

Supplementary Material for Genomic Analysis of Amphioxus Reveals a Wide Range of Fragments Homologous to Viral sequences

Authors: Qiao Du ^{1,†}, Fang Peng ^{1,†}, Qing Xiong ^{2,3,†}, Kejin Xu ¹, Kevin Yi Yang ^{2,3}, Mingqiang Wang ^{2,4}, Zhitian Wu ¹, Shanying Li ¹, Xiaorui Cheng ¹, Xinjie Rao ¹, Yuyouye Wang ¹, Stephen Kwok-Wing Tsui ^{2,3} and Xi Zeng ^{1,*}

Affiliations:

¹ Agricultural Bioinformatics Key Laboratory of Hubei Province and 3D Genomics Research Centre, College of Informatics, Huazhong Agricultural University, Wuhan 430070, China

² School of Biomedical Sciences, The Chinese University of Hong Kong, Hong Kong, China

³ Hong Kong Bioinformatics Centre, The Chinese University of Hong Kong, Hong Kong, China

⁴ Stanford Cardiovascular Institute, Stanford University School of Medicine, Stanford, CA 94305, USA

* Correspondence: zengxi@mail.hzau.edu.cn; Tel.: +86-15818755917

† These authors contributed equally to this work.

This file includes:

Table S1 to S4

Figure S1 to S4

Table S1. The 17 viruses with viral HFs

Virus accession	Virus name	Number of HFs in viruses	Number of HFs in amphioxus	Number of HFs in viral CDS
NC_019491.1	<i>Cyprinid herpesvirus 1</i>	10	2	0
NC_021858.1	<i>Pandoravirus dulcis</i>	7	78	7
NC_026440.1	<i>Pandoravirus inopinatum</i>	7	73	3
NC_008168.1	<i>Choristoneura occidentalis granulovirus</i>	5	30	1
NC_006639.1	<i>Cotesia congregata bracovirus</i>	4	42	4
NC_026421.1	<i>Equid gammaherpesvirus 5</i>	3	10	0
NC_028045.1	<i>Tadarida brasiliensis circovirus 1</i>	2	37	0
NC_028094.1	<i>Chrysochromulina ericina virus</i>	2	8	0
NC_001716.2	<i>Human betaherpesvirus 7</i>	2	1	0
NC_001550.1	<i>Mason-Pfizer monkey virus</i>	1	29	0
NC_022098.1	<i>Pandoravirus salinus</i>	1	21	1
NC_026141.2	<i>Adelie penguin polyomavirus</i>	1	5	1
NC_008603.1	<i>Paramecium bursaria Chlorella virus FR483</i>	1	4	0
NC_008724.1	<i>Acanthocystis turfacea Chlorella virus 1</i>	1	4	1
NC_000852.5	<i>Paramecium bursaria Chlorella virus 1</i>	1	1	1
NC_008094.1	<i>Y73 sarcoma virus</i>	1	1	1
NC_023006.1	<i>Pseudomonas phage PPpW-3</i>	1	1	1

Table S2. The HFs in CDS regions of viral genomes

Virus accession	Virus name	Number of HFs in CDS	CDS annotation
NC_021858.1	<i>Pandoravirus dulcis</i>	7	Histone H2B domain containing protein
NC_006639.1	<i>Cotesia congregata bracovirus</i>	4	Histone (histone H4 like)
NC_026440.1	<i>Pandoravirus inopinatum</i>	3	hypothetical protein
NC_008094.1	<i>Y73 sarcoma virus</i>	1	protein-tyrosine kinase
NC_008168.1	<i>Choristoneura occidentalis granulovirus</i>	1	similar to Cydia pomonella granulovirus orf35 (hypothetical protein)
NC_008724.1	<i>Acanthocystis turfacea chlorella virus 1</i>	1	hypothetical protein
NC_000852.5	<i>Paramecium bursaria Chlorella virus 1</i>	1	hypothetical protein
NC_022098.1	<i>Pandoravirus salinus</i>	1	Histone H2B domain
NC_023006.1	<i>Pseudomonas phage PPpW-3</i>	1	PPpW-3_ORF-54 (hypothetical protein)
NC_026141.2	<i>Adelie penguin polyomavirus</i>	1	large T antigen

Table S3. The 10 amphioxus genes with the most HFs

Gene name (description)	Number of HFs	Product of gene
<i>Histone H2B 1/2</i>	78	Histone H2B 1/2
<i>Histone H4</i>	50	Histone H4
<i>Late histone H2B.2.1</i>	30	Histone H2B
<i>Transposon TX1 uncharacterized 149 kDa protein</i>	14	Transposon TX1 uncharacterized 149 kDa protein
<i>Histone H2B (Fragments)</i>	8	Histone H2B fragments
<i>DCST1 (E3 ubiquitin-protein ligase DCST1)</i>	7	E3 ubiquitin-protein ligase DCST1
<i>H2BC13 (Histone H2B type 1-L)</i>	7	Histone H2B type 1-L
<i>DCST2 (DC-STAMP domain- containing protein 2)</i>	6	DC-STAMP domain-containing protein 2
<i>hist2h2l (Histone H2B 3)</i>	6	Histone H2B 3
<i>Histone H2B</i>	4	Histone H2B

Table S4. The amphioxus genes with HFs in CDS regions

Gene name	Number of HFs in CDS	Number of HFs in gene
<i>Histone H2B 1/2</i>	77	78
<i>Histone H4</i>	45	51
<i>Late histone H2B.2.1</i>	30	30
<i>H2BC13 (Histone H2B type 1-L)</i>	7	7
<i>Histone H2B (Fragments)</i>	7	8
<i>hist2h2l (Histone H2B 3)</i>	6	6
<i>Histone H2B</i>	4	4
<i>AGAP012199 (Histone H2B)</i>	2	2
<i>DCST1 (E3 ubiquitin-protein ligase DCST1)</i>	2	7
<i>H2AJ (Histone H2A.J)</i>	2	2
<i>H2bc3 (Histone H2B type 1-B)</i>	2	2
<i>Histone H3</i>	2	2
<i>H4DEKL (Histone H4)</i>	1	1
<i>Maob (Amine oxidase [flavin-containing] B)</i>	1	1
<i>Src42A (Tyrosine-protein kinase Src42A)</i>	1	1
<i>vit-6 (Vitellogenin-6)</i>	1	1
<i>VPS52 (Vacuolar protein sorting-associated protein 52 homolog)</i>	1	1

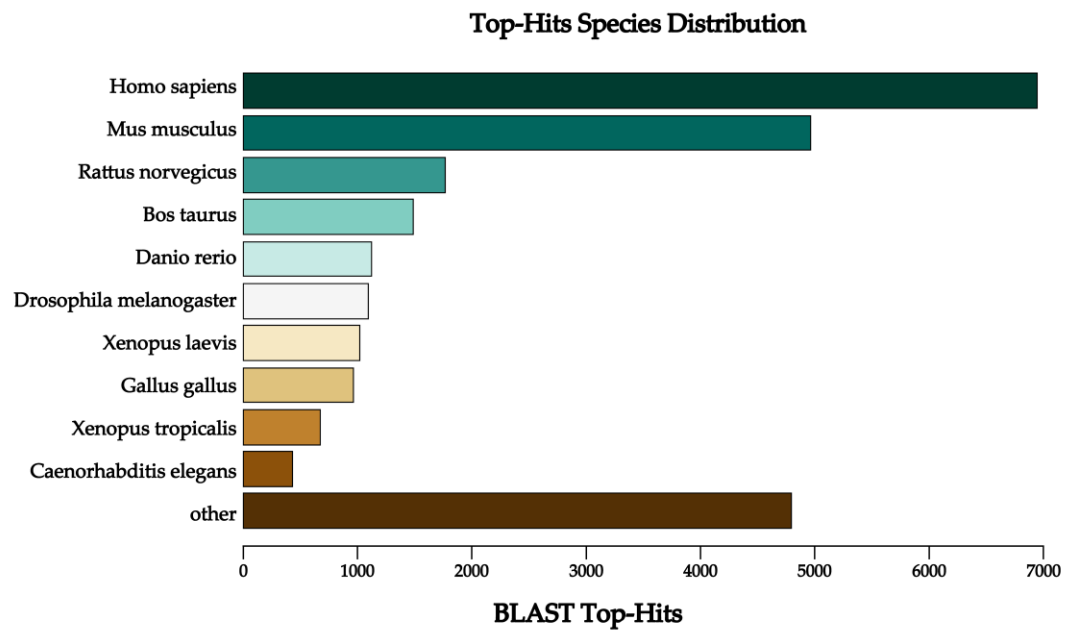


Figure S1. The species of top BLAST hits in functional annotation of the amphioxus genome.

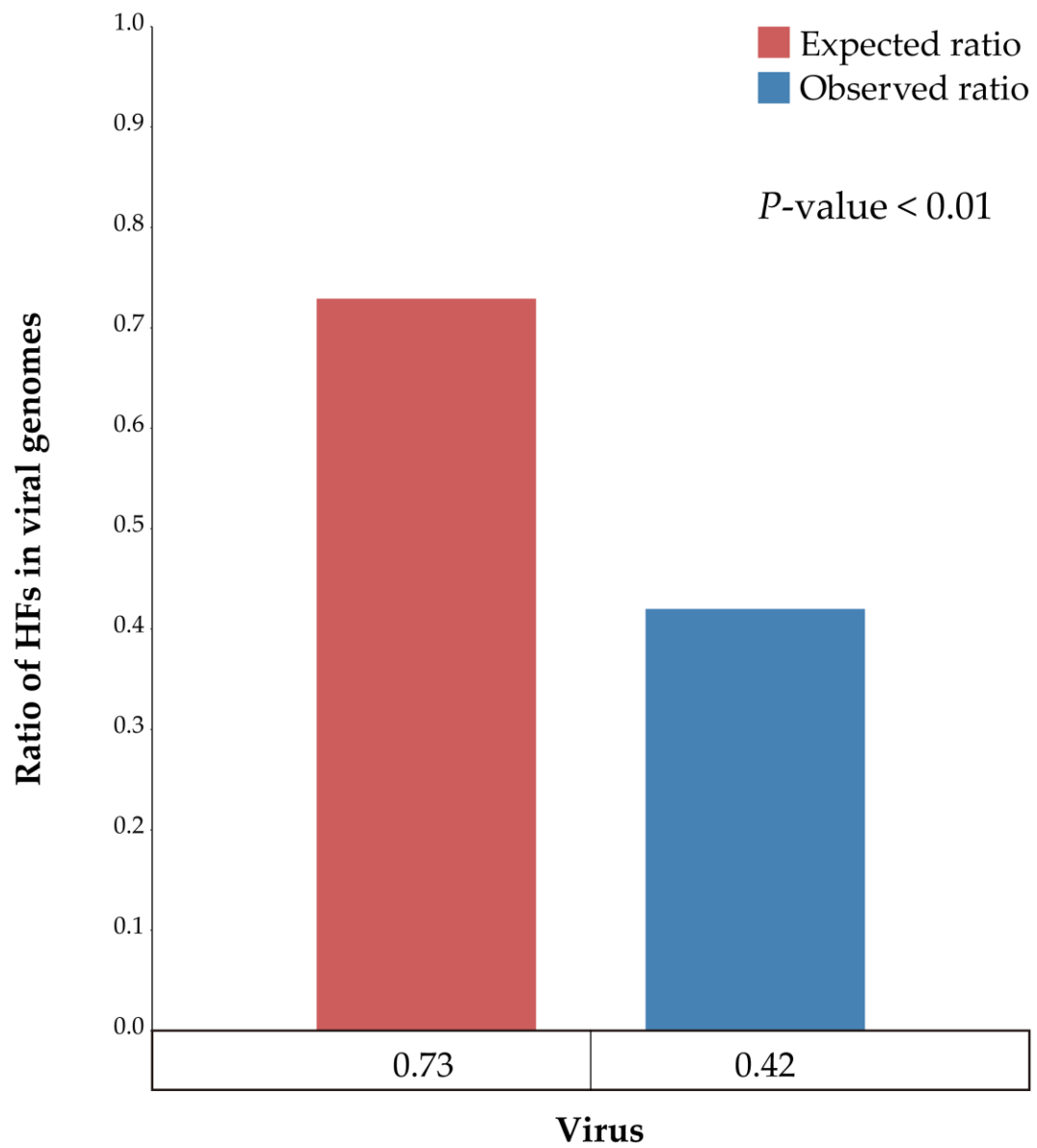


Figure S2. The distribution of HFs in CDS regions of viral genomes. The expected (random distribution, red) and the observed (actual numbers, blue) ratios of HFs in CDS regions are shown. The p value was calculated by χ^2 test.

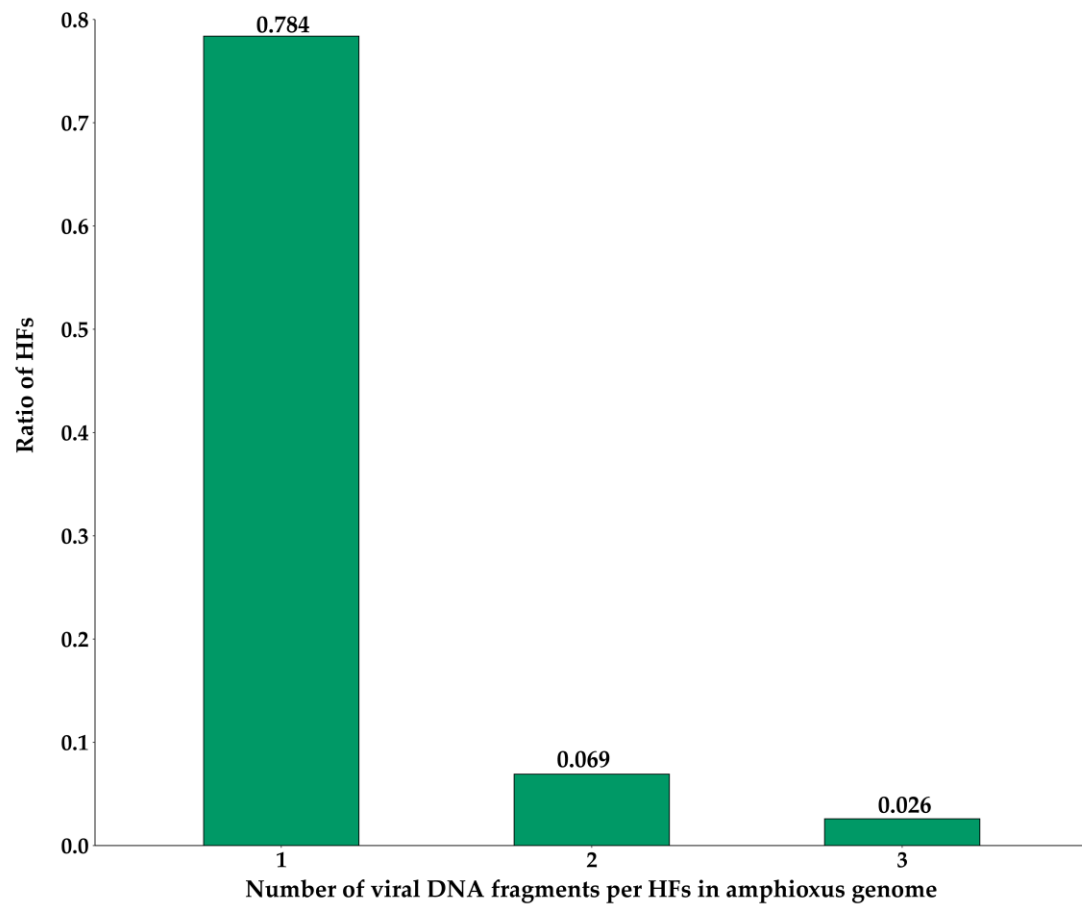


Figure S3. The number of viral DNA fragments which were homologous to each HF in the *B. belcheri* beihai genome. One HF's of the *B. belcheri* beihai genome could be homologous to more than one DNA fragment of viral genomes.

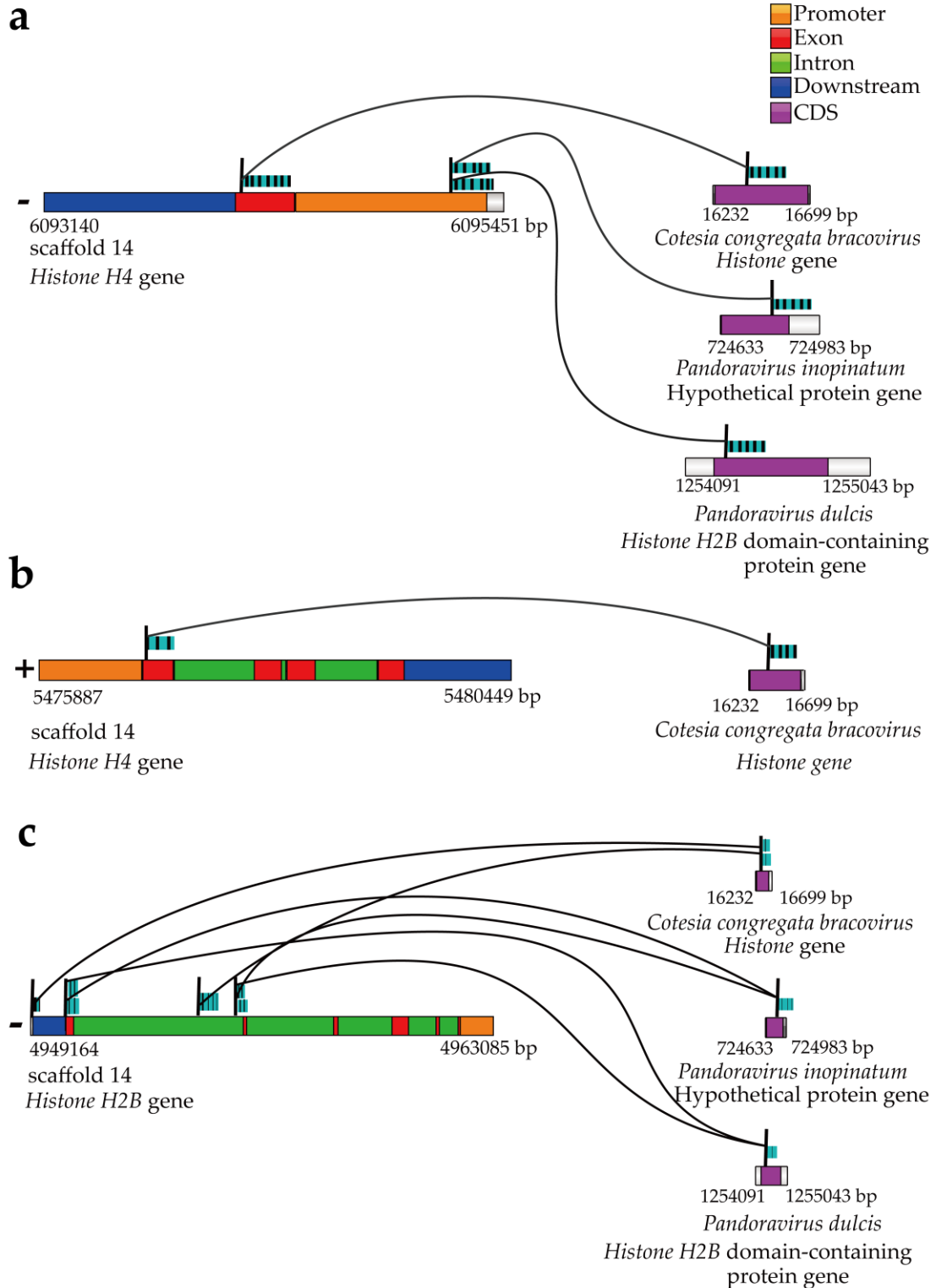


Figure S4. The HFs of amphioxus genes containing the HFs with the highest confidence of long length and high alignment quality. The 5 HFs with the highest confidence were located in *Histone H4* and *Histone H2B* of the amphioxus and were homologous to *C. congregata bracovirus*. In addition to the HFs with the top longest length and the highest alignment quality, other HFs within the gene were also shown. There were multiple viral DNA fragments homologous to one DNA fragment of amphioxus and one viral DNA fragment were possibly homologous to multiple DNA fragments of amphioxus. The bars on the left represent the amphioxus genome; different colors represent different genomic regions of

amphioxus; the dark green bars with vertical line represent the HFs of amphioxus. The bars on the right represent viral genomes; different colors represent different genomic regions of virus; the dark green bars with vertical line represent viral HFs. Yellow represents promoter regions in the amphioxus genome; red represents exon regions in the amphioxus genome; green represents intron regions in the amphioxus genome; blue represents downstream regions in the amphioxus genome; purple represents CDS regions in viral genomes. (a) The HFs of *Histone H4* on scaffolds 14. (b) The HFs of *Histone H4* on scaffolds 14. (c) The HFs of *Histone H2B* on scaffolds 14.