that arises from changes in
7 TC and BSB alleles, *cis*- and/or *trans*-regulatory patterns were established based on significant
8 differences between TC and BSB in parents and hybrids, as described in McManus *et al.*
9 (McManus et al. 2010). In the above analysis, binomial exact tests and Fisher's exact tests, with
10 an FDR of 5%, were used to identify genes with significant differences. In addition, we further
11 sorted the *cis*- and/or *trans*-regulatory results based on the plus-minus log₂-transformed ratio
12 of allelic expression of TC to BSB in parents or hybrids.

13

14 **3.8 Correlation analysis of cis- and trans-regulatory expression, K_a/K_s and allelic**
15 **recombinant genes**

16 To investigate the underlying regulatory mechanism of *cis*- and *trans*-regulatory
17 expression, correlation analysis was performed using Pearson's rank correlation coefficients
18 and Student's *t* test in GraphPad Prism (v 7.0) between ω and $|\log_2(\text{TC/BSB})|$ in both the parents
19 and the hybrids. Furthermore, another correlation analysis based on the Spearman method was
20 used to obtain the correlation coefficient between the absolute value of the difference in CDS
21 length and $|\log_2(\text{TC/BSB})|$ in both the parents and the hybrids. To further investigate *cis*- and
22 *trans*-regulatory expression of genes less likely than others to be influenced by allelic
23 recombinant events, a chi-square test was used to assess the correlation between the distribution
24 of allelic recombinant events and the "cis only" gene in each sample.

25

26 **References**

27 Boetzer M, Henkel CV, Jansen HJ, Butler D, Pirovano W. 2011. Scaffolding pre-assembled
28 contigs using SSPACE. *Bioinformatics* **27**: 578-579. doi:
29 [10.1093/bioinformatics/btq683](https://doi.org/10.1093/bioinformatics/btq683)

30 Branchetti E, Poggio P, Sainger R, Shang E, Grau JB, Jackson BM, Lai EK, Parmacek MS,
31 Gorman RC, Gorman JH et al. 2013. Oxidative stress modulates vascular smooth
32 muscle cell phenotype via CTGF in thoracic aortic aneurysm. *Cardiovasc Res* **100**: 316-
33 324. doi: [10.1093/cvr/cvt205](https://doi.org/10.1093/cvr/cvt205)

34 Burge C, Karlin S. 1997. Prediction of complete gene structures in human genomic
35 DNA. Edited by F. E. Cohen. *J Mol Biol* **268**: 78-94. doi: [10.1006/jmbi.1997.0951](https://doi.org/10.1006/jmbi.1997.0951)

36 Burton JN, Adey A, Patwardhan RP, Qiu R, Kitzman JO, Shendure J. 2013. Chromosome-scale
37 scaffolding of de novo genome assemblies based on chromatin interactions. *Nat*

1 *Biotechnol* **31**: 1119. doi: 10.1038/nbt.2727

2 Conesa A, Götz S, García-Gómez JM, Terol J, Talón M, Robles M. 2005. Blast2GO: a universal
3 tool for annotation, visualization and analysis in functional genomics research.
4 *Bioinformatics* **21**: 3674-3676. doi: 10.1093/bioinformatics/bti610

5 Dillies M-A, Rau A, Aubert J, Hennequet-Antier C, Jeanmougin M, Servant N, Keime C, Marot
6 G, Castel D, Estelle J. 2013. A comprehensive evaluation of normalization methods for
7 Illumina high-throughput RNA sequencing data analysis. *Brief Bioinform* **14**: 671-683.
8 doi: 10.1093/bib/bbs046

9 Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, Batut P, Chaisson M, Gingeras
10 TR. 2013. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**: 15-21. doi:
11 10.1093/bioinformatics/bts635

12 Edgar RC. 2004. MUSCLE: multiple sequence alignment with high accuracy and high
13 throughput. *Nucleic Acids Res* **32**: 1792-1797. doi: 10.1093/nar/gkh340

14 Francesc P, Andy B, Jeanclaude FR, Martin FH, Pierrick H, Lorenzo C. 2009. Polyploid fish
15 and shellfish: Production, biology and applications to aquaculture for performance
16 improvement and genetic containment. *Aquaculture* **293**: 125-156. doi:
17 10.1016/j.aquaculture.2009.04.036

18 Gibson G, Rileyberger R, Harshman L, Kopp A, Vacha S, Nuzhdin S, Wayne M. 2004.
19 Extensive sex-specific nonadditivity of gene expression in *Drosophila melanogaster*.
20 *Genetics* **167**: 1791-1799. doi: 10.1534/genetics.104.026583

21 Gnerre S, Maccallum I, Przybylski D, Ribeiro FJ, Burton JN, Walker BJ, Sharpe T, Hall G,
22 Shea TP, Sykes S. 2011. High-quality draft assemblies of mammalian genomes from
23 massively parallel sequence data. *Proc Natl Acad Sci USA* **108**: 1513-1518. doi:
24 10.1073/pnas.1017351108

25 Gong DB, Xu LH, Wu C, Wang S, Liu QF, Cao L, Mao ZW, Wang YD, Hu FZ, Zhou R et al.
26 2019. Two types of gynogenetic blunt snout bream derived from different sperm.
27 *Aquaculture* **511**: 734250. doi: 10.1016/j.aquaculture.2019.734250

28 Grabherr MG, Haas BJ, Yassour M, Levin JZ, Thompson DA, Amit I, Adiconis X, Fan L,
29 Raychowdhury R, Zeng Q. 2011. Full-length transcriptome assembly from RNA-Seq
30 data without a reference genome. *Nat Biotechnol* **29**: 644-652. doi: 10.1038/nbt.1883

31 Haas BJ, Salzberg SL, Wei Z, Pertea M, Allen JE, Orvis J, White O, Buell CR, Wortman JR.
32 2008. Automated eukaryotic gene structure annotation using EVidenceModeler and the
33 Program to Assemble Spliced Alignments. *Genome Biol* **9**: R7. doi: 10.1186/gb-2008-
34 9-1-r7

35 Harris RS. 2007. *Improved pairwise alignment of genomic DNA*. PhD thesis, Pennsylvania
36 State University, University Park, PA.

37 Huang J, Elicker J, Bowens N, Liu X, Cheng L, Cappola TP, Zhu X, Parmacek MS. 2012.

1 Myocardin regulates BMP10 expression and is required for heart development. *J Clin*
2 *Invest* **122**: 3678-3691. doi: 10.1172/JCI63635

3 Huang J, Min Lu M, Cheng L, Yuan LJ, Zhu X, Stout AL, Chen M, Li J, Parmacek MS. 2009.
4 Myocardin is required for cardiomyocyte survival and maintenance of heart function.
5 *Proc Natl Acad Sci USA* **106**: 18734-18739. doi: 10.1073/pnas.0910749106

6 Huang J, Parmacek MS. 2012. Modulation of smooth muscle cell phenotype: the other side of
7 the story. *Circ Res* **111**: 659-661. doi: 10.1161/CIRCRESAHA.112.277368

8 Johnson AD, Handsaker RE, Pilit SL, Nizzari MM, O'donnell CJ, De Bakker PI. 2008. SNAP:
9 a web-based tool for identification and annotation of proxy SNPs using HapMap.
10 *Bioinformatics* **24**: 2938-2939. doi: 10.1093/bioinformatics/btn564

11 Keilwagen J, Wenk M, Erickson JL, Schattat MH, Grau J, Hartung F. 2016. Using intron
12 position conservation for homology-based gene prediction. *Nucleic Acids Res* **44**: e89.
13 doi: 10.1093/nar/gkw092

14 Kent WJ. 2002. BLAT--the BLAST-like alignment tool. *Genome Res* **12**: 656-664. doi:
15 10.1101/gr.229202

16 Krzywinski MI, Schein JE, Birol I, Connors J, Gascoyne R, Horsman D, Jones SJ, Marra MA.
17 2009. Circos: an information aesthetic for comparative genomics. *Genome Res.* doi:
18 10.1101/gr.092759.109

19 Li H, Durbin R. 2010. Fast and accurate long-read alignment with Burrows-Wheeler transform.
20 *Bioinformatics* **26**: 589-595. doi: 10.1093/bioinformatics/btp698

21 Li J, Bowens N, Cheng L, Zhu X, Chen M, Hannenhalli S, Cappola TP, Parmacek MS. 2012.
22 Myocardin-like protein 2 regulates TGFbeta signaling in embryonic stem cells and the
23 developing vasculature. *Development* **139**: 3531-3542. doi: 10.1242/dev.082222

24 Li L, Stoeckert CJ, Roos DS. 2003. OrthoMCL: identification of ortholog groups for eukaryotic
25 genomes. *Genome Res* **13**: 2178-2189. doi: 10.1101/gr.1224503

26 Liu D, Ma C, Hong W, Huang L, Liu M, Liu H, Zeng H, Deng D, Xin H, Song J. 2014.
27 Construction and analysis of high-density linkage map using high-throughput
28 sequencing data. *Plos One* **9**: e98855. doi: 10.1371/journal.pone.0098855

29 Liu S, Luo J, Chai J, Ren L, Zhou Y, Huang F, Liu X, Chen Y, Zhang C, Tao M et al. 2016.
30 Genomic incompatibilities in the diploid and tetraploid offspring of the goldfish x
31 common carp cross. *Proc Natl Acad Sci USA* **113**: 1327-32. doi:
32 10.1073/pnas.1512955113

33 Lowe TM, Chan PP. 2016. tRNAscan-SE On-line: integrating search and context for analysis
34 of transfer RNA genes. *Nucleic Acids Res* **44**: W54-W57. doi: 10.1093/nar/gkw413

35 Luo R, Liu B, Xie Y, Li Z, Huang W, Yuan J, He G, Chen Y, Qi P, Liu Y. 2012. SOAPdenovo2:
36 an empirically improved memory-efficient short-read de novo assembler. *GigaScience*
37 **1**: 18-18. doi: 10.1186/2047-217X-1-18

1 Majoros WH, Pertea M, Salzberg SL. 2004. TigrScan and GlimmerHMM: two open source ab
2 initio eukaryotic gene-finders. *Bioinformatics* **20**: 2878-2879. doi:
3 10.1093/bioinformatics/bth315

4 McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, Garimella K,
5 Altshuler D, Gabriel S, Daly M et al. 2010. The Genome Analysis Toolkit: a
6 MapReduce framework for analyzing next-generation DNA sequencing data. *Genome
7 Res* **20**: 1297-1303. doi: 10.1101/gr.107524.110

8 McManus CJ, Coolon JD, Duff MO, Eipper-Mains J, Graveley BR, Wittkopp PJ. 2010.
9 Regulatory divergence in *Drosophila* revealed by mRNA-seq. *Genome Res* **20**: 816-
10 825. doi: 10.1101/gr.102491.109

11 Nawrocki EP, Eddy SR. 2013. Infernal 1.1: 100-fold faster RNA homology searches.
12 *Bioinformatics* **29**: 2933-2935. doi: 10.1093/bioinformatics/btt509

13 Nguyen AT, Gomez D, Bell RD, Campbell JH, Clowes AW, Gabbiani G, Giachelli CM,
14 Parmacek MS, Raines EW, Rusch NJ et al. 2013. Smooth muscle cell plasticity: fact or
15 fiction? *Circ Res* **112**: 17-22. doi: 10.1161/CIRCRESAHA.112.281048

16 Parmacek MS, Epstein JA. 2009. Cardiomyocyte renewal. *N Engl J Med* **361**: 86-88. doi:
17 10.1056/NEJM McB0903347

18 Parmacek MS., Epstein JA. 2013. An epigenetic roadmap for cardiomyocyte differentiation.
19 *Circ Res* **112**: 881-883. doi: 10.1161/CIRCRESAHA.113.301134

20 Parra G, Blanco E, Guigo R. 2000. GeneID in *Drosophila*. *Genome Res* **10**: 511-515. doi:
21 10.1101/gr.10.4.511

22 Quackenbush J. 2002. Microarray data normalization and transformation. *Nat Genet* **32**: 496-
23 501. doi: 10.1038/ng1032

24 Rhoads A, Au KF. 2015. PacBio Sequencing and Its Applications. *Genomics Proteomics
25 Bioinformatics* **13**: 278-289. doi: 10.1016/j.gpb.2015.08.002

26 Robinson MD, McCarthy DJ, Smyth GK. 2010. edgeR: a Bioconductor package for differential
27 expression analysis of digital gene expression data. *Bioinformatics* **26**: 139-140. doi:
28 10.1093/bioinformatics/btp616

29 Rong L, Liu J, Qi Y, Graham AM, Parmacek MS, Li S. 2012. GATA-6 promotes cell survival
30 by up-regulating BMP-2 expression during embryonic stem cell differentiation. *Mol
31 Biol Cell* **23**: 3754-3763. doi: 10.1091/mbc.E12-04-031

32 Sabio G, Davis RJ. 2010. cJun NH2-terminal kinase 1 (JNK1): roles in metabolic regulation of
33 insulin resistance. *Trends Biochem Sci* **35**: 490. doi: 10.1016/j.tibs.2010.04.004

34 Schaeafke B, Emerson JJ, Wang TY, Lu MY, Hsieh LC, Li WH. 2013. Inheritance of gene
35 expression level and selective constraints on trans- and cis-regulatory changes in yeast.
36 *Mol Biol Evol* **30**: 2121-2133. doi: 10.1093/molbev/mst114

37 Shen D, Li J, Lepore JJ, Anderson TJ, Sinha S, Lin AY, Cheng L, Cohen ED, Roberts JD, Jr.,

1 Dedhar S et al. 2011. Aortic aneurysm generation in mice with targeted deletion of
2 integrin-linked kinase in vascular smooth muscle cells. *Circ Res* **109**: 616-628. doi:
3 10.1161/CIRCRESAHA.110.239343

4 Stamatakis A. 2014. RAxML version 8: a tool for phylogenetic analysis and post-analysis of
5 large phylogenies. *Bioinformatics* **30**: 1312-1313. doi: 10.1093/bioinformatics/btu033

6 Stanke M, Waack S. 2003. Gene prediction with a hidden Markov model and a new intron
7 submodel. *Bioinformatics* **19 Suppl 2**: ii215-225. doi: 10.1093/bioinformatics/btu033

8 Varet H, Brillet-Gueguen L, Coppee JY, Dillies MA. 2016. SARTools: A DESeq2- and EdgeR-
9 Based R Pipeline for Comprehensive Differential Analysis of RNA-Seq Data. *PLoS
One* **11**: e0157022. doi: 10.1371/journal.pone.0157022

10 Wang Y, Tang H, DeBarry JD, Tan X, Li J, Wang X, Lee T-h, Jin H, Marler B, Guo H. 2012.
11 MCSanX: a toolkit for detection and evolutionary analysis of gene synteny and
12 collinearity. *Nucleic Acids Res* **40**: e49-e49. doi: 10.1093/nar/gkr1293

13 Xiao J, Hu F, Luo K, Li W, Liu S. 2016. Unique nucleolar dominance patterns in distant hybrid
14 lineage derived from *Megalobrama Amblycephala* \times *Culter Alburnus*. *BMC Genet* **17**:
15 150. doi: 10.1186/s12863-016-0457-3.

16 Xu H, Luo X, Qian J, Pang X, Song J, Qian G, Chen J, Chen S. 2012. FastUniq: A Fast De
17 Novo Duplicates Removal Tool for Paired Short Reads. *Plos One* **7**: e52249. doi:
18 10.1371/journal.pone.0052249

19 Yang Z. 2007. PAML 4: phylogenetic analysis by maximum likelihood. *Mol Biol Evol* **24**: 1586-
20 1591. doi: 10.1093/molbev/msm088

21

22