# Supplemental Methods

*Collection and Preparation*

After extraction, the DNA samples were sheared with a 1.5 blunt end needles (Jensen Global, Santa Barbara, CA, USA) and run through pulse field gel electrophoresis, in order to separate large DNA molecules, for 16 hours. Samples were checked with a spectrophotometer (Synergy H1 Hybrid Reader, BioTek Instruments Inc., Winooski, VT) and Qubit® fluorometer for quantification of our genomic DNA.

*Illumina and Pacbio Hybrid Assembly, Genome Size Estimation, and Quality Assessment*

All computational and bioinformatics analyses were conducted on the High Performance Computing (HPC) Cluster located that University of California, Irvine.  Sequence data generated from two lanes of Illumina HiSeq 2500 were concatenated and raw sequence reads were assembled through PLATANUS v1.2.1 (*1*), which accounts for heterozygous diploid sequence data. Parameters used for PLATANUS was -m 256 (memory) and -t 48 (threads) for this initial assembly. Afterwards, contigs assembled from PLATANUS and reads from 40 SMRT cells of PacBio sequencing were assembled with a hybrid assembler DBG2OLC v1.0 (*2*). We used the following parameters in DBG2OLC: k 17 KmerCovTh 2 MinOverlap 20 AdaptiveTh 0.01 LD1 0 and RemoveChimera 1 and ran pbdagcon with default parameters. Without Illumina sequence reads, we also conducted a PacBio reads only assembly with FALCON v0.3.0 (https://github.com/PacificBiosciences/FALCON) with default parameters in order to assemble PacBio reads into contiguous sequences. The parameters we used were input type = raw, length_cutoff = 4000, length_cutoff_pr = 8000, with different cluster settings 32, 16, 32, 8, 64,

and 32 cores, concurrency setting jobs were 32, and the remaining were default parameters. After our FALCON assembly, we used the outputs from FALCON and DBG2OLC as input for QUICKMERGE v1.0 (*3*), a metassembler and assembly gap filler developed for long molecule-based assemblies. Several different parameters in QUICKMERGE v1.0 were conducted as suggested by the authors until an optimal assembly was obtained with HCO 10 C 3.0 -l 2400000 -ml 5000.  With the final genome assembly, we processed our genome through REPEATMASKER v.4.0.6 (*4*) to mask repetitive elements with the parameter -species teleostei.  We also estimated transposable elements (TE) content by using the fugu repeat database in REPEATMASKER.  Genome size was estimated with only Illumina sequences and tandem repeats were detected with bioinformatic tools (electronic supplemental material). We used BUSCO v3 (*5*) to estimate the completeness of our genome assembly with the vertebrata and Actinopterygii gene set [consists of 2,586 Benchmarking Universal Single-Copy Orthologs (BUSCOs)] to estimate completeness of our *C. violaceus* genome.

*RNA-Seq Tissue Extraction and Sequencing*

Five individual *C. violaceus* were collected during the fall of 2015 for our transcriptomic analyses and annotation of the *C. violaceus* genome. We extracted brain, gill, gonads (testes), heart, liver, mid intestine, proximal intestine, pyloric caeca, and spleen from each individual and preserved the tissues in RNAlater® (Ambion, Austin, TX, USA). All individuals had digesta in their guts during dissection (i.e., they had all eaten), and digesta was removed prior to tissue fixation in RNAlater®. Total RNA was extracted using a Trizol protocol, and sample quality (RNA Integrity Number $\geq 8.2$) confirmed using an Agilent bioanalyzer 2100 (RNA nano chip; Agilent Technologies). We used an Illumina TruSeq Sample Preparation v2 (Illumina) kit with

AMPure XP beads (Beckman Coulter Inc.) and SuperScript™ III Reverse Transcriptase
(Invitrogen) to prepare our tissue samples for Illumina sequencing.  See Supplemental Table S8
for adaptor indexes used for the transcriptomic sequence data, which was performed using
multiplexed samples at 10 nM in 10 μl, and sequenced on two lanes on an Illumina HiSeq 2500
(100 bp Paired Ends) at UC Irvine's GHTF.

*Transcript Assembly, Annotation, and Heatmap Generation of All Tissues and Genes Associated
with Diet*

Reads were normalized, assembled, annotated and gene expression of all transcripts were
measured from all nine tissue types and TRINITY v2.3.2 (*6*) was used to identify how many reads
mapped back to the *C. violaceus* genome (Supplemental Figure S12 and Table S8). Differentially
expressed genes for all tissue types were viewed with a heatmap that was generated with the
CUMMERBUND  R package (http://compbio.mit.edu/cummeRbund/; Supplemental Figure S13).
Candidate genes which pertained to glycolytic, lipid metabolism/gluconeogenesis, ketone
degradation, glucosidases (both α and β), proteases, and lipases were identified in the *C.*
*violaceus* transcriptome by scanning the annotation of CUFFLINKS assembled transcripts and used
to generate our heatmap.

*Genome Size Estimation and Identification of Tandem Repeats*

The c-value has been estimated for *Cebidichthys violaceus* (*7*), which is 0.81. Based on
this c-value, the estimate of the genome size is ~792 Mb. In addition, we estimated the genome
size using only Illumina sequences by using JELLYFISH v2.2.0 (*8*). We selected multiple k-mers
(25, 27, 29, 31) for counting and generating a histogram of the k-mer frequencies. We used a perl

script (written by Joseph Ryan) to estimate genome size based on k-mer sizes and peak values determined from histograms generated in JELLYFISH.

We used tandem repeats finder (trf v4.07b; *9*) to identify tandem repeats throughout the unmasked genome. We used the following parameters in trf 1 1 2 80 5 200 2000 -d -h to identify repeats. Once the largest repeats were identified, we used period size of the repeats multiplied by the number of copies of the repeat to generate the largest fragments. This method was used to identify repetitive regions which can possibly represent centromere or telomere regions of the *C. violaceus* genome.

*Transcript Assembly for all Tissues and Annotation*

The following pipeline was used to assemble and measure expression of all transcripts from all nine tissue types (Supplemental Figure S12). Prior to assembly, all raw reads were trimmed with TRIMMOMATIC v0.35 (*10*). Afterwards, trimmed reads were normalized using a perl script provided by TRINITY v r2013-02-16 (*11*). Prior to aligning transcriptomic reads to the genome, the final masked assembled genome was prepared with BOWTIE2-BUILD v2.2.7 (*12*) for a BOWTIE index and then all (normalized) reads from each tissue type were mapped using TOPHAT v2.1.0 (*13*) to our assembled masked genome using the following parameters -I 1000 -i 20 -p 4. Afterwards, aligned reads from each tissue were indexed with SAMTOOLS v1.3 (*14*) as a BAM file. Once indexed through SAMTOOLS, transcripts were assembled by using CUFFLINKS v2.2.1 (*15*) with an overlap-radius 1. All assemblies were merged using CUFFMERGE and then differential expression was estimated with CUFFDIFF, both programs are part of the CUFFLINKS package. All differential expression analyses and plots were produced in R (https://www.r-

project.org/) using CUMMERBUND tool located on the bioconductor website

(https://www.bioconductor.org/). Once all transcripts were assembled, we ran REPEATMASKER

v.4.0.6 with the parameter -species teleostei to mask repetitive elements within our

transcriptomes.

All masked transcripts were annotated with the trinotate annotation pipeline

(https://trinotate.github.io/), which uses Swiss-Prot (*16*), Pfam (*17*), eggNOG (*18*), Gene

Ontology (*19*), SignalP (*20*), and Rnammer (*21*). We also processed our transcripts through

BLASTX against the UniProt database (downloaded on September 26th, 2017) with the following

parameters: num_threads 8, evalue 1e-20, and max_target_seqs 1. The BLASTX output was

processed through trinity analyze_blastPlus_topHit_coverage.pl script to count the amount of

transcripts of full length or near full length. To provide a robust number of full-length transcripts,

the assembled genome was processed through AUGUSTUS v3.2.1 (*22*) without hints using default

parameters for gene predictions using a generalized hidden Markov model in order to identify

genes throughout the genome, and predicted transcripts were also masked for repetitive elements

through REPEATMASKER.

*Comparative Analysis for Syntenic Regions across Teleosts fishes*

We selected the following fish genomes: zebrafish (*Danio rerio*), stickleback

(*Gasterosteus aculeatus*), spotted gar (*Lepisosteus oculatus*), and the Japanese medaka (*Oryzias

latipes*) to identify syntenic regions with our *C. violaceus* genome assembly. All genomes were

masked using REPEATMASKER v.4.0.6 with the parameter -species teleostei. To identify syntenic

regions we used contigs from our *C. violaceus* which represent 1MB or larger and then

concatenated the remaining contigs. We used SATSUMA v3.1.0 (*23*) with the following
parameters -n 4 -m 8 for identifying syntenic regions between zebrafish, stickleback, spotted gar,
Japanese medaka and the *C. violaceus* genome. Afterwards, we developed circos plots to view
syntenic regions shared between species using CIRCOS v0.63-4 (*24*).


*Identification of Orthologs Across Teleost Fishes and Identification of Syntenic Regions*

The following teleost and non-teleost genomes were taken from ENSEMBL release 89 for our
comparative analysis of orthologs and phylogeny of fishes; *Poecilia formosa* (Amazon molly),
*Astyanax mexicanus* (blind cave fish), *Gadus morhua* (Atlantic Cod), *Takifugu rubripes* (Fugu),
*Oryzias latipes* (Japanese medaka), *Xiphophorus maculatus* (platyfish), *Lepisosteus oculatus*
(spotted gar), *Gasterosteus aculeatus* (stickleback), *Tetraodon nigroviridis* (green spotted
puffer), *Oreochromis niloticus* (Nile tilapia), *Danio rerio* (zebrafish), *Latimeria chalumnae*
(coelacanth).

We used INPARANOID v4.0 (*25*) to conduct 78 possible pairwise comparisons, where N is number
of taxa [(N(N-2))/2= possible pairwise comparisons]. From the outputs of INPARANOID, we used
QUICKPARANOID (http://pl.postech.ac.kr/QuickParanoid/) to identify orthologous clusters from all
13 species.


# Supplemental Results and Discussion


*I - Genome and Transcriptome Assembly*

From PacBio sequencing, we were able to generate ~29,700 Mb of sequence data from
40 SMRT cells, this represents ~37X coverage based on the c-value estimated for *C. violaceus*

(792 Mb; *7*). From two lanes of Illumina sequencing we were able to generate 84,539 Mb which represents ~107X coverage of the genome. From the PLATANUS assembly with Illumina only sequence data, our N50 was 2,760 bp. When combining PLATANUS assembled contigs and 40 SMRT cells of PacBio for a hybrid assembly in DBG2OLC, we managed to obtain an N50 of 2.21Mb. Afterwards, when using FALCON (PacBio only reads) we managed to obtain an N50 of 2.45Mb. We used the FALCON assembly and the DBG2OLC assembly and through QUICKMERGE, we obtained an N50 of 6.69 Mb. Following two rounds of QUIVER, and PILON, we assembled our final draft genome which composed of 467 contigs, and obtained N25, N50, and N75 values of 15.17, 6.71, and 1.85 (Mb) respectively which were composed of 593,001,491 base pairs (Supplemental Figure S4). Genomic regions were masked and a summary of the identities of repetitive elements are present in Supplemental Table S4.

*II - Estimation of Genome Size, Completeness Assessment, and Tandem Repeats Throughout the Genome*

By using JELLYFISH we estimated the genome size based on an average of four k-mer size counts (25, 27, 29, 31) 656,598,967 base pairs based with a standard deviation of 4,138,853 base pairs (Supplemental Figure S5, Table S7). Through BUSCO v3 there were 97% (2,508 genes out of 2,586) complete orthologs detected which included 1.3% duplicated orthologs. In addition, there were 1.1% (28 BUSCOs) partial orthologs present and 1.9% (50 BUSCOs) of orthologs were not detected in the *C. violaceus* genome (Supplemental Table S3).

In identifying large genomic regions of tandem repeats, we were able to identify about 38,448 repeated loci where the repetitive sequence (period) range from 1 to 1,983 and number of

repeats identified were 1.8 to 14,140.8bps. The largest repeat locus (period size multiplied by the repeat amount) was 109,161bps with a period of 90bp and repeat amount of 1,212.9bps. We saw an increase in size for 35 loci which had repeat locus the size of 32,594.1bps and greater (Supplemental Figure S6-S7) as compared to any other repetitive locus.

*III – Assembly and Annotation of Nine Tissue Transcriptomes*

From five individuals that we selected for our transcriptomic analyses, the total reads mapped back to the genome ranged from 67.37% (liver) to 84.8% (heart; Supplemental Table S8). The range of transcripts present in each tissue type ranged from 20,008 (liver) to 78,629 (gill) transcripts (Supplemental Table S8).

From the TUXEDO package, there were 101,922 transcripts estimated from the nine tissues. When evaluating Fragments Per Kilobase per Million mapped Reads (FPKM) for each of the nine transcriptomes, we see the lowest median with the liver and the highest median with the gill tissue (Supplemental Figure S14). All other tissues appeared to have a similar profile. In addition, we see that the gill tissue had a short Quartile group 2 as compared to the other tissue types (except liver tissue; Supplemental Figure S14). When we look at the differentially expressed genes across all tissue types, we see that there are a cluster of genes highly expressed in the liver as compared to some of the tissues (Supplemental Figure S13). By using getSig in the CUMMERBUND package we evaluated which genes are significantly regulated, with the highest number (1,383) between liver and brain, as these tissues are highly specialized. When we look at Jensen-Shannon Distances (Supplemental Figures S16-S17), we see that pyloric caeca and

middle intestine have similar expression profiles, brain and gonad have similar profiles, and the proximal intestine has a very different expression profile.

From the assembled transcripts in TRINITY and predicted transcripts from AUGUSTUS, there were a total of 105,167,222 and 44,120,550 bases with 0.08% and 2.15% of bases masked respectively (Supplemental Tables S9-10). When annotated, there were 65,535 transcripts in TRINOTATE and there were only 26,356 transcripts with a BLASTX hit identified. When conducting a BLASTX on our transcripts to identify full length transcripts in our dataset, we were only able to obtain 5,199 transcripts which had an 80% hit coverage (Supplemental Tables S12). When using only using AUGUSTUS (without hints) and identified 29,525 genes. There were 29,485 genes which had 60 amino acids or greater in our transcriptomic dataset.

*IV - Expression Profiles of Candidate Genes for Digestion and Metabolism*

We were able to identify candidate genes associated with digestion, fermentation, ketone degradation in our transcriptomic assembly (Fig. 2a & b in the main manuscript) and view differential gene expression patterns across the nine tissues where we have transcriptomic dataset (Supplemental Figure S13). We were not able to distinguish between *amy2a* and *amy2b* in our transcriptomic assembly. Therefore, we used the AUGUSTUS gene prediction from the genome as a transcriptome reference to detect the *amy2a* and *amy2b* genes and mapped our transcriptome reads back to this reference transcriptome dataset. From this dataset, we were able to detect *amy2a* and *amy2b* gene expression profiles. As expected, we see high expression profiles for genes associated with digestion and metabolism in the pyloric caeca, proximal intestine, mid intestine, and liver (Fig. 2a & b).

*V - Evaluation of Candidate Genes Associated with Digestion*

From our MUMMER and BLAST search for pancreatic amylase, we have identified three tandem copies of amylase (*amy2a*) and (*amy2b*), as opposed to the six haploid copies detected in the German et al. (*26*). We identified amylase on contig 440 and we see two hypothetical proteins between the three amylase genes and a transposase near *amy2b* (Fig. 3b; Supplemental Figure S26). In addition, each amylase gene is preceded by a 4.3K20bp DNA element encoding a transposase (Fig. 3b, Supplemental Figure S26). The three tandem amylase loci differ from the estimated six haploid copies (based on gene dosage curves using RT-qPCR) proposed to be present in the *C. violaceus* genome by German et al. (*26*). Upon further inspection, we have reached the conclusion that the per cell gene count of the German et al. (*26*) study is the diploid copy number, not the haploid copy number (*C. violaceus* is a diploid, vertebrate). Hence, the copy number based on gene dosage curves is three for amylase in general, with roughly two copies for *amy2a* and one copy for *amy2b* (*26*) which agrees completely with what is observed in our genomic assembly. Although there is the possibility for copy number variation amongst individuals within a population, as there is for human salivary amylase (*27*) and dog pancreatic amylase (*28*), that is not what was observed by German et al. (*26*), as that would entail different methodology and more robust sampling of *C. violaceus* individuals.

We only selected one gene copy of *amy2a* because they are identical, and *amy2b* when estimating selection in DATAMONKEY. When testing all 11 branches for seven taxa in aBSREL, we see only one branch under episodic diversifying selection (*C. violaceus*, *amy2b*; Fig. 3c). We do not see this pattern of positive selection in any of the other branch with a significant p-value.

In MEME, we identified three sites with episodic positive selection with a p-value threshold of 0.05 (sites: 41, 256, and 279; Fig. 3d) and in GARD there was no evidence of recombination.

Our analyses of aminopeptidase (also known as alanyl aminopeptidase, E.C. 3.4.11.2) genes have revealed some interesting results. Fishes appear to have five aminopeptidase genes, which varies from the one (aminopeptidase N) seen in mammals. As we probed genomes of sufficient quality (e.g., using http://ensembl.org), it became clear that fish aminopeptidase loci show signatures of retention following whole genome duplication (WGD) events (*29-32*). The website http://ohnologs.curie.fr lists aminopeptidase a (Supplemental Figure S21) and aminopeptidase b (Supplemental Figure S22b & c) in *Danio rerio* and *Oryzias latipes* as being ohnologs from the vertebrate WGD (*31-32*). This same website then lists aminopeptidase b and aminopeptidase N (Supplemental Figure S22c) as ohnologs from the Teleost-specific WGD event (*29-30*). Our synteny maps for other fishes, including *C. violaceus*, support this contention, especially among aminopeptidase b and aminopeptidase N, as there are other shared genes (e.g., *svp2b*) among the separate loci for aminopeptidase b and aminopeptidase N (Supplemental Figure S22). Indeed, our limited phylogenetic analysis suggests that aminopeptidase a is sister to all other aminopeptidase genes (Supplemental Figure S24). Aminopeptidase b and Ey are more related (Supplemental Figure S22b), and they are sister to a clade that includes aminopeptidase N and Ey-like (Supplemental Figures S22 and S24). The evolutionary history of aminopeptidases clearly requires more work, but our preliminary analysis suggests that the history of aminopeptidases in fishes may involve WGD events. What this means for digestion in fishes with different diets should be explored, as all of the aminopeptidase genes show elevated gut expression, except aminopeptidase N (Fig. 2), which is the name of the alanyl aminopeptidase in mammals.  Aminopeptidase activity can be plastic in fishes fed different diets (*33-35*), and it

doesn't always appear to only be elevated in fishes consuming more protein. Thus, each aminopeptidase protein necessitates investigation into how their functions may vary in the gut and how this may matter for fishes with different diets.

When examining aminopeptidase a (*anpepa*) codons for positive selection, we did not identify any branches under episodic diversifying selection and identified four sites under episodic positive selection (sites: 194, 412, 445, and 593; Supplemental Figure S21).

We found no branches under episodic diversifying selection in aminopeptidase b, N, and Ey (Supplemental Figure S23). For aminopeptidase Ey-like (*anpep-Ey-like*) we found one branch which leads to *C. violaceus* and *A. purpurescens* under episodic diversifying selection and one site with episodic positive selection with a p-value threshold of 0.05 (site: 38; Supplemental Figure S23, Table S13). We also found one site under episodic positive selection in *anpepb* with a threshold of 0.05 (site: 156) and one site in *anpep N* (site: 351).

For Phospholipase B1, plb1-1, we found no branches under selection and six sites under selection and a p-value threshold of 0.05 (sites: 66, 97, 289, 438, 800, and 821; Supplemental Figure S27, Table S13). As for *plb1-2*, we found one branch which leads to *C. violaceus* and *A. purpurescens* under episodic diversifying selection with three sites episodic positive selection with a p-value threshold of 0.05 (sites: 183, 230 and 476). Lastly, we did not see any branches or sites under selection for plb1-3. When evaluating Phospholipase Group 12 B (*pg12b1* & 2; Supplemental Figure S28) we found no branches under positive diversifying selection. Only for pg12b-2, we found nine sites under selection with a p-value threshold of 0.05 (sites: 25, 31, 32, 33, 37, 83, 142, 146, and 185). When evaluating *cel* and *cel*-like and we did not detect any branches under episodic diversifying selection. We only found three sites under episodic positive

selection for cel-1 with a p-value threshold of 0.05 (sites: 64, 258, and 355). For chymotrypsin A (*chymo A*), we did not detect any branches under episodic diversifying selection or sites under episodic positive selection. For chymotrypsin B (*ctrb*), we see two branches with episodic diversifying selection (*C. violaceus* and *A. purpurescens*). There is one site with episodic positive selection (site: 112; Supplemental Figure S18, Table S13). For chymotrypsin-like, we did not detect any branches under episodic diversifying selection or sites under episodic positive selection (Supplemental Figure S19). As for trypsin (*tryp3-1 & -2*), we did not detect any branches under episodic diversifying selection or sites under episodic positive selection except for *tryp3-2* which has one site under episodic positive selection with a threshold of 0.05 (site: 91; Supplemental Figure S20). All candidate genes were evaluated for recombination with gard, which some loci had putative breakpoints, but with the KH (Kishino–Hasegawa test, at $P = 0.1$) test, there were 0 breakpoints with significant topological incongruence (Supplemental Table S13).

We focused on the loci that encode for *cel* and were able to identify the four tandem copies of *cel* on contig 445 (Supplemental Figure S29). We estimated selection and only see evidence episodic selection of sites on *cel*-1 genes (Fig. 4e in the main manuscript).

With regards to elevated lipase activity in the taxa consuming more fiber in their diets, it is known that in industry settings, adding microcrystalline fiber to lipolytic reactions acts to stabilize the lipase proteins and can increase lipase activities (*36-37*). However, these fibrous compounds added directly to the reactive environment. In our case, when measuring lipase activities, we homogenize the tissues separate from gut contents, and then centrifuge the homogenates to get the supernatant, which would contain only those enzymes that are soluble and not bound to larger molecules, like fiber. Hence, the elevated lipase activities measured *in*

*vitro* in our investigations cannot be coming from the potential effects of fiber on the lipase proteins themselves, even if these interactions might act to aid lipolytic action *in vivo* within the gut environment. We are, therefore, confident that the increased lipase activities in the algae-eating fishes are due to the molecular differences in CEL proteins (Fig. 4; Supplemental Figure S29-S31), and any differences in gene expression.

*VIII - Orthologs and Phylogenetic Analyses*

We constructed a phylogenetic tree using maximum likelihood using thirteen fish taxa including the *Cebidichthys violaceus* which included 30 loci and 33,508 bases with 1,000 bootstrap replicates (Fig. 1; Supplemental Table S5). We used JMODELTEST v2.1.0 and with AICc we detected that GTR+I+G was the best model selected for our phylogenetic analyses. All 30 loci used for our phylogenetic analyses were extracted from orthologs detected in INPARANOID (Supplemental Table S6).

*VI - Syntenic regions across multiple fish species*

We have 114 contigs which have a 1MB or greater, and we pooled all contigs which had less than 1MB (353 contigs) were merged together in our synteny analyses. When we compare our assembled genome to the *G. aculeatus* genome, we see multiple homologous regions between the two species, and multiple loci from each linkage group of the *G. aculeatus* genome represented in the *C. violaceus* genome; Supplemental Figure S8). When comparing the *O. latipes, D. rerio, L. oculatus*, and genomes to our *C. violaceus* genome (Supplemental Figures

S9-S11), we also see each chromosome/linkage group represented in the *C. violaceus* genome and strong synteny between *O. latipes* and *C. violaceus* whereas we less homologous strands between the *C. violaceus* genome and the *D. rerio* or *L. oculatus* genomes.


*VIII - Opsin Gene Copies and Selection*

After reviewing our BUSCO analyses for gene duplicates present in our *C. violaceus* genome, we identified three Opsin Short Wave Sensitive (*opn1sw*) genes in tandem on contig 443 (Supplemental Figure S32). We find this interesting because *C. violaceus* endures a period of time out of water during low tide, in which Horn and Riegle (*38*) showed that a large *C. violaceus* (~24 cm SL; 92 g) can survive out of water for 37 hours. The ability to survive out of water may require adaptations of vision when exposed to air during low tides, in which we further evaluated the *opn1sw* gene copies for signatures of positive selection. In addition, we identified and compared gene copy numbers of short-wave opsin genes from *Danio rerio*, *Oreochromis niloticus*, and *Gasterosteus aculeatus* genomes, which have one or two gene copies present (Supplemental Figure S32). In addition, we estimated selection by using the DATAMONKEY server v2.0 by using GARD, aBSREL, and MEME. With GARD, we observed evidence of recombination breakpoints, (locations: 196 and 319) but there are 0 breakpoints with significant topological incongruence (p=0.01). From our aBSREL analysis, we observed one branch of episodic diversifying selection out of 11 leading to *opn1sw2a* along with an *opn1sw2* identified in *G. aculeatus* with a corrected p-value of 0.0172. In addition, we detected episodic positive/diversifying selection at 7 sites. The sites under selection with a p-value less than 0.05 were the following: sites 4, 95, 165, 202, 224, and 326, and there was one site under selection

with a p-value less than 0.01 (site 337; Supplemental Figure S32). From this analysis, evaluation of *opn1sw* gene sequences from subtidal and intertidal stichaeids can elucidate how vision may play an important role for intertidal prickleback species.

# References

1. Kajitani R et al. 2014 Efficient de novo assembly of highly heterozygous genomes from whole-genome shotgun short reads. Gen. Res. 24, 1384–1395. (doi:10.1101/gr.170720.113)

2. Ye C, Hill CM, Wu S, Ruan J, Ma Z. 2016 DBG2OLC: Efficient Assembly of Large Genomes Using Long Erroneous Reads of the Third Generation Sequencing Technologies. Sci. Rep. 6. (doi:10.1038/srep31900)

3. Chakraborty M, Baldwin-Brown JG, Long AD, Emerson JJ. 2016 Contiguous and accuratede novoassembly of metazoan genomes with modest long read coverage. Nucl. Acid. Res., gkw654. (doi:10.1093/nar/gkw654)

4. Smit AF. 2004 Repeat-Masker Open-3.0. http://www. repeatmasker. org.

5. Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. 2015 BUSCO: assessing genome assembly and annotation completeness with single-copy orthologs. Bioinformat. 31, 3210–3212. (doi:10.1093/bioinformatics/btv351)

6. Haas BJ et al. 2013 De novo transcript sequence reconstruction from RNA-seq using the Trinity platform for reference generation and analysis. Nature Prot. 8, 1494–1512. (doi:10.1038/nprot.2013.084)

7. Hinegardner R, Rosen DE. 1972 Cellular DNA content and the evolution of teleostean fishes. American Nat.106, 621-644

8. Marçais G, Kingsford C. 2011 A fast, lock-free approach for efficient parallel counting of occurrences of k-mers. Bioinformatics 27, 764–770. (doi:10.1093/bioinformatics/btr011)

9. Benson G. 1999 Tandem repeats finder: a program to analyze DNA sequences. Nucleic Acids Research 27, 573–580. (doi:10.1093/nar/27.2.573)

10. Bolger AM, Lohse M, Usadel B. 2014 Trimmomatic: a flexible trimmer for Illumina sequence data. Bioinformatics 30, 2114–2120. (doi:10.1093/bioinformatics/btu170)

11. Grabherr MG et al. 2011 Full-length transcriptome assembly from RNA-Seq data without a reference genome. Nature Biotechnology 29, 644–652. (doi:10.1038/nbt.1883)

12. Langmead B, Salzberg SL. 2012 Fast gapped-read alignment with Bowtie 2. Nature Methods 9, 357–359. (doi:10.1038/nmeth.1923)

13. Kim D, Pertea G, Trapnell C, Pimentel H, Kelley R, Salzberg SL. 2013 TopHat2: accurate alignment of transcriptomes in the presence of insertions, deletions and gene fusions. Genome Biology 14, R36. (doi:10.1186/gb-2013-14-4-r36)

14. Li H et al. 2009 The Sequence Alignment/Map format and SAMtools. Bioinformatics 25, 2078–2079. (doi:10.1093/bioinformatics/btp352)

15. Trapnell C et al. 2012 Differential gene and transcript expression analysis of RNA-seq experiments with TopHat and Cufflinks. Nature Protocols 7, 562–578. (doi:10.1038/nprot.2012.016)

16. Boeckmann B. 2003 The SWISS-PROT protein knowledgebase and its supplement TrEMBL in 2003. Nucleic Acids Research 31, 365–370. (doi:10.1093/nar/gkg095)

17. Bateman A. 2004 The Pfam protein families database. Nucleic Acids Research 32, 138D–141. (doi:10.1093/nar/gkh121)

18. Powell S et al. 2013 eggNOG v4.0: nested orthology inference across 3686 organisms. Nucleic Acids Research 42, D231–D239. (doi:10.1093/nar/gkt1253)

19. Ashburner M et al. 2000 Gene Ontology: tool for the unification of biology. Nature Genetics 25, 25–29. (doi:10.1038/75556)

20. Petersen TN, Brunak S, von Heijne G, Nielsen H. 2011 SignalP 4.0: discriminating signal peptides from transmembrane regions. Nature Methods 8, 785–786. (doi:10.1038/nmeth.1701)

21. Lagesen K, Hallin P, Rødland EA, Stærfeldt H-H, Rognes T, Ussery DW. 2007 RNAmmer: consistent and rapid annotation of ribosomal RNA genes. Nucleic Acids Research 35, 3100–3108. (doi:10.1093/nar/gkm160)

22. Stanke M, Keller O, Gunduz I, Hayes A, Waack S, Morgenstern B. 2006 AUGUSTUS: ab initio prediction of alternative transcripts. Nucleic Acids Research 34, W435–W439. (doi:10.1093/nar/gkl200)

23. Grabherr MG, Russell P, Meyer M, Mauceli E, Alfoldi J, Di Palma F, Lindblad-Toh K. 2010 Genome-wide synteny through highly sensitive sequence alignment: Satsuma. Bioinformatics 26, 1145–1151. (doi:10.1093/bioinformatics/btq102)

24. Krzywinski M, Schein J, Birol I, Connors J, Gascoyne R, Horsman D, Jones SJ, Marra MA. 2009 Circos: An information aesthetic for comparative genomics. Genome Research 19, 1639–1645. (doi:10.1101/gr.092759.109)

25. O'Brien KP. 2004 Inparanoid: a comprehensive database of eukaryotic orthologs. Nucl. Acid. Res. 33, D476–D480. (doi:10.1093/nar/gki107)

26. German DP, Foti DM, Heras J, Amerkhanian H, Lockwood BL. 2016 Elevated Gene Copy Number Does Not Always Explain Elevated Amylase Activities in Fishes. Physiological and Biochemical Zoology 89, 277–293. (doi:10.1086/687288)

27. Perry GH et al. 2007 Diet and the evolution of human amylase gene copy number variation. Nature Genetics 39, 1256–1260. (doi:10.1038/ng2123)

28. Axelsson E et al. 2013 The genomic signature of dog domestication reveals adaptation to a starch-rich diet. Nature 495, 360–364. (doi:10.1038/nature11837)

29. Christoffels A, Koh EGL, Chia J, Brenner S, Aparicio S, Venkatesh B. 2004 Fugu Genome Analysis Provides Evidence for a Whole-Genome Duplication Early During the Evolution of Ray-Finned Fishes. Molecular Biology and Evolution 21, 1146–1151. (doi:10.1093/molbev/msh114)

30. Glasauer SMK, Neuhauss SCF. 2014 Whole-genome duplication in teleost fishes and its evolutionary consequences. Molecular Genetics and Genomics 289, 1045–1060. (doi:10.1007/s00438-014-0889-2)

31. Kasahara M. 2007 The 2R hypothesis: an update. Current Opinion in Immunology 19, 547–552. (doi:10.1016/j.coi.2007.07.009)

32 Ohno S. 1970 Evolution by gene duplication. New York, Springer-Verlag.

33. German DP, Horn MH, Gawlicka A. 2004 Digestive Enzyme Activities in Herbivorous and Carnivorous Prickleback Fishes (Teleostei: Stichaeidae): Ontogenetic, Dietary, and Phylogenetic Effects. Physiological and Biochemical Zoology 77, 789–804. (doi:10.1086/422228)

34 Harpaz S, Uni Z. 1999 Activity of intestinal mucosal brush border membrane enzymes in relation to the feeding habits of three aquaculture fish species. Comparative Biochemistry and Physiology Part A: Molecular & Integrative Physiology 124, 155–160. (doi:10.1016/s1095-6433(99)00106-3)

35. Leigh SC, Nguyen-Phuc B-Q, German DP. 2017 The effects of protein and fiber content on gut structure and function in zebrafish (Danio rerio). Journal of Comparative Physiology B 188, 237–253. (doi:10.1007/s00360-017-1122-5)

36. Cai W, Liang X, Yuan X, Liu L, He S, Li J, Li B, Xue M. 2018 Different strategies of grass carp (Ctenopharyngodon idella) responding to insufficient or excessive dietary carbohydrate. Aquaculture 497, 292–298. (doi:10.1016/j.aquaculture.2018.07.042)

37. Kim JH, Park S, Kim H, Kim HJ, Yang Y-H, Kim YH, Jung S-K, Kan E, Lee SH. 2017 Alginate/bacterial cellulose nanocomposite beads prepared using *Gluconacetobacter xylinus* and their application in lipase immobilization. Carbohydrate Polymers 157, 137–145. (doi:10.1016/j.carbpol.2016.09.074)

38. Horn MH, Riegle KC. 1981 Evaporative water loss and intertidal vertical distribution in relation to body size and morphology of stichaeoid fishes from California. Journal of Experimental Marine Biology and Ecology 50, 273–288. (doi:10.1016/0022-0981(81)90054-x)

**Supplemental Figure S1: Theoretical activity level of a digestive enzyme in an animal's gut as a function of ingested concentrations of substrate for that enzyme.** Two hypotheses are commonly invoked in the literature to explain patterns of digestive enzyme activities observed frequently in nature. The Adaptive Modulation Hypothesis (Karasov, 1992) suggests a positive correlation between substrate concentration and enzyme activities; an abundant substrate should invoke more enzyme activity ensure adequate digestion of that substrate. This is commonly seen for carbohydrases, like amylase (German et al. 2016). The Nutrient Balancing Hypothesis (Clissold et al. 2010) suggests that an animal should invest in elevated enzyme activities against limiting nutrients that are low in concentration to ensure acquisition of these important nutrients. This is seen in transporter density for some minerals (e.g., iron), and for lipases to digest lipids in herbivorous fishes (German et al. 2004, 2015).

**Supplemental Figure S2: Phylogenetic relationships of the polyphyletic family Stichaeidae based on 2,100 bp of *cytb*, *16s*, and *tomo4c4* genes (Kim et al. 2014).** Bayesian posterior probabilities are indicated on nodes. *Cebidichthys violaceus* is bolded, and photos of *C. violaceus* and other studied taxa are shown with their digestive systems beneath their bodies. Note the differences in gut size. H=herbivory, O=omnivory, C=carnivory. Evolution of herbivory (— — — —) and omnivory (...........) are shown. Numbers in parentheses show number of taxa evaluated at that branch. Boxes highlight alleged families or subfamilies within the polyphyletic family Stichaeidae, with Cebidichthyidae (top), Xiphisterinae (middle), and Alectriinae (bottom) all highlighted. Hindgut short chain fatty acid (SCFA) concentrations are mean ± standard deviation, and were compared with ANOVA ($F_{3,33}$ = 127.92; $P <$ 0.001). SCFA data sharing a superscript letter are not significantly different (from German et al. 2015).

**Supplemental Figure S3: Flowchart of our final genome assembly from Illumina (two lanes of PE 100bp) and Pacific Biosciences (40 single molecule real-time [SMRT] cells) sequence reads.** Light blue boxes indicate raw sequence reads and dark blue boxes indicate bioinformatic programs used for the assembly. (…………) indicates the type of sequence information that was used for the assembly method. Arrows indicate the next step taken to proceed in the genome assembly of *Cebidichthys violaceus*.

**Supplemental Figure S4: Summarization of contig lengths of the *Cebidichthys violaceus* genome.** N25 (blue bars), N50 (blue + green bars), and N75 (blue + green + orange bars) values was estimated for the *C. violaceus* genome. There are 66 contigs which represent the N75 of the *C. violaceus* genome. Contig ID is labeled along the x-axis.

**Supplemental Figure S5: K-mer frequency for *Cebidichthys violaceus*.** Use of raw Illumina (only) reads from *C. violaceus* gDNA to estimate the *C. violaceus* genome size. K-mer sizes of 25, 27, 29, and 31 were selected to generate histograms.

**Supplemental Figure S6: Period size and copies of pattern of tandem repeats identified in the *Cebidichthys violaceus* genome.** The length of the tandem repeat sequence and the size of the repeat found.

**Supplemental Figure S7: Histogram of period size multiplied by number of copies of repeat.** Total of the 200 longest repeats identified, where blue indicates values of 30,877 base pairs (bps) or less. Green indicates values greater than 30,877 bps.

**Supplemental Figure S8: Circos plot showing synteny between the assembled genome**
*Cebidichthys violaceu*s and *Gasterosteus aculeatus* **(three-spined stickleback).** There are 21
chromosomes (green boxes) which represent the three-spined stickleback genome. There are 114
blue boxes which are 1 MB or greater that represent the *C. violaceus* genome. There are 353
contigs that are less than 1 MB which were concatenated into box labeled as 115. Gray strands
indicate syntenic regions between the two genomes. Both *G. aculeatus* and *C. violaceus*
illustrations were drawn by Andrea Dingeldein.

**Supplemental Figure S9: Circos plot showing synteny between the assembled genome *Cebidichthys violaceus* and *Oryzias latipes* (Japanese rice fish).** There are 24 chromosomes (red boxes) which represent the Japanese rice fish genome. There are 114 blue boxes which are 1 MB or greater that represent the *C. violaceus* genome. There are 353 contigs that are less than 1 MB which were concatenated into box labeled as 115. Gray strands indicate syntenic regions between the two genomes. Both *O. latipes* and *C. violaceus* illustrations were drawn by Andrea Dingeldein.

**Supplemental Figure S10: Circos plot showing synteny between the assembled genome** *Cebidichthys violaceus* **and** *Danio rerio* **(Zebrafish).** There are 25 chromosomes (yellow boxes) which represent the zebrafish genome. There are 114 blue boxes which are 1 MB or greater that represent the *C. violaceus* genome. There are 353 contigs that are less than 1 MB which were concatenated into box labeled as 115. Gray strands indicate syntenic regions between the two genomes. Both *D. rerio* and *C. violaceus* illustrations were drawn by Andrea Dingeldein.

**Supplemental Figure S11:  Circos plot showing synteny between the assembled genome**
*Cebidichthys violaceus* **and** *Lepisosteus oculatus* **(Spotted gar).** There are 29 linkage groups
(orange boxes) which represent the Spotted gar genome. There are 114 blue boxes which are 1
MB or greater that represent the *C. violaceus* genome. There are 353 contigs that are less than 1
MB which were concatenated into box labeled as 115. Gray strands indicate syntenic regions
between the two genomes. *Lepisosteus oculatus* photo was taken by David Solomon and the *C.*
*violaceus* illustration was drawn by Andrea Dingeldein.

**Supplemental Figure S12: Flowchart of our genome guided transcriptome assembly from nine tissues using Illumina (two lanes of PE 100bp).**
Light blue boxes indicate raw sequence reads for nine tissues: liver, heart, gill, pyloric caeca (PC), proximal intestine (PI), middle intestine (MI), spleen, gonad (testes), and brain.
Dark blue boxes indicate the bioinformatic program used in the pipeline for trimming/cleaning reads, normalizing reads, assembling transcripts with our assemble genome as a reference, and estimating differential gene expression (DEGs) and analysis of DEGs. (…………) indicates the type of sequence information that was used for the start of the assembly method. Arrows indicate the next step taken to proceed in the transcriptome assembly and analysis.

**Supplemental Figure S13: Heatmap showing Differentially Expressed Genes for nine tissues of *Cebidichthys violaceus*.** Heatmap was generated with the csHeatmap feature in cummerbund, where dark blue represents a high FPKM value and white indicates a low FPKM value. There were 15,490 differentially expressed genes across all nine tissue types.

**Supplemental Figure S14: Boxplots visualized for all nine tissue types displaying summary statistics of Fragments Per Kilobase of transcript per Million mapped reads (FPKM).** Boxplots were generated with the csBoxplot function in cummerbund for nine tissues.

**Supplemental Figure S15: Significant features of genes for all nine tissue types. Significant features (alpha value 0.01) of genes between tissue types were estimated by using the sigmatrix function in CUMMERBUND.** The darker green shades indicate a higher significant features of genes identified, whereas a lighter green/white shade indicates a less significant features of genes identified.

**Supplemental Figure S16: Dendrogram of all nine tissue transcriptomes to determine relationships of each tissue type.** A dendrogram was constructed of all nine tissue types by using Jensen-Shannon (JS) distances as shown.

**Supplemental Figure S17: Distance matrix of all nine tissue transcriptomes based on Jensen-Shannon (JS).** Plot was constructed with csDistHeat function in CUMMERBUND. Dark red indicates an increased JS distance between the pairwise comparison. Lighter red/white indicates less distance between the pairwise comparison.

**Supplemental Figure 18: Phylogenetic relationship of chymotrypsin in stichaeids, gene copy number, and molecular evolution of chymotrypsin.**

**a**, Synteny map for chymotrypsin genes from *Danio rerio*, *Oryzias latipes*, *Gasterosteus aculeatus,* and *C. violaceus*. **c-d**, An adaptive branch-site Random Effects Likelihood (aBSREL) test for episodic diversification phylogenetic tree constructed for chymotrypsin genes from *C. violaceus* (*chymo A* and *ctrb*) and other intertidal stichaeid species. Branches thicker than the other branches have a P<0.05 (corrected for multiple comparisons) to reject the null hypothesis of all ω on that branch (neutral or negative selection only). A thick branch is considered to have experienced diversifying positive selection. **e-f**, The output of Mixed Effects Model of Evolution (MEME) to detect episodic positive/diversifying selection at sites. $\beta_+$ is the non-synonymous substitution rate at a site for the positive/neutral evolution throughout the sequence of the chymotrypsin A gene. **e**, MEME output for chymotrypsin B gene. ** is an indication that the positive/diversifying site is statistically significant with a p-value < 0.01.

**Supplemental Figure 19: Phylogenetic relationship of chymotrypsin-like gene in stichaeids, gene copy number, and molecular evolution of chymotrypsin-like (*ctrl*).**

**a**, A maximum likelihood tree was generated for chymotrypsin-like genes from *Cebidichthys violaceus* and other intertidal stichaeid species: *Anoplarchus purpurescens*, *Phytichthys chirus, Xiphister mucosus*, and *Xiphister atropurureus*. *Gasterosteus aculeatus* chymotrypsin gene was used as an outgroup for our phylogenetic analysis. **b**, Synteny map for amylase genes from *Danio rerio*, *Oryzias latipes*, *Gasterosteus aculeatus,* and *C. violaceus*. **c**, An adaptive branch-site Random Effects Likelihood (aBSREL) test for episodic diversification phylogenetic tree constructed for the chymotrypsin gene from *C. violaceus* (*ctrl*) and other intertidal stichaeid species. Branches thicker than the other branches have a P<0.05 (corrected for multiple comparisons) to reject the null hypothesis of all ω on that branch (neutral or negative selection only). A thick branch is considered to have experienced diversifying positive selection. **d**, The output of Mixed Effects Model of Evolution (MEME) to detect episodic positive/diversifying selection at sites. β+ is the non-synonymous substitution rate at a site for the positive/neutral evolution throughout the sequence of the gene. There were zero sites under positive/diversifying selection.

**Supplemental Figure 20: Phylogenetic relationship of trypsin in stichaeids, gene copy number, and molecular evolution of trypsin (*try3*).**

**a-b**, Synteny map for trypsin genes from *Danio rerio*, *Oryzias latipes*, *Gasterosteus aculeatus,* and *C. violaceus*. Clear boxes indicate multiple loci present in this region. **c-d**, An adaptive branch-site Random Effects Likelihood (aBSREL) test for episodic diversification phylogenetic tree constructed for trypsin genes from *C. violaceus* (*try3*) and other intertidal stichaeid species. Branches thicker than the other branches have a P<0.05 (corrected for multiple comparisons) to reject the null hypothesis of all ω on that branch (neutral or negative selection only). A thick branch is considered to have experienced diversifying positive selection. **e-f**, The output of Mixed Effects Model of Evolution (MEME) to detect episodic positive/diversifying selection at sites. β+ is the non-synonymous substitution rate at a site for the positive/neutral evolution throughout the sequence of the gene. * is an indication that the positive/diversifying site is statistically significant with a p-value < 0.05.
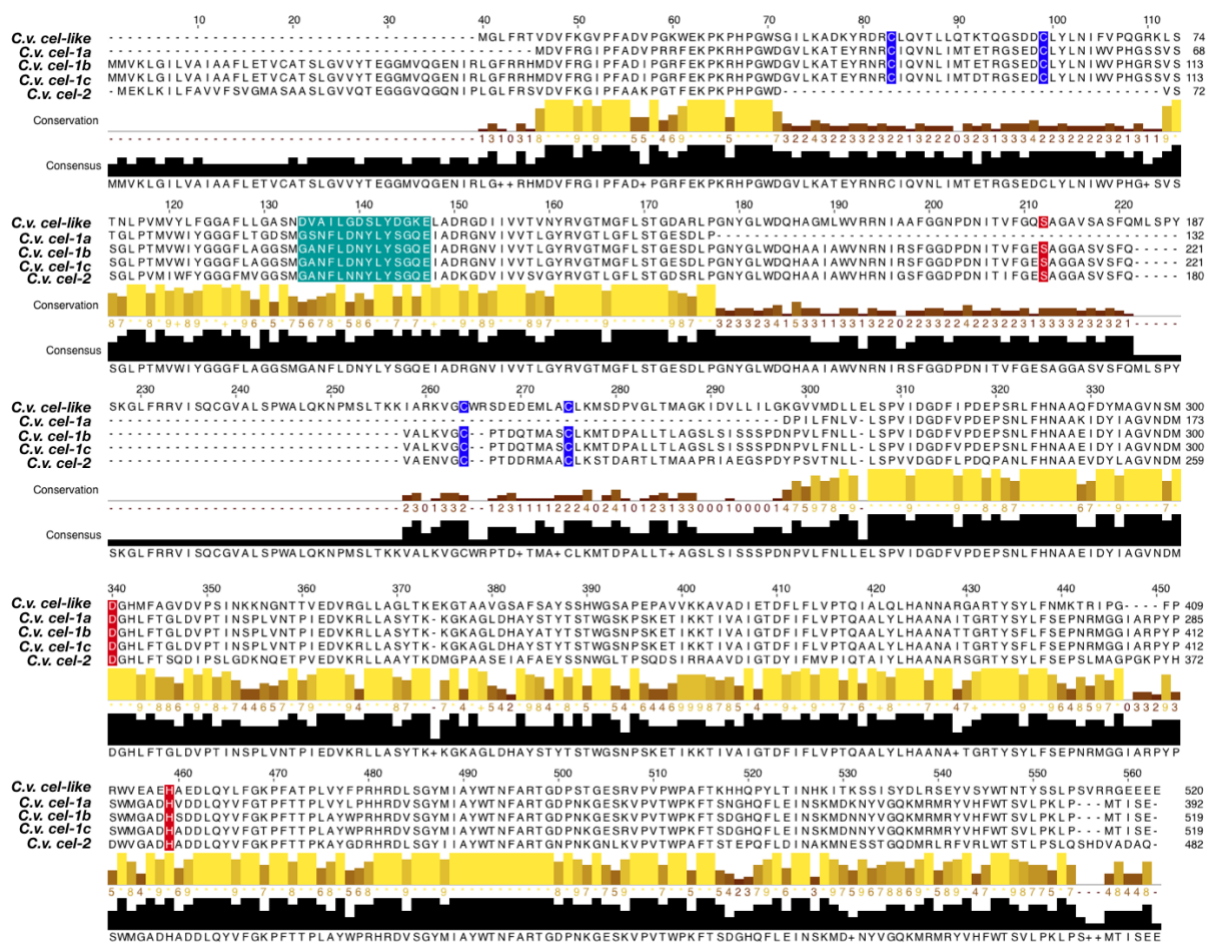
**Supplemental Figure 21: Gene copy number and molecular evolution of alanyl aminopeptidase a (*anpepa*). a**, Synteny map for *anpepa* genes from *Danio rerio*, *Oryzias latipes*, *Gasterosteus aculeatus*, and *C. violaceus*. **b**, An adaptive branch-site Random Effects Likelihood (aBSREL) test for episodic diversification phylogenetic tree constructed for the *anpepa* gene from *C. violaceus* and other intertidal stichaeid species. ω is the ratio of nonsynonymous to synonymous substitutions. The color gradient represents the magnitude of the corresponding ω. Branches thicker than the other branches have a P<0.05 (corrected for multiple comparisons) to reject the null hypothesis of all ω on that branch (neutral or negative selection only). A thick branch is considered to have experienced diversifying positive selection. **c**, The output of Mixed Effects Model of Evolution (MEME) to detect episodic positive/diversifying selection at sites. $\beta_+$ is the non-synonymous substitution rate at a site for the positive/neutral evolution throughout the sequence of the gene. ** is an indication that the positive/diversifying site is statistically significant with a p-value < 0.01 and * is for p-value < 0.05.

**Supplemental Figure 22: Enzyme activity and gene copy number of alanyl aminopeptidase (*anpepb*, *anpep Ey*, *anpep N*, and *anpep Ey-like*). a**, Total gut standardized aminopeptidase activity for *Cebidichthys violaceus* (Cv) and other intertidal stichaeid species: *Phytichthys chirus* (Pc), *Xiphister mucosus* (Xm), *Xiphister atropurpureus* (Xa), and *Anoplarchus purpurescens* (Ap). H = herbivory, O = Omnivory, and C = Carnivory. Values are mean ± standard deviation with n = 6 for Cv, Xm, Xa, and Ap; and n = 9 for Pc (German et al. 2015). Interspecific comparisons were made for aminopeptidase with ANOVA, where circles that share a letter are not significantly different. **b-c**, Synteny maps for aminopeptidase genes from *Danio rerio*, *Oryzias latipes*, *Gasterosteus aculeatus,* and *C. violaceus*.

# Aminopeptidase B, Ey, N, and Ey-like



anpepb

**a**



anpep Ey

**b**



anpep N

**c**



anpep Ey-like

**d**

e — anpep b

f — anpep Ey





g — anpep N

h — anpep Ey-like

**Supplemental Figure 23: Estimation of selection analyses of alanyl aminopeptidase (*anpepb*, *anpep Ey*, *anpep N*, and *anpep Ey-like*) genes in stichaeids. a-d**, An adaptive branch-site Random Effects Likelihood (aBSREL) test for episodic diversification phylogenetic tree constructed for *anpep* genes from *C. violaceus* and other intertidal stichaeid species. Branches thicker than the other branches have a $P<0.05$ (corrected for multiple comparisons) to reject the null hypothesis of all $\omega$ on that branch (neutral or negative selection only). A thick branch is considered to have experienced diversifying positive selection. **e-h**, The output of Mixed Effects Model of Evolution (MEME) to detect episodic positive/diversifying selection at sites. ** is an indication that the positive/diversifying site is statistically significant with a p-value < 0.01 and * is for p-value < 0.05.

**Supplemental Figure 24: Phylogenetic relationship of alanyl aminopeptidase genes in fishes (including *Cebidichthys violaceus*).** A maximum likelihood (ML) tree was constructed with 1,000 bootstrap replicates in PhyML v3.1 based alanyl aminopeptidase sequences from *C. violaceus* and other fish species (e.g. *Danio rerio*, *Gasterosteus aculeatus*, *Oryzias latipes*, and *Petromyzon marinus*). Alanyl aminopeptidase sequences from *P. marinus* were used as an outgroup.

**Supplemental Figure S25: Pacific Biosciences (PacBio) reads mapped to *Cebidichthys violaceus* genome assembly on contig 440.** PacBio reads with read ID numbers labeled which span regions of the three amylase loci (two *amy2a* loci and the *amy2b*) on contig 440. AUGUSTUS gene predictions are used to reference where *amy2* loci are located on the *C. violaceus* genome.

**Supplemental Figure S26: Repetitive elements identified adjacent to amylase loci.**
Light green bars are repetitive elements identified on contig 440 and blue bars represent exons of amylase loci (AUGUSTUS gene prediction).

# Phospholipase B1



a

b

c

*plb1-1*

d

*plb1-2*

e

*plb1-3*

f

*plb1-1*

g

*plb1-2*

h

*plb1-3*

**Supplemental Figure 27: Gene copy number, and molecular evolution of Phospholipase B1 (*plb*).**
**a-b**, Synteny map for phospholipase B1 (*plb1*) genes from *Danio rerio*, *Oryzias latipes*, *Gasterosteus aculeatus,* and *C. violaceus*. **c-e**, An adaptive branch-site Random Effects Likelihood (aBSREL) test for episodic diversification phylogenetic tree constructed for *plb1* genes from *C. violaceus* and other intertidal stichaeid species. Branches thicker than the other branches have a P<0.05 (corrected for multiple comparisons) to reject the null hypothesis of all ω on that branch (neutral or negative selection only). A thick branch is considered to have experienced diversifying positive selection. **f-h,** The output of Mixed Effects Model of Evolution (MEME) to detect episodic positive/diversifying selection at sites. $\beta_+$ is the non-synonymous substitution rate at a site for the positive/neutral evolution throughout the sequence of the gene. ** is an indication that the positive/diversifying site is statistically significant with a p-value < 0.01 and * is for p-value < 0.05.

**Supplemental Figure 28: Phylogenetic relationships, gene copy number, and molecular evolution of secretory phospholipase Group 12B (*pg12b*) a-b**, Synteny map for genes from *Danio rerio*, *Oryzias latipes*, *Gasterosteus aculeatus*, and *C. violaceus*. **c-d**, An adaptive branch-site Random Effects Likelihood (aBSREL) test for episodic diversification phylogenetic tree constructed for *pg12b* genes from *C. violaceus* and other intertidal stichaeid species. Branches thicker than the other branches have a P<0.05 (corrected for multiple comparisons) to reject the null hypothesis of all ω on that branch (neutral or negative selection only). A thick branch is considered to have experienced diversifying positive selection. **e-f**, The output of Mixed Effects Model of Evolution (MEME) to detect episodic positive/diversifying selection at sites. β+ is the non-synonymous substitution rate at a site for the positive/neutral evolution throughout the sequence of the gene. ** is an indication that the positive/diversifying site is statistically significant with a p-value < 0.01.

**Supplemental Figure S29: Repetitive elements identified adjacent to carboxyl ester lipase (*cel*) loci.**
Light green bars are repetitive elements identified on contig 445 and blue bars represent exons of *cel* loci (AUGUSTUS gene prediction).

**Carboxyl Ester Lipase-like**

**Supplemental Figure 30: Molecular analyses of carboxyl ester lipase-like genes in stichaeids a,** A maximum likelihood (ML) phylogenetic tree of carboxyl ester lipase genes using *Danio rerio*, *Oryzias latipes*, *Gasterosteus aculeatus, Ctenopharyngodon idella, Eptatretus burgeri*, and *C. violaceus.* This tree was constructed with 1,000 bootstrap replicates and a TIM2+I+G model. **b,** gene copy number of *cel-like* loci. **c,** estimation of lineage-specific selection of stichaeids *cel-like* loci using aBSREL. Branches thicker than the other branches have a P<0.05 (corrected for multiple comparisons) to reject the null hypothesis of all ω on that branch (neutral or negative selection only). **d,** estimation of site-level episodic selection in stichaeids *cel-like* loci using MEME**.** There were zero sites under positive/diversifying selection.

**Supplemental Figure S31: Alignment of carboxyl ester lipase (*cel*) amino acid sequences.**
*Cebidichthys violaceus cel* sequences viewed in JALVIEW v2.10.5 (http://www.jalview.org/).
Amino acids highlighted in blue indicates disulfide bonds, red indicates active sites present, and
teal indicates the bile  salt-binding site.

**Opsin Short Wave Sensitive**

**Supplemental Figure S32: Gene copy number and molecular evolution of Opsin Short Wave Sensitive (*opn1sw*) genes. a**, Synteny map for *opn1sw* genes from *Danio rerio*, *Oreochromis niloticus*, *Gasterosteus aculeatus*, and *Cebidichthys violaceus*. *D. rerio*, *G. aculeatus*, and *C. violaceus* were drawn by Andrea Dingeldein. *O. niloticus* illustration was drawn by Milton Tan. **b**, An adaptive Branch-Site Random Effects Likelihood (aBSREL) test for episodic diversification was estimated and represented as a phylogenetic tree for *opn1sw* genes from *C. violaceus* and three other fishes. Branches thicker than the other branches have a P<0.05 (corrected for multiple testing) to reject the null hypothesis of all ω on that branch (neutral or negative selection only). A thick branch is considered to have experienced diversifying positive selection. c, The output of Mixed Efffects Model of Evolution (MEME) to detect episodic positive/diversifying selection at sites. β+ is the non-synonymous substitution rate at a site for the positive/neutral evolution throughout the sequence of the gene. ** is an indication that the positive/diversifying site is statistically significant with a p-value < 0.01 and * is for p-value < 0.05.

Supplemental Table S1. Fish genomes available online, including all relevant data about each submission.  Natural diets indicated by color with red being carnivore, purple omnivores, and green herbivore.

| Organism | Name | Link | Submitter | Date | Genome representation | Assembly level | Version status | RefSeq category | Total sequence length | Total assembly gap length | Gaps between scaffolds | Number of scaffolds | Scaffold N50 | Scaffold L50 | Number of contig | Contig N50 | Contig L50 | Total number of chromosomes and plasmids | Number of component sequences (WGS or clone) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

Carnivore
Omnivore
Herbivore

**Supplemental Table S2: Genome Sequencing Information**

| | |
|---|---|
| **Pacific Biosciences** | |
| Number of SMRT Cells | 40 |
| Amount Polymerase in picomolar (pM)* | 150 (17); 200 (1); 300 (13); 400 (9) |
| Total Number of Reads | 2,421,941 |
| Average N50 (bp) | 17,102.78 |
| Total data (Mb) | 29,700 |
| | |
| **Illumina (100 Paired End Sequencing)** | |
| Number of Lanes | 2 |
| Number of Reads from the 1st Lane (Both Reads 1 and 2) | 422,313,916 |
| Number of Reads from the 2nd Lane (Both Reads 1 and 2) | 423,075,242 |
| Total data (Mb) | 84,539 |
| * values within parentheses indicate the amount of SMRT cells used | |

**Supplemental Table S3: BUSCO v3 Estimation on the *Cebidichthys violaceus* genome**

|  | BUSCO V3 |
|---|---|
| Complete BUSCOs | 2508 |
| Complete BUSCOs and single-copy BUSCOs | 2474 |
| Complete BUSCOs and duplicated BUSCOs | 34 |
| Fragmented BUSCOs | 28 |
| Missing BUSCOs | 50 |
| Total BUSCO groups searched | 2586 |
|  |  |
|  |  |
|  |  |

**Supplemental Table S4: RepeatMasker for the *C. violaceus* assembled genome**

| | Number of Elements | Length Occupied (bp) | Percentage of Sequence (%) | | |
|---|---|---|---|---|---|
| Retroelements | 11715 | 5340689 | 0.9 | | |
| SINEs: | 1985 | 223655 | 0.04 | | |
| Penelope | 60 | 24362 | 0 | | |
| LINEs: | 8905 | 4191047 | 0.71 | | |
| CRE/SLACS | 0 | 0 | 0 | | |
| L2/CR1/Rex | 4980 | 2049678 | 0.35 | | |
| R1/LOA/Jockey | 0 | 0 | 0 | | |
| R2/R4/NeSL | 85 | 24828 | 0 | | |
| RTE/Bov-B | 3316 | 1768649 | 0.3 | | |
| L1/CIN4 | 336 | 205561 | 0.03 | | |
| LTR elements: | 825 | 925987 | 0.16 | | |
| BEL/Pao | 20 | 36895 | 0.01 | | |
| Ty1/Copia | 22 | 20106 | 0 | | |
| Gypsy/DIRS1 | 703 | 846315 | 0.14 | | |
| Retroviral | 79 | 22637 | 0 | | |
| | | | | | |
| DNA transposons | 10562 | 2582768 | 0.44 | | |
| hobo-Activator | 4384 | 730619 | 0.12 | | |
| Tc1-IS630-Pogo | 4812 | 1690957 | 0.29 | | |
| En-Spm | 0 | 0 | 0 | | |
| MuDR-IS905 | 0 | 0 | 0 | | |
| PiggyBac | 352 | 30854 | 0.01 | | |
| Tourist/Harbinger | 178 | 26872 | 0 | | |
| Other (Mirage, P-element, Transib) | 0 | 0 | 0 | | |
| | | | | | |
| Rolling-circles | 0 | 0 | 0 | | |
| | | | | | |
| Unclassified: | 203 | 21708 | 0 | | |
| | | | | | |
| Total interspersed repeats: | | 7945165 | 1.34 | | |
| | | | | | |
| | | | | | |
| Small RNA: | 2430 | 215618 | 0.04 | | |
| | | | | | |
| Satellites: | 4 | 1112 | 0 | | |
| Simple repeats: | 468887 | 27486227 | 4.64 | | |
| Low complexity: | 37947 | 2434577 | 0.41 | | |

| Ortholog cluster ID | Annotation | Alignment Length | | |
|---|---|---|---|---|
| | **Supplemental Table S5: Annotation of the 30 loci used for the phylogeny in figure 1** | | | |
| 5 | amyloid beta precursor protein (cytoplasmic tail) binding protein 2 (APPBP2) | 1767 | | |
| 12 | anaphase promoting complex subunit 7 (ANAPC7) | 1806 | | |
| 25 | aminoadipate-semialdehyde dehydrogenase (AASDH) | 3860 | | |
| 39 | tyrosyl-tRNA synthetase (YARS) | 1634 | | |
| 123 | PRP18 pre-mRNA processing factor 18 homolog (prpf18) | 1083 | | |
| 131 | sphingomyelin phosphodiesterase 2, neutral membrane (neutral sphingomyelinase) (SMPD2) | 1456 | | |
| 141 | UbiA prenyltransferase domain containing 1 (UBIAD1) | 1086 | | |
| 163 | dual oxidase 1-like | 5184 | | |
| 187 | phosphatidylinositol glycan anchor biosynthesis class M (pigm) | 1296 | | |
| 191 | methyltransferase like 9 (METTL9) | 1194 | | |
| 198 | chromosome 22 C6orf62 homolog (c22h6orf62) | 701 | | |
| 214 | uncharacterized protein F13E9.13, mitochondrial-like | 875 | | |
| 218 | mediator complex subunit 7 (MED7) | 831 | | |
| 220 | transmembrane protein 98 (TMEM98) | 687 | | |
| 223 | prolactin regulatory element binding (preb) | 1520 | | |
| 230 | Shwachman-Bodian-Diamond syndrome (SBDS) | 783 | | |
| 235 | ATP synthase, H+ transporting, mitochondrial F1 complex, O subunit (ATP5O) | 639 | | |
| 236 | sodium channel modifier 1 (SCNM1) | 906 | | |
| 241 | cysteine dioxygenase type 1 (CDO1) | 787 | | |
| 242 | GINS complex subunit 2 (Psf2 homolog) (GINS2) | 688 | | |
| 252 | mediator complex subunit 22 (MED22) | 686 | | |
| 269 | HD domain containing 3 (HDDC3) | 550 | | |
| 271 | mago homolog, exon junction complex subunit (MAGOH) | 447 | | |
| 281 | optic atrophy 3 (autosomal recessive, with chorea and spastic paraplegia) (OPA3) | 495 | | |
| 288 | splicing factor 3b subunit 6 (sf3b6) | 378 | | |
| 290 | polymerase (RNA) II (DNA directed) polypeptide F (POLR2F) | 503 | | |
| 301 | mitochondrial ribosomal protein S17 (mrps17) | 489 | | |
| 311 | C1D nuclear receptor corepressor (c1d) | 484 | | |
| 315 | replication protein A3, 14kDa (RPA3) | 369 | | |
| 320 | NADH:ubiquinone oxidoreductase subunit S5 (ndufs5) | 324 | | |

**Supplemental Table S6: Pairwise Comparison of Orthologs of *Cebidichthys violaceus* and fish genomes deposited on Ensembl**

| | Latimeria chalumnae | Lepisosteus oculatus | Danio rerio | Astyanax mexicanus | Gadus morhua | Takifugu rubripes | Tetraodon nigroviridis | Gasterosteus aculeatus | Cebidichthys violaceus | Xiphophorus maculatus | Poecilia formosa | Oryzias latipes | Oreochromis niloticus |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Latimeria chalumnae** | | | | | | | | | | | | | |
| **Lepisosteus oculatus** | 12563 | | | | | | | | | | | | |
| **Danio rerio** | 12374 | 7237 | | | | | | | | | | | |
| **Astyanax mexicanus** | 11795 | 13247 | 16214 | | | | | | | | | | |
| **Gadus morhua** | 9052 | 110 | 1594 | 12385 | | | | | | | | | |
| **Takifugu rubripes** | 11346 | 6281 | 14060 | 2246 | 13671 | | | | | | | | |
| **Tetraodon nigroviridis** | 10745 | 11808 | 13207 | 1399 | 12851 | 15335 | | | | | | | |
| **Gasterosteus aculeatus** | 4717 | 12652 | 12624 | 13279 | 6725 | 15429 | 2245 | | | | | | |
| **Cebidichthys violaceus** | 9472 | 10483 | 5357 | 11000 | 7886 | 11805 | 10965 | 12022 | | | | | |
| **Xiphophorus maculatus** | 11730 | 13157 | 14824 | 13864 | 14081 | 2059 | 14318 | 15867 | 12249 | | | | |
| **Poecilia formosa** | 12271 | 13926 | 4201 | 14610 | 14478 | 6240 | 14865 | 8343 | 13241 | 18391 | | | |
| **Oryzias latipes** | 10795 | 11868 | 3796 | 12418 | 13040 | 14465 | 5278 | 14674 | 11024 | 14610 | 15367 | | |
| **Oreochromis niloticus** | 11858 | 4274 | 7487 | 7576 | 14102 | 15732 | 14502 | 15897 | 3913 | 5828 | 17602 | 14761 | |

| Supplemental Table S7: Genome Size Estimation with Jellyfish v2.2.0 (Marçais and | | |
|---|---|---|
| **Kmer Size** | Peak | Genome Estimate Size |
| **31** | 45 | 657,524,901 |
| **29** | 46 | 661,608,906 |
| **27** | 48 | 651,654,143 |
| **25** | 49 | 655,607,918 |
| | | |
| **Average Genome** | | **656,598,967** |
| **Standard Deviation** | | 4,138,853 |

**Supplemental Table S8: Transcriptomic Sequencing Information and Trinity Assembly**

| | TUBE ID | RIN Values | Barcode ID | Number of Reads 1st lane | Number of Reads 2nd lane | Trimmomatic | Normalized | Mapped (Overall) | Aligned Pairs | Transcripts Present |
|---|---|---|---|---|---|---|---|---|---|---|
| **Liver** | CV100 | 8.6 | ACAGTG | 7,954,337 | 8,119,254 | 14,705,062 | 1,330,380 | 71.5% | 849,686 | 20,008 |
| **Brain** | CV98 | 9.4 | GCCAAT | 12,832,744 | 13,197,868 | 25,319,484 | 7,389,717 | 83.7% | 5,711,530 | 60,430 |
| **Heart** | CV97,98,99 | 8.5 | CAGATC | 17,371,046 | 17,402,450 | 33,817,846 | 3,961,155 | 84.8% | 3,129,930 | 35,570 |
| **Gill** | CV96 | 8.4 | CGATGT | 12,365,345 | 11,647,444 | 22,903,000 | 4,518,995 | 80.5% | 3,354,450 | 78,629 |
| **Pyloric Caeca** | CV96 | 8.8 | CTTGTA | 16,098,871 | 18,297,873 | 33,400,321 | 4,396,658 | 83.3% | 3,400,229 | 40,201 |
| **Proximal Intestine** | CV97,98 | 9.7 | TGACCA | 10,420,112 | 9,421,382 | 15,949,732 | 2,621,189 | 69.7% | 1,602,595 | 37,277 |
| **Middle Intestine** | CV97,98 | 8.6 | AGTTCC | 16,528,911 | 16,500,293 | 31,957,043 | 4,805,496 | 83.8% | 3,740,247 | 41,978 |
| **Spleen** | CV97,98,99,100 | 8.2 | ATGTCA | 31,903,535 | 39,356,506 | 69,076,597 | 6,354,519 | 81.5% | 4,796,120 | 48,270 |
| **Gonad (Testes)** | CV99,100 | 9.7 | AGTCAA | 14,000,194 | 16,098,169 | 29,305,941 | 9,061,164 | 67.4% | 5,640,612 | 60,487 |

**Supplemental Table S9: RepeatMasker for the Genome Guided (TRINITY) Transcriptome**

| | Number of Elements | Length Occupied (bp) | Percentage of Sequence (%) | |
|---|---|---|---|---|
| Retroelements | 345 | 30165 | 0.03% | |
| SINEs | 18 | 1278 | 0.00% | |
| Penelope | 2 | 92 | 0.00% | |
| LINEs | 164 | 14962 | 0.01% | |
| L2/CR1/Rex | 104 | 9935 | 0.01% | |
| R1/LOA/Jockey | 8 | 789 | 0.00% | |
| R2/R4/NeSL | 4 | 394 | 0.00% | |
| RTE/Bov-B | 7 | 613 | 0.00% | |
| L1/CIN4 | 32 | 2704 | 0.00% | |
| LTR elements | 163 | 13925 | 0.01% | |
| BEL/Pao | 14 | 733 | 0.00% | |
| Ty1/Copia | 0 | 0 | 0.00% | |
| Gypsy/DIRS1 | 110 | 9224 | 0.01% | |
| Retroviral | 25 | 2683 | 0.00% | |
| | | | | |
| DNA transposons | 513 | 36648 | 0.03% | |
| hobo-Activator | 180 | 12759 | 0.01% | |
| Tc1-IS630-Pogo | 27 | 1952 | 0.00% | |
| PiggyBac | 4 | 421 | 0.00% | |
| Tourist/Harbinger | 23 | 2569 | 0.00% | |
| Other (Mirage, P-element, Transib) | 0 | 0 | 0.00% | |
| Total interspersed repeats | | 75248 | 0.07% | |
| | | | | |
| Small RNA | 3 | 264 | 0.00% | |
| Satellites | 13 | 948 | 0.00% | |
| Simple repeats | 69 | 6138 | 0.01% | |
| Low complexity | 10 | 1450 | 0.00% | |

| Supplemental Table S10: RepeatMasker for AUGUSTUS Predicted Genes | | | |
|---|---|---|---|
| | Number of Elements | Length Occupied (bp) | Percentage of Sequence (%) |
| Retroelements | 451 | 82999 | 0.19% |
| SINEs | 13 | 827 | 0.00% |
| Penelope | 9 | 2950 | 0.01% |
| LINEs | 211 | 36107 | 0.08% |
| L2/CR1/Rex | 128 | 15756 | 0.04% |
| R1/LOA/Jockey | 10 | 813 | 0% |
| R2/R4/NeSL | 5 | 1512 | 0% |
| RTE/Bov-B | 12 | 1606 | 0% |
| L1/CIN4 | 22 | 7877 | 0.02% |
| LTR elements | 227 | 46065 | 0.1% |
| BEL/Pao | 11 | 1225 | 0% |
| Ty1/Copia | 19 | 4233 | 0.01% |
| Gypsy/DIRS1 | 102 | 18320 | 0.04% |
| Retroviral | 80 | 19852 | 0.04% |
| | | | |
| DNA transposons | 877 | 123135 | 0.28% |
| hobo-Activator | 431 | 65743 | 0.15% |
| Tc1-IS630-Pogo | 84 | 12063 | 0.03% |
| PiggyBac | 36 | 6698 | 0.02% |
| Tourist/Harbinger | 31 | 4641 | 0.01% |
| Other (Mirage, P-element, Transib) | 166 | 13032 | 0.03% |
| Total interspersed repeats | | 240853 | 0.55% |
| | | | |
| Small RNA | 4 | 250 | 0% |
| Satellites | 33 | 16963 | 0.04% |
| Simple repeats | 10737 | 482867 | 1.09% |
| Low complexity | 3249 | 206240 | 0.47% |

| Supplemental Table S11: Estimation of Total Genes from Cuffmerge and Augustus | | | | |
|---|---|---|---|---|
| **Total Estimated Transcripts from Cuffmerge** | 101,922 | | | |
| **Trinotate Annotation** | 65,535 | | | |
| **Top BLASTX hit** | 26,356 | | | |
| **Total Estimates from Augustus (de novo)** | 29,525 | | | |
| **80% Hit Coverage from Cuffmerge Assembly Uniprot** | 5,199 | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |

**Supplemental Table S12: Estimation of Full Length of Transcripts from all nine transcriptomes**

| Hit percent coverage bin | Count in bin | >Bin below |
|---|---|---|
| 100 | 2692 | 2692 |
| 90 | 1258 | 3950 |
| 80 | 1249 | 5199 |
| 70 | 1358 | 6557 |
| 60 | 1598 | 8155 |
| 50 | 1991 | 10146 |
| 40 | 2411 | 12557 |
| 30 | 2912 | 15469 |
| 20 | 2966 | 18435 |
| 10 | 1132 | 19567 |

**Supplemental Table S13: Candidate Genes (Digestive Enzymes) Selection Analyses**

| Gene Name | Contig Location | Alignment Length (Nucleotide) | Alignment Length (AA) | % of Swissprot Annotation Hit | Gard (sites) | Absrel | Meme (sites) |
|---|---|---|---|---|---|---|---|
| Aminopeptidase A (*anpepa*) | 86 | 2127 | 709 | 73.40% | 111, 1039 | No evidence | 194*,412**,445**,593** |
| Aminopeptidase B (*anpepb*) | 3 | 1497 | 499 | 51.34% | 801 | No evidence | 156* |
| Aminopeptidase Ey (*anpep Ey*) | 3 | 1500 | 500 | 51.71% | No recombination | No evidence | No Sites |
| Aminopeptidase Ey-like (*anpep Ey-like*) | 78 | 927 | 358 | 37.02% | 309 | 1 branch | 38* |
| Aminopeptidase N (*anpep N*) | 78 | 1449 | 483 | 50.00% | 616 | No evidence | 351** |
| Amylase (*amy2a* and *amy2b*) | 440 | 1536 | 512 | 100.79% | No recombination | 1 branch | 41*, 256*, 279* |
| Carboxyl Ester Lipase 1 (*cel-1a,b,c*) | 445 | 1530 | 510 | 85.43% | No recombination | No evidence | 64**, 258**, 355** |
| Carboxyl Ester Lipase 2 (*cel-2*) | 445 | 1053 | 351 | 58.79% | 241, 303, 975 | No evidence | No sites |
| Carboxyl Ester Lipase-like (*cel-like*) | 138 | 1560 | 520 | 87.10% | No recombination | No evidence | No sites |
| Chymotrypsin A (*ctra-1 & ctra-2*) | 55 | 789 | 263 | 100.00% | 99 | No evidence | No sites |
| Chymotrypsin B (*ctrb*) | 55 | 792 | 264 | 107.76% | 81, 336, 581 | 2 branches | 112** |
| Chymotrypsin-like (*ctrl*) | 427 | 789 | 263 | 99.62% | No recombination | No evidence | No sites |
| Phospholipase B1 (*plb1-1*) | 442 | 2697 | 899 | 60.99% | 202 | No evidence | 66**, 97*, 289*, 438**, 800*, 821* |
| Phospholipase B1 (*plb1-2*) | 434 | 2004 | 668 | 45.32% | 532 | 1 branch | 183**, 230*,476* |
| Phospholipase B1 (*plb1-3*) | 356 | 582 | 194 | 13.16% | 95, 124, 282 | No evidence | No sites |
| Phospholipase B12 (*pg12b-1*) | 428 | 606 | 202 | 103.59% | No recombination | No evidence | No sites |
| Phospholipase B12 (*pg12b-2*) | 413 | 582 | 194 | 99.49% | 95, 124, 282 | No evidence | 25**,31**,32**,33**,37**,83**,142**,146**,185** |
| Trypsin-3_1 (*try3-1*) | 427 | 750 | 250 | 105.04% | 21 | No evidence | No sites |
| Trypsin-3_2 (*try3-2*) | 435 | 732 | 244 | 102.52% | 141 | No evidence | 91* |

** is an indication that the positive/diversifying site is statistically significant with a p-value < 0.01 and * is for p-value < 0.05.

| Supplementary Table S14: Ensembl IDs used for *anpep* and *cel* | | | | |
|---|---|---|---|---|
| **Species Name** | **Gene Name** | **Ensembl ID** | | |
| *Danio rerio* | *anpepb* | ENSDARG00000103878 | | |
| *Oryzias latipes* | *anpepb* | ENSORLG00020014580 | | |
| *Gasterosteus aculeatus* | *anpepb* | ENSGACG00000014140 | | |
| *Danio rerio* | *si:ch211* | ENSDARG00000097285 | | |
| *Oryzias latipes* | *anpep Ey* | ENSORLG00020014549 | | |
| *Danio rerio* | *anpep* | ENSDARG00000089706 | | |
| *Oryzias latipes* | *anpep N* | ENSORLG00000014691 | | |
| *Gasterosteus aculeatus* | *anpep-201* | ENSGACG00000014748 | | |
| *Oryzias latipes* | *anpep Ey-like* | ENSORLG00000029229 | | |
| *Gasterosteus aculeatus* | *anpep-202* | ENSGACG00000014748 | | |
| *Danio rerio* | *anpepa* | ENSDARG00000041083 | | |
| *Oryzias latipes* | *anpepa* | ENSORLG00000019272 | | |
| *Gasterosteus aculeatus* | *anpepa* | ENSGACG00000002363 | | |
| *Petrus marinus* | *anpep* | ENSPMAG00000003227 | | |
| *Petrus marinus* | *anpep* | ENSPMAG00000009142 | | |
| *Petrus marinus* | *anpep* | ENSPMAG00000009172 | | |
| | | | | |
| *Danio rerio* | *cel.2-201* | ENSDARG00000029822 | | |
| *Danio rerio* | *cel.1-202* | ENSDARG00000017490 | | |
| *Oryzias latipes* | *cel-1a* | ENSORLG00000014439 | | |
| *Oryzias latipes* | *cel-1b* | ENSORLG00000014464 | | |
| *Gasterosteus aculeatus* | *cel-2* | ENSGACG00000018127 | | |
| *Gasterosteus aculeatus* | *cel-1* | ENSGACG00000018130 | | |
| *Oryzias latipes* | *cel-like* | ENSORLG00000016428 | | |
| *Eptatretus burgeri* | *cel* | ENSEBUG00000006718 | | |
| *Ctenopharyngodon* | *cel* | CI01000006 (scaffold) | | |
| *C. idella cel* gene was taken from the Grass Carp Genome Database | | | | |
| (http://bioinfo.ihb.ac.cn/gcgd). | | | | |