

# SCIENTIFIC REPORTS

OPEN

## Signatures of selection in tilapia revealed by whole genome resequencing

Jun Hong Xia<sup>1,2</sup>, Zhiyi Bai<sup>3</sup>, Zining Meng<sup>2</sup>, Yong Zhang<sup>2</sup>, Le Wang<sup>1</sup>, Feng Liu<sup>1</sup>, Wu Jing<sup>4</sup>, Zi Yi Wan<sup>1</sup>, Jiale Li<sup>3</sup>, Haoran Lin<sup>2</sup> & Gen Hua Yue<sup>1,5,6</sup>

Received: 20 January 2015

Accepted: 18 August 2015

Published: 16 September 2015

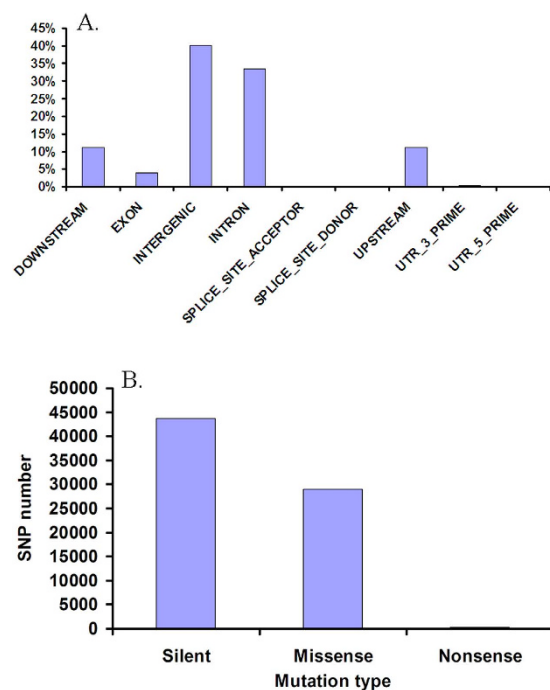
Natural selection and selective breeding for genetic improvement have left detectable signatures within the genome of a species. Identification of selection signatures is important in evolutionary biology and for detecting genes that facilitate to accelerate genetic improvement. However, selection signatures, including artificial selection and natural selection, have only been identified at the whole genome level in several genetically improved fish species. Tilapia is one of the most important genetically improved fish species in the world. Using next-generation sequencing, we sequenced the genomes of 47 tilapia individuals. We identified a total of 1.43 million high-quality SNPs and found that the LD block sizes ranged from 10–100 kb in tilapia. We detected over a hundred putative selective sweep regions in each line of tilapia. Most selection signatures were located in non-coding regions of the tilapia genome. The Wnt signaling, gonadotropin-releasing hormone receptor and integrin signaling pathways were under positive selection in all improved tilapia lines. Our study provides a genome-wide map of genetic variation and selection footprints in tilapia, which could be important for genetic studies and accelerating genetic improvement of tilapia.

Genetic selection in which the best individuals are selected as parents of the next generation is the principal tool to improve crops and livestock<sup>1</sup>. The heritable variation, which is the basis for genetic improvement in agriculture, reflects the potential of a population to respond to selection<sup>2</sup>. Adaptation in response to selection on polygenic phenotypes may occur via subtle allele frequency shifts at many loci<sup>3</sup>. The search for genomic variations and selection signatures across the entire genome can bring new insights into genes that contribute most positively to the agronomic phenotypes of animals. A number of statistical methods for identifying selection signatures have been developed, such as Fay and Wu's *H* Test<sup>4</sup>, allele frequency differences<sup>5–7</sup>, EHH<sup>8</sup>, iHS (Integrated Haplotype Score)<sup>9</sup> and Rsb<sup>10</sup>. With the advance of cost-effective and high throughput genotyping/sequencing methods, a number of genome-wide scans for selection signatures have already identified hundreds of regions targeted by recent positive selection in humans<sup>8,11</sup> and livestock, such as cattle<sup>12</sup>, chicken<sup>13</sup> and pig<sup>14</sup>. However, only a few studies on selection signatures have been conducted in aquacultured species, e.g., Channel catfish<sup>15</sup>, carp<sup>16</sup>, Atlantic salmon<sup>17</sup>, even though fish is one of the most important protein sources for humans.

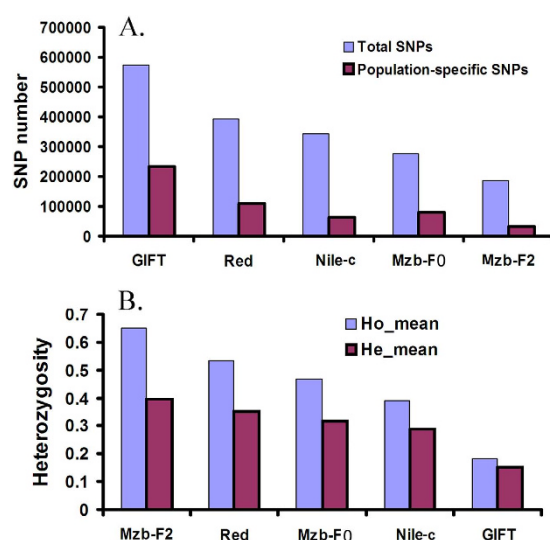
<sup>1</sup>Molecular Population Genetics and Breeding Group, Temasek Life Sciences Laboratory, 1 Research Link, National University of Singapore, 117604 Singapore. <sup>2</sup>State Key Laboratory of Biocontrol, Institute of Aquatic Economic Animals and Guangdong Provincial Key Laboratory for Aquatic Economic Animals, College of Life Sciences, Sun Yat-Sen University, Guangzhou 510275, PR China. <sup>3</sup>Key Laboratory of Exploration and Utilization of Aquatic Genetic Resources, Shanghai Ocean University, Ministry of Education, Shanghai 201306, China. <sup>4</sup>Key Laboratory of Freshwater Fisheries and Germplasm Resources Utilization, Ministry of Agriculture, Freshwater Fisheries Research Center, Chinese Academy of Fishery Sciences, Wuxi 214081, China. <sup>5</sup>Department of Biological Sciences, National University of Singapore, Singapore 117543, Singapore. <sup>6</sup>School of Biological Sciences Nanyang Technological University, 60 Nanyang Drive, Singapore, 637551, Singapore. Correspondence and requests for materials should be addressed to J.H.X. (email: xiajunh3@mail.sysu.edu.cn) or G.H.Y. (email: genhua@tll.org.sg)







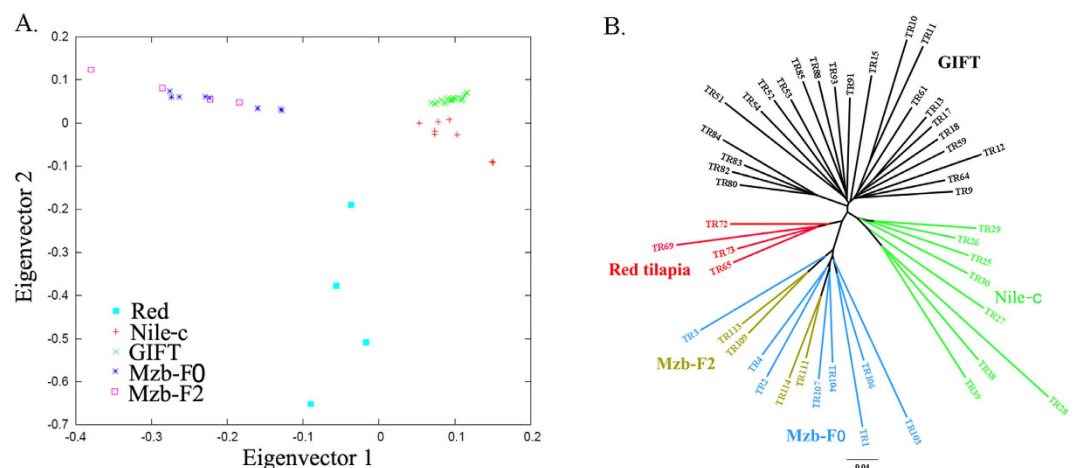
**Figure 1.** Percentage of SNPs in different locations in the tilapia genome and number of SNPs in each category of mutations in exons. (A) Percentage of SNPs in different locations in the tilapia genome; (B) Number of SNPs in each category of mutations in exons.



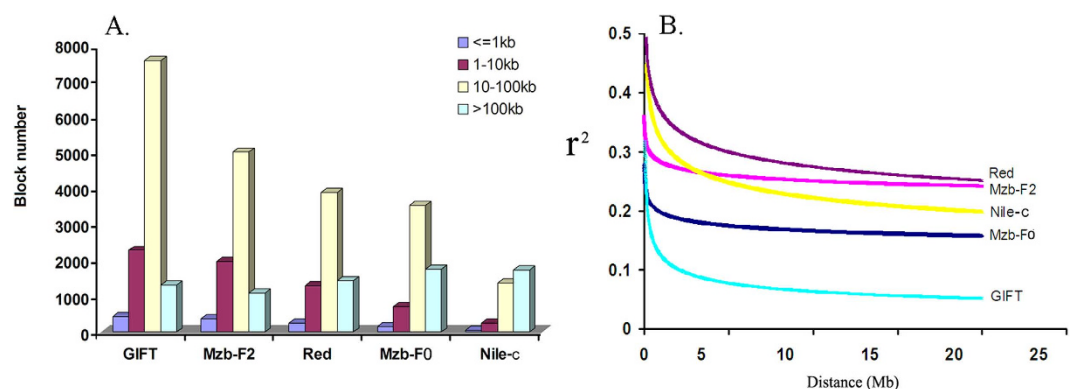
**Figure 2.** Number and heterozygosities of SNPs identified in five populations of tilapia. (A) Number of SNPs identified in five populations of tilapia; (B) Heterozygosities of SNPs in five populations of tilapia.

based on the high-quality SNPs data. Using the first and second eigenvectors, the PCA clearly divided the 47 samples into three groups (Red tilapia, Mozambique tilapia and Nile tilapia) (Fig. 3A). All the Mzb-F0 and Mzb-F2 individuals were very closely related, showing apparent differentiation from Nile-c tilapia and Red tilapia, which is consistent with our previous work<sup>26</sup>. The GIFT and Nile-c populations were closely related but clearly separated from one another, suggesting relatively little subsequent gene flow between the two tilapia lines. This genetic difference could be because these two lines have been selectively bred to adapt to different environmental conditions. PCA showed a more dispersed population substructure in the Red tilapia strain, which may be due to the cross between Mozambique and Nile tilapia during the breeding of the red tilapia<sup>26</sup>.

A NJ tree was used to cluster samples based on average genetic distances (Fig. 3B). The phylogenetic analysis showed a similar population structure as the one generated with the PCA. The NJ tree



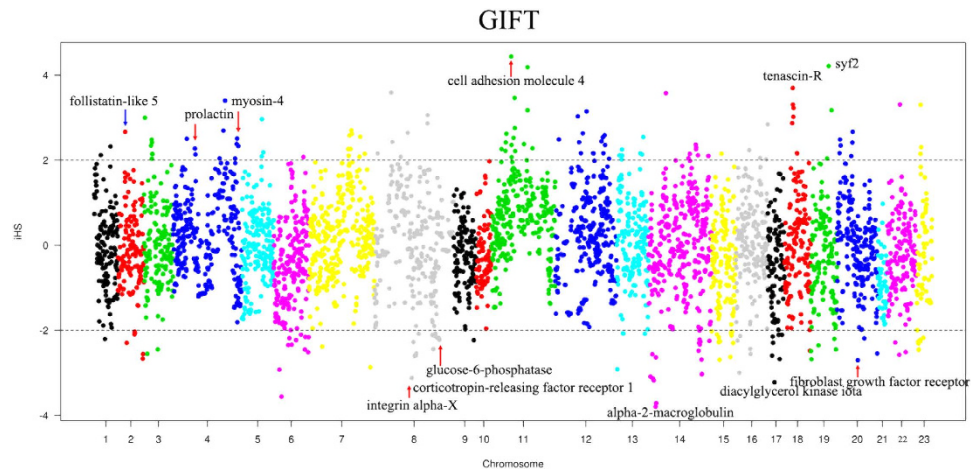
**Figure 3. Population substructure in tilapia revealed by principal component analysis and phylogenetic tree analysis.** (A) Principal component analysis. The five populations are shown in different colors and symbols; (B) Phylogenetic tree analysis. The five populations and clades are shown in different colors. The population ID and sample ID are presented.



**Figure 4. Number of different size LD blocks (A) and LD decay detected in the five populations of tilapia (B).**

contained four major groups, corresponding to Red tilapia, Mozambique tilapia, GIFT and Nile-c tilapia. As expected, the phylogenetic analysis incorporated the Mzb-F0 and Mzb-F2 into one clade. The GIFT strains were separate from Nile-c tilapia, indicating independent domestication/breeding, which is in agreement with the actual breeding process of the two lines<sup>18</sup>. The four Red samples were closely related to Mozambique tilapia and were clustered into one clade, suggesting a different domestication/breeding event from Nile and Mzb populations. These results are consistent with our previous study based on EST-originated SNP markers<sup>26</sup> where significant population structure in tilapia was detected through phylogenetic and population structure analysis. Taken together, our data suggest that there are significant population structures among the tilapia populations.

**Analysis of linkage disequilibrium (LD) decay in tilapia.** LD is useful for determining the number of markers that are needed for an association study. To estimate the LD patterns in the different tilapia groups, we calculated squared correlations of allele frequencies against genome distance ( $r^2$ ) between pairs of SNPs using Haploview<sup>39</sup>. We found slight differences among the populations at the LD block size  $>100$  kb. However, at the block size of  $<100$  kb, the highest LD block number was detected in GIFT strains and the lowest LD block number in Nile-c tilapia (Fig. 4A). Except the Nile-c tilapia population, most of the LD block sizes in the tilapia populations ranged from 10–100 kb. The plot of LD decay pattern against the genome distances is shown in Fig. 4B. Our study revealed that the LD level decreased faster in Nile-c and GIFT compared to other populations. In many previous studies, domesticated strains often showed longer LD than wild populations<sup>40</sup>. It is interesting that the GIFT population had the lowest LD level while the Red population had the highest LD level. This may suggest that the LD level in a population is more complex, and may be associated with the genetic background and breeding history of the strains. During the selection of the GIFT strain, eight founder populations from both wild and



**Figure 5. The genome-wide sweep analysis based on iHS statistics for the GIFT tilapia population.** The iHS value at a locus on each chromosome are shown in the figure. Some genes with significant SNPs when using an iHS threshold of  $\pm 2$  are shown. The LG8\_24, LG16\_21 and Unk1 scaffold are shown as Chromosome 8, 16 and 21 for clarity.

farmed populations were used<sup>19</sup>. The samples used in this study consisted of small groups collected from Shanghai and Guangzhou of China and Singapore, which have been selected for growth for over 15 generations since 1987. However, the detailed cross-breeding history of the GIFT strains used in the study is not known. To address the LD issues clearly, we have to characterize more samples in the future. Nevertheless, the information about LD patterns in each tilapia group would be useful in selecting SNPs for genome-wide association studies to identify markers associated with economically important traits on the whole genome level.

**Signatures of positive selection revealed by iHS statistics.** Population genomics has offered a new paradigm for detecting signatures of selection. To detect selective sweeps driven by artificial selection and natural selection in cultured tilapia, we used the iHS statistic<sup>9</sup>, which detects the evidence of recent positive selection at a locus by searching for genomic regions with excess homozygosity. We conducted a sweep analysis in four tilapia breeding populations. By using the iHS test, we found a total of 1120 extreme iHS values exceeding the empirical threshold level of 2 across the four populations. There were 163, 362, 243 and 352 outliers detected in the Mzb-F2, Nile-c, GIFT and Red populations, respectively (Supplementary Tables S3, S4, S5 and S6). We annotated the SNPs with the software SnpEff<sup>35</sup> by using the tilapia genome annotation information to identify candidate genes that underwent sweeps. Annotation of the regions harboring clustered iHS signals revealed 115, 243, 151 and 219 candidate genes in the Mzb-F2, Nile-c, GIFT and Red populations, respectively. Some of these candidate genes (Supplementary Tables S3, S4, S5 and S6), which were under artificial or natural selection, are novel and have not been functionally annotated yet.

The GIFT population has experienced intense artificial selection. The GIFT strains, when compared to the wild Nile tilapia populations, have gained some significant advantages in aquaculture, such as resistance to handling conditions and fast growth. In the sequencing data of the GIFT strains, a substantial number (135 out of 243 cases) of SNPs with extreme iHS values were found to be positioned in intergenic regions. Among them, 101 were located in introns, 22 in upstream and 26 in the downstream regions (Supplementary Table S5). Nonsense mutations, which are obvious candidates of functional significance, may have contributed to the rapid evolution in domestic animals<sup>41</sup>. However, in the GIFT tilapia, only one SNP was a non-synonymous mutation, which was located at 21 Mb on chromosome LG14 annotated as extracellular calcium-sensing receptor-like. A panel of interesting candidate genes containing SNPs with extreme iHS values responding to the artificial selection was identified, such as prolactin and myosin-4 in the GIFT population (Fig. 5). Previous studies showed significant associations of their homologous genes with economic traits in animals<sup>42,43</sup>, supporting the notion that these genes could be selected during breeding for genetic improvement. These candidate genes under selection identified in our study could provide useful information for rapidly identifying DNA markers associated with economically important traits in Tilapia.

In the other tilapia populations, most (more than 80%) SNPs under selection could be clustered into intergenic regions and introns. Only 3, 6 and 7 non-synonymous SNPs were detected in the Red, Mzb-F2 and Nile-c populations, respectively (Supplementary Tables S3, S4 and S6). We found non-synonymous SNPs in genes OR11A1, EPS15L1 and EFNA2 in the red population; OR11A1, hypothetical protein LOC100700781, CDAN1, MR1, hypothetical protein LOC100690504 and zwilch homolog in the Mzb-F2 population; and IL13RA2, WDFY3, GPRC6A, PDK1, Slc9a3, Anpep and HPS5 in the Nile-c







54. Patel, R. K. & Jain, M. NGS QC Toolkit: a toolkit for quality control of next generation sequencing data. *Plos One* **7**, e30619 (2012).
55. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat Methods* **9**, 357–359 (2012).
56. Li, H. *et al.* The sequence alignment/map format and SAMtools. *Bioinformatics* **25**, 2078–2079 (2009).
57. Price, A. L. *et al.* Principal components analysis corrects for stratification in genome-wide association studies. *Nat Genet.* **38**, 904–909 (2006).
58. Lee, T. H., Guo, H., Wang, X., Kim, C. & Paterson, A. H. SNPhylo: a pipeline to construct a phylogenetic tree from huge SNP data. *BMC Genomics* **15**, 162 (2014).
59. Jeanmougin, F., Thompson, J. D., Gouy, M., Higgins, D. G. & Gibson, T. J. Multiple sequence alignment with Clustal X. *Trends Biochem Sci.* **23**, 403–405 (1998).
60. Gautier, M. & Vitalis, R. rehh: An R package to detect footprints of selection in genome-wide SNP data from haplotype structure. *Bioinformatics* **28**, 1176–1177 (2012).

## Acknowledgments

This research is supported by the National Research Foundation, Prime Minister's Office, Singapore under its Competitive Research Program (CRP Award No. NRF-CRP002-001). We thank the Broad Institute for providing us the assembled genome sequence of Nile tilapia. We are very grateful to our colleagues Ms. May Lee and Mr. Baoqing Ye for English editing, as well as members of the aquaculture team of our institute for taking care of our tilapia breeding stocks.

## Author Contributions

G.H.Y. coordinated and supervised the project. J.H.X., Z.Y.B., Z.N.M., Z.Y., Z.Y.W., L.W., F.L., W.J., H.R.L., J.L.L. and G.H.Y. designed the research and performed experiments. J.H.X. analyzed data and wrote the paper. G.H.Y. finalized the paper.

## Additional Information

**Supplementary information** accompanies this paper at <http://www.nature.com/srep>

**Competing financial interests:** The authors declare no competing financial interests.

**How to cite this article:** Hong Xia, J. *et al.* Signatures of selection in tilapia revealed by whole genome resequencing. *Sci. Rep.* **5**, 14168; doi: 10.1038/srep14168 (2015).



This work is licensed under a Creative Commons Attribution 4.0 International License. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in the credit line; if the material is not included under the Creative Commons license, users will need to obtain permission from the license holder to reproduce the material. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>