















Horseshoe crab genomes reveal the evolution of genes and microRNAs after three rounds of whole genome duplication

Wenyan Nong ^{1,19}, Zhe Qu^{1,19}, Yiqian Li^{1,19}, Tom Barton-Owen^{1,19}, Annette Y. P. Wong^{1,19}, Ho Yin Yip¹, Hoi Ting Lee¹, Satya Narayana¹, Tobias Baril ², Thomas Swale³, Jianquan Cao ¹, Ting Fung Chan ⁴, Hoi Shan Kwan⁵, Sai Ming Ngai⁴, Gianni Panagiotou ^{6,7}, Pei-Yuan Qian ⁸, Jian-Wen Qiu ⁹, Kevin Y. Yip ¹⁰, Noraznawati Ismail¹¹, Siddhartha Pati ^{12,13,14}, Akbar John ¹⁵, Stephen S. Tobe ¹⁶, William G. Bendena ¹⁷, Siu Gin Cheung¹⁸, Alexander Hayward ² & Jerome H. L. Hui ¹✉

Whole genome duplication (WGD) has occurred in relatively few sexually reproducing invertebrates. Consequently, the WGD that occurred in the common ancestor of horseshoe crabs ~135 million years ago provides a rare opportunity to decipher the evolutionary consequences of a duplicated invertebrate genome. Here, we present a high-quality genome assembly for the mangrove horseshoe crab *Carcinoscorpius rotundicauda* (1.7 Gb, N50 = 90.2 Mb, with 89.8% sequences anchored to 16 pseudomolecules, $2n = 32$), and a resequenced genome of the tri-spine horseshoe crab *Tachypleus tridentatus* (1.7 Gb, N50 = 109.7 Mb). Analyses of gene families, microRNAs, and synteny show that horseshoe crabs have undergone three rounds (3R) of WGD. Comparison of *C. rotundicauda* and *T. tridentatus* genomes from populations from several geographic locations further elucidates the diverse fates of both coding and noncoding genes. Together, the present study represents a cornerstone for improving our understanding of invertebrate WGD events on the evolutionary fates of genes and microRNAs, at both the individual and population level. We also provide improved genomic resources for horseshoe crabs, of applied value for breeding programs and conservation of this fascinating and unusual invertebrate lineage.

Polyploidy provides new genetic raw material for evolutionary diversification, as gene duplication can lead to the evolution of new gene functions and regulatory networks¹. Nevertheless, whole-genome duplication (WGD) is a relatively rare occurrence in animals when compared to fungi and plants². Two rounds of ancient WGD occurred in the last common ancestor of the vertebrates, with additional rounds in some teleost fish lineages²⁻⁴. Fixation of these WGD events (i.e., ‘polyploidization’) is considered a major force in shaping the evolutionary success of vertebrate lineages, by facilitating fundamental changes in physiology and morphology, leading to the origin of new adaptations^{5,6}. Among the invertebrates, horseshoe crabs⁷⁻⁹, spiders, and scorpions¹⁰ represent the only sexually reproducing lineages that are known to have undergone WGD (Fig. 1a).

Horseshoe crabs are considered to be ‘living fossils’. The oldest actual fossils of horseshoe crabs date to the Ordovician period ~450 million years ago (Mya)¹¹, and remarkably, extant species remain relatively unchanged morphologically since this extremely ancient date. However, despite their long history, there are only four extant species of horseshoe crabs worldwide: the Atlantic horseshoe crab (*Limulus polyphemus*) from the Atlantic East

Coast of North America, and the mangrove horseshoe crab (*Carcinoscorpius rotundicauda*), the Indo-Pacific horseshoe crab (*Tachypleus gigas*), and the tri-spine horseshoe crab (*Tachypleus tridentatus*), from South and East Asia¹². All extant horseshoe crabs are estimated to have diverged from a common ancestor that existed ~135 Mya¹³, and they share an ancestral WGD⁹. A high-quality genome assembly was recently announced as a genomic resource for *T. tridentatus*^{14,15}, leaving an exciting research opportunity to analyse the genomes of other horseshoe crab species to understand how WGD events reshape the genome and rewire genetic regulatory networks in invertebrates.

In the present study, we provide the first high quality genome of the mangrove horseshoe crab (*C. rotundicauda*), and a resequenced genome of the tri-spine horseshoe crab (*T. tridentatus*). Importantly, we present evidence for the number of rounds of WGD that have occurred in these genomes, and investigate if these represent a shared event with spiders. We also examine the evolutionary fate of genes and microRNAs at both the individual and population level in these genomes. Collectively, this study highlights the evolutionary consequences of a unique invertebrate WGD, while at the same time providing detailed genetic insights of utility for diverse genomic, biomedical, and conservation applications.

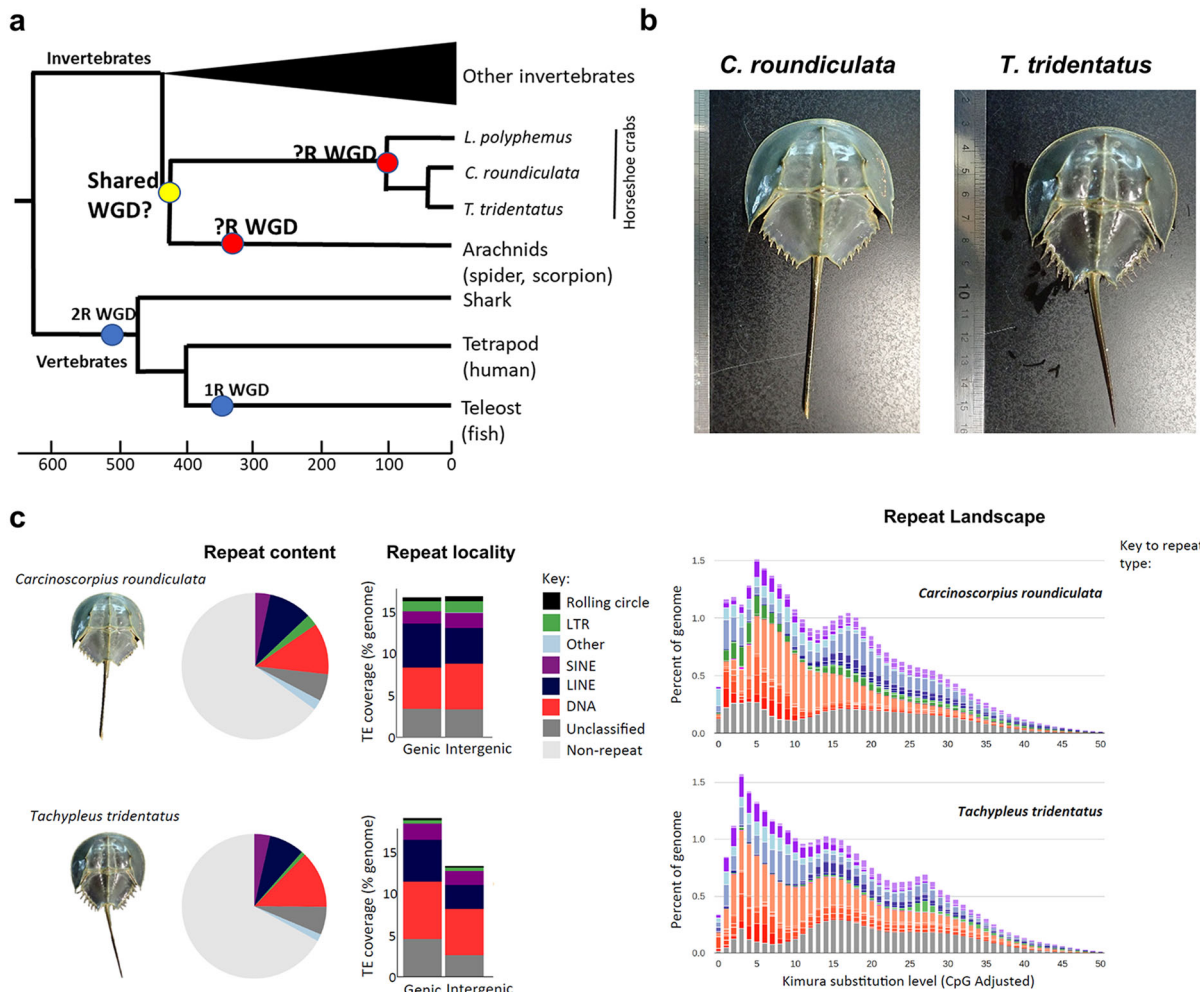


Fig. 1 Horseshoe crabs *C. rotundicauda* and *T. tridentatus*. **a** Schematic diagram illustrating the current knowledge of whole-genome duplication (WGD) in animals. ‘?R’ denotes unknown rounds of whole-genome duplication; **b** pictures of horseshoe crabs *C. rotundicauda* and *T. tridentatus*; **c** Repeat content for the two horseshoe crab genomes, *C. rotundicauda* and *T. tridentatus*: Pie charts illustrating repeat content as a proportion of total genomic content; Repeat content present in genic versus intergenic regions; and Repeat landscape plots illustrating transposable element activity in each horseshoe crab genome. Source data reveals these figures can be found in Supplementary Data 8.

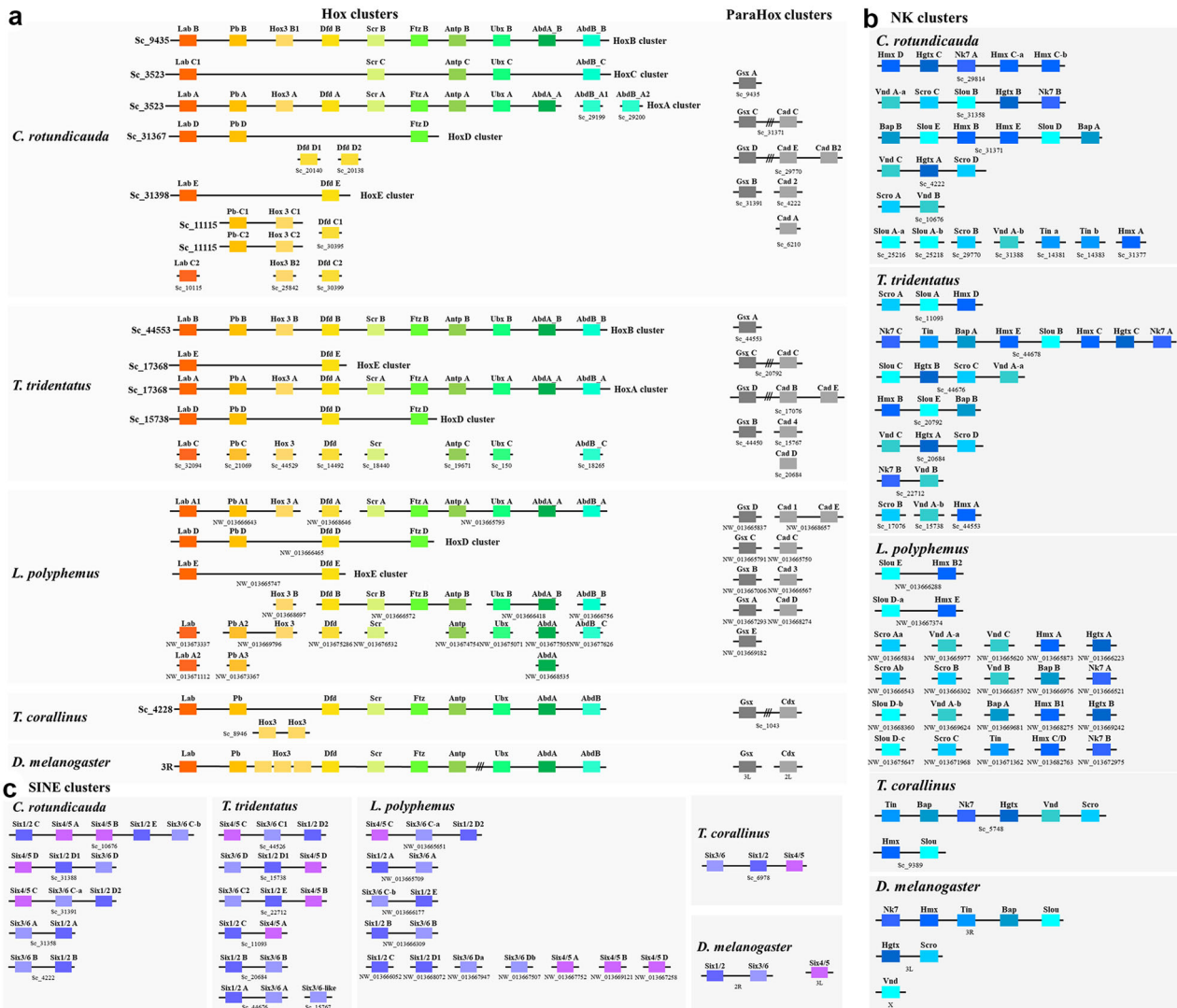


Fig. 2 Homeobox gene organisation. **a** Genomic organisation of the Hox (left) and ParaHox (right) cluster genes in the horseshoe crab genomes. **b** Genomic organisation of the NK and **c** SINE cluster genes in the horseshoe crab genomes. Note that the *L. polyphemus* genome assembly (GCF_000517525) was obtained from Battelle et al. [8].

first close to chromosomal-level genome assembly for this species. Consequently, the two high-quality horseshoe crab genomes presented in this study provide an unprecedented opportunity to address the issue of invertebrate WGD and its evolutionary consequences.

An important outstanding question is how many rounds of WGD occurred in the last common ancestor of horseshoe crabs, or alternatively, if all rounds of WGD occurred already in the ancestor of arachnids and horseshoe crabs (Fig. 1a)? To address this question, we first investigated the number and genomic location of *Hox* cluster genes, which have played the role of a ‘Rosetta stone’ for understanding animal evolution¹⁸. For example, the genome of the cephalochordate amphioxus contains only a single *Hox* gene cluster with 15 *Hox* genes, while the mouse genome contains four *Hox* gene clusters with 39 *Hox* genes, providing evidence that two rounds of WGD occurred between the most recent common ancestor of amphioxus and human^{19,20}. In our horseshoe crab genomes for *C. rotundicauda* and *T. tridentatus*, the number of *Hox* genes was found to be 43 and 36, respectively (Fig. 2a, Supplementary Data 1). In *C. rotundicauda*, we found that there are five *Hox* clusters, with other *Hox* genes located on additional small scaffolds; while in *T. tridentatus*, there

are three *Hox* clusters, again with other *Hox* genes scattered across different scaffolds (Fig. 2a). This situation is similar to that for the genome assembly of the Atlantic horseshoe crab *L. polyphemus*⁸, where our analyses demonstrated that there are four *Hox* clusters with additional *Hox* genes located on different scaffolds. In a recent study of the *T. tridentatus* re-sequenced genome, the authors could only detect two *Hox* clusters, and could not identify the *Ftz* gene inside these clusters¹⁴. In contrast, our results suggest that there are in fact three *Hox* clusters (including *Ftz*), and thus more than one round of WGD occurred in the lineage leading to extant horseshoe crabs.

We next investigated the sister cluster of the *Hox* genes—the *ParaHox* cluster genes, which are also highly clustered in bilaterians^{21–23}. Similar to the *Hox* cluster genes, the cephalochordate amphioxus contains only a single *ParaHox* gene cluster in its genome, while the *ParaHox* cluster genes are located on four chromosomes in human¹⁹. In comparison, both the horseshoe crab genomes for *C. rotundicauda* and *T. tridentatus* contain two *ParaHox* clusters, composed of *Gsx* and *Cdx*, with other *ParaHox* genes located on three additional scaffolds. Meanwhile, in the genome assembly of *L. polyphemus*⁸, perhaps due to the lower sequence continuity of the genome (i.e. low

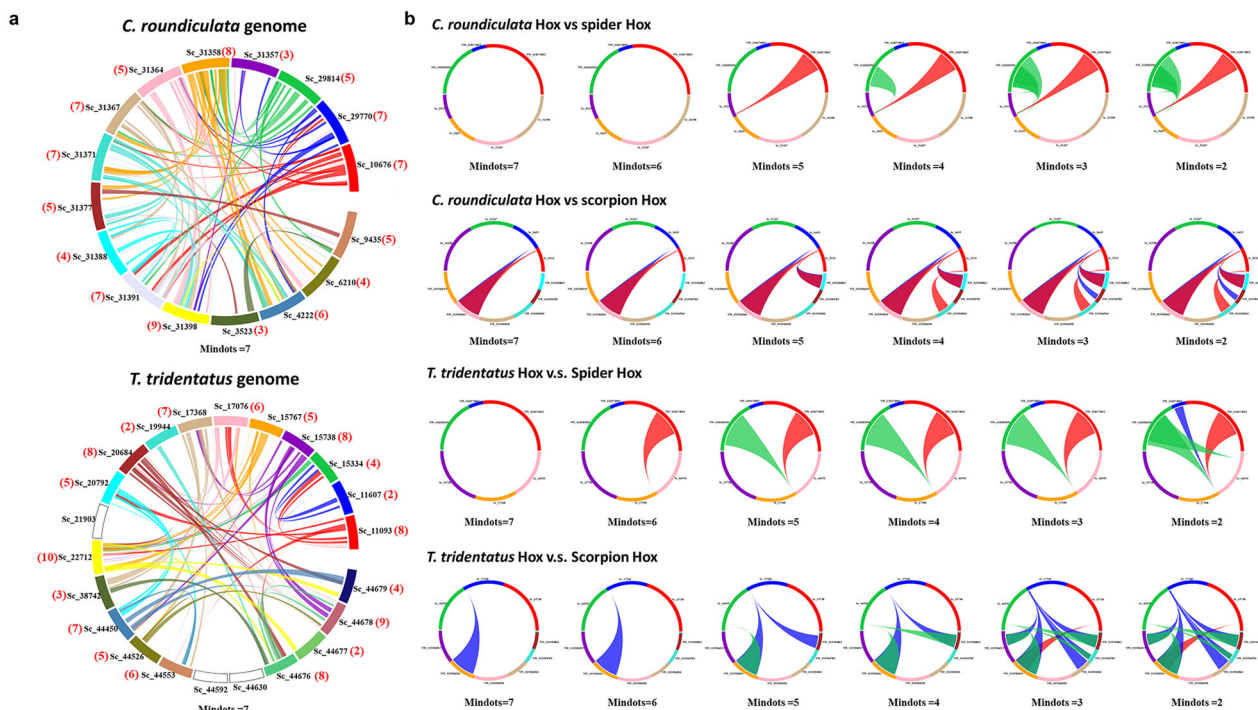


Fig. 3 Syntenic regions between chelicerate genomes. **a** Synteny between different chromosomes of *C. roundiculata* and *T. tridentatus*. Note that the bracketed numbers highlighted in red refer to the numbers of chromosomes that syntenic blocks with that chromosome (counting include its own copy). **b** Syntenic relationships of Hox scaffolds of (Upper panel): *C. roundiculata*, spider and scorpion; (Lower panel) and *T. tridentatus*, spider, and scorpion.

scaffold N50), only a single *ParaHox* cluster for *Cdx* was identified, with the other *ParaHox* genes located on eight additional scaffolds (Fig. 2a). For other well-known homeobox gene clusters, including the *NK* cluster and *SINE* clusters, multiple clusters were revealed, as above (Fig. 2b, c). In *C. rotundicauda* and *T. tridentatus*, five and seven *SINE* clusters were found respectively, while in the genome assembly of *L. polyphemus*⁸, four *SINE* clusters were revealed, with the other six genes located elsewhere in the genome.

Using genome-wide analyses of homeobox gene content in three horseshoe crab genomes, we find that many homeobox genes are present in more than four copies (details are provided in Supplementary Table 6, Supplementary Figs. 3–7). These results suggest that at least two rounds (2R), and likely three rounds (3R) of WGD have occurred. So the question then becomes, exactly how many rounds of WGD did occur in horseshoe crabs? To address this question, we carried out additional genome-wide synteny analyses. We found that, using a default set of a minimum of seven genes to define a syntenic block, most of the chromosomes of *C. rotundicauda* exhibit synteny with other chromosomes (between 4–8 including its own copy) (Fig. 3a). Thus, we propose that three rounds of whole-genome duplication (3R WGD) occurred in horseshoe crabs.

Shared or independent duplications with spiders? Another major unresolved question relating to horseshoe crab genomes is whether the reported cases of WGD in chelicerates constitute shared or independent events. Gene family analyses of spider and scorpion genomes have suggested that an ancient WGD is shared between them, independent of the further WGDs that occurred in the horseshoe crab lineage¹⁰. Using the two horseshoe crab genome assemblies generated here, this study addressed this important question from two different perspectives by: (1) performing analyses of synteny, and (2) considering recent new evidence on phylogenetic relationships within the Chelicerata.

We first carried out analyses of synteny between the *Hox* scaffolds of *C. rotundicauda* and published spider and scorpion genomes¹⁰ (Fig. 3b). Despite no clear shared duplication event between *C. rotundicauda* and spider *Hox* clusters, surprisingly, we observed syntenic relationships between two *Hox* scaffolds when using a minimum of five genes to define a syntenic block (Fig. 3b). Similarly, during synteny comparisons of *Hox* scaffolds of *T. tridentatus* and published spider and scorpion genomes, we observed syntenic relationships between two different *Hox* scaffolds when using a minimum of five genes to define a syntenic block (Fig. 3b). Under less stringent criteria using a minimum of two genes to define a syntenic block, we additionally observed syntenic relationships for two other *Hox* scaffolds between *T. tridentatus* and spider genomes (Fig. 3b).

An important consideration necessary to fully understand WGD events identified for horseshoe crab and chelicerate genomes is the phylogenetic relationships between these lineages. Horseshoe crabs have long been regarded as a monophyletic group (Xiphosura) and the sister group to the terrestrial chelicerate clade that includes spiders, scorpions, ticks, mites, harvestmen, and solifuges (Arachnida). However, in a recent phylogenetic analysis using publicly available data, including three xiphosurans, two pycnogonids, and 34 arachnids, it was suggested that the horseshoe crabs represent a group of marine arachnids²⁴. Conversely, another group of researchers recovered the Xiphosura as the sister group to the Arachnida²⁵, suggesting a single terrestrialisation event occurred after the last common ancestor of arachnids and horseshoe crabs diverged. Consequently, additional analyses and data are needed to differentiate between these scenarios, and fully demonstrate whether a shared WGD event occurred at the common ancestor of horseshoe crabs, spiders, and scorpions; or if an ancestral WGD occurred in the ancestral lineage to chelicerates and xiphosurans, followed by massive gene losses in some lineages, such as ticks and mites; or multiple WGDs originated independently in and arachnopholpnates and xiphosurans.

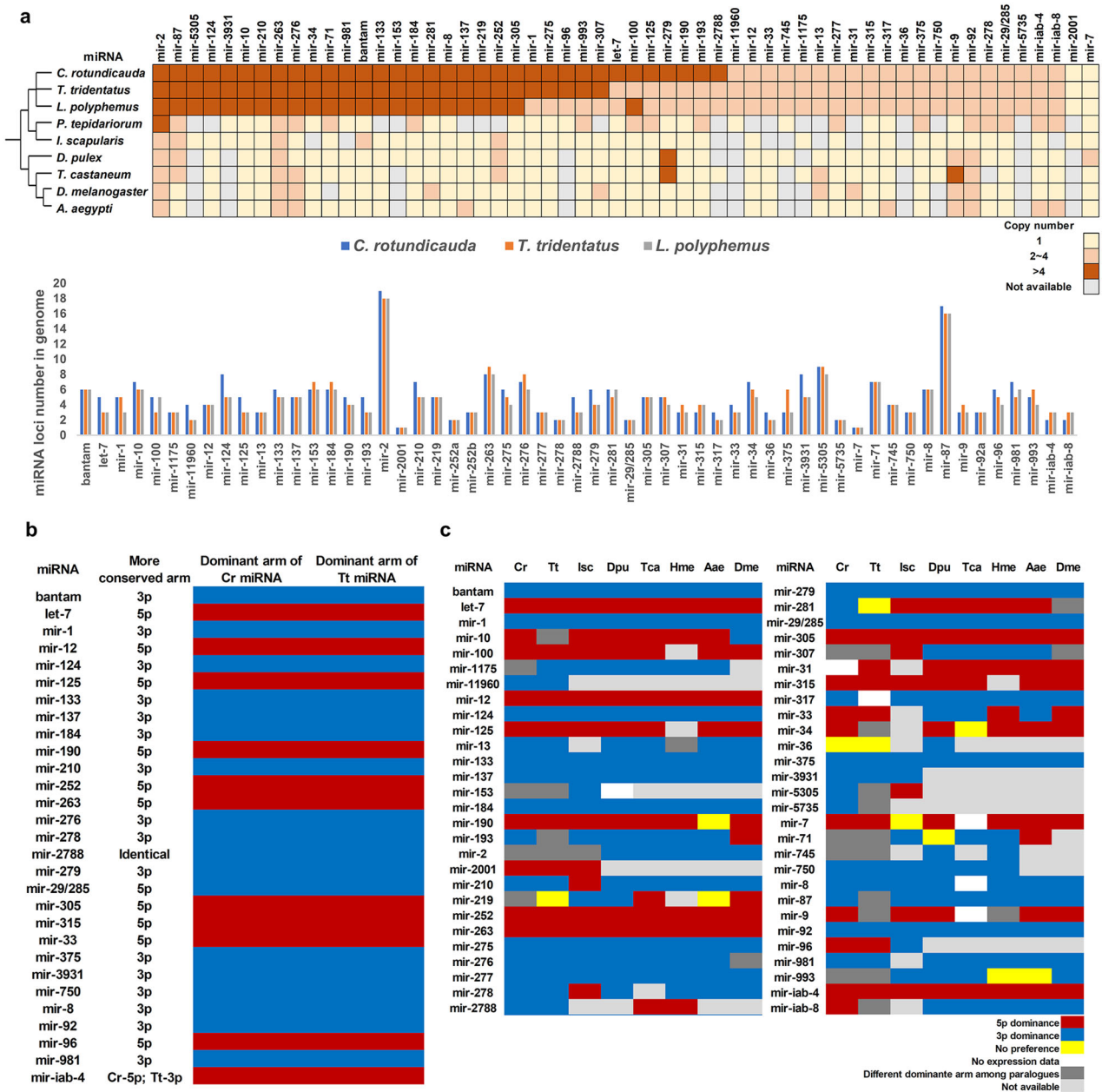


Fig. 4 MicroRNA copies and dominant arm usage. **a** Number of gene copies of conserved microRNAs in the arthropod genomes; **b** Sequence conservation and arm switching of horseshoe crab microRNAs; **c** Comparison of microRNA arm preference among different arthropod species. *Isc* *Xodes scapularis*, *Dpu* *Daphnia pulex*, *Tca* *Tribolium castaneum*, *Hme* *Heliconius melpomene*, *Aae* *Aedes aegypti*, *Dme* *Drosophila melanogaster*. Arm preference: blue—3p dominance, red—5p dominance, yellow—no preference, white—no expression. Source data reveals Fig. 4a can be found in Supplementary Data 8.

Duplicated fates of noncoding microRNAs. The availability of new transcriptomic data, especially the first small RNA transcriptomic data for both species of horseshoe crabs (Supplementary Tables 7 and 8), enabled us to analyse the evolutionary consequences of small noncoding RNAs after the WGD events in both *C. rotundicauda* and *T. tridentatus*. To reveal if duplicated microRNAs can also provide insights into the number of rounds of WGD that occurred, we first examined the number of paralogues for the bilaterian conserved set of 57 microRNAs, across three horseshoe crab genomes (Fig. 4a, Supplementary Data 2–5). Of these microRNAs, 27 and 34 have more than 4 copies in *T. tridentatus* and *C. rotundicauda*, respectively (Fig. 4a, Supplementary Fig. 8). These data further support the hypothesis that 3R WGD occurred in horseshoe crabs.

To understand the fates of microRNA paralogues, we first analysed sequence divergence in 41 conserved microRNA families and 4 chelicerate-specific microRNAs by aligning their sequences (Supplementary Figs. 9–21, Supplementary Data 6). We found that the paralogues always have greater sequence conservation in one arm (rather than showing similar conservation for both arms across paralogues) after WGD (Supplementary Figs. 9–21). An example is illustrated for the microRNA *bantam*, where the sequence of the 5p arm is less conserved than the 3p arm between paralogues (Supplementary Fig. 22).

To explore whether greater conservation in the 3p microRNA arm correlates with a greater expression level, we mapped small RNA reads to different paralogues. By eliminating microRNA species which have different arm usage between their paralogues or between horseshoe crab species, we found that out of the 29

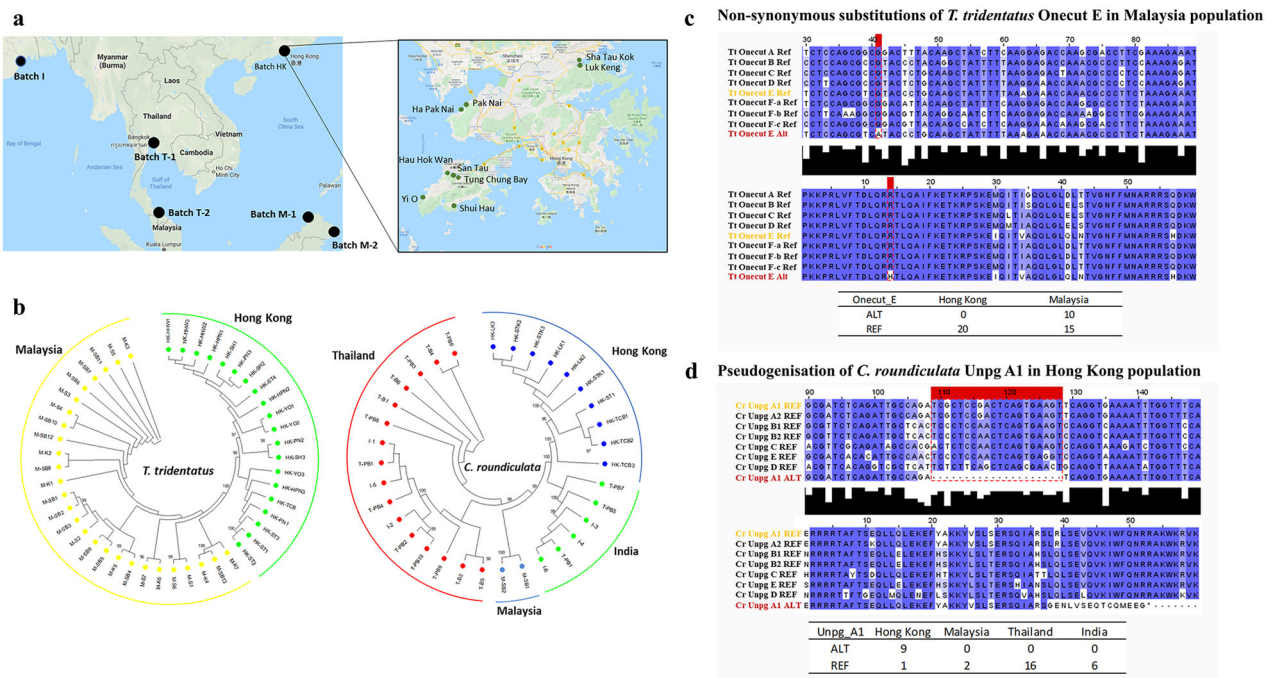


Fig. 5 Population genomics and evolutionary fates of paralogues. **a** Geographical distribution of *C. roundiculata* and *T. tridentatus* collected samples; **b** Phylogenetic trees of the collected samples. **c** Non-synonymous substitutions of *T. tridentatus* *Oncut-E* genes in individuals collected in Malaysia; **d** Pseudogenisation of *C. roundiculata* *Unpg-A1* gene in individuals collected in Hong Kong population.

assessed microRNAs, 26 show a higher expression level/dominant arm usage for the conserved arm (Fig. 4b, Supplementary Data 2, 6). For example, the 3p arm shows more sequence conservation between *bantam* paralogues in horseshoe crabs, and they also show higher expression levels than corresponding 5p arms (Fig. 4b). The 26 conserved microRNAs identified as showing higher expression levels for the conserved arm serve as the first example correlating expression level and conservation of mature microRNA sequences in paralogues following WGD.

In addition to relatively ancient conserved microRNAs, we also investigated novel microRNAs specific to a certain horseshoe crab species, to understand the impact of WGD on these. We found that 7 xiphosuran novel miRNAs are conserved in all three horseshoe crab genomes, while another 13 miRNAs are conserved in both *C. rotundicauda* and *T. tridentatus* (Supplementary Figs. 23–25). Considering that these xiphosuran miRNAs paralogues are located on different scaffolds in their respective genomes, they are unlikely to be generated via tandem duplications. The identified novel microRNAs show higher sequence conservation between orthologues than paralogues (e.g., Supplementary Figs. 23–25), suggesting these horseshoe crab-specific novel microRNAs originate in the horseshoe crab ancestor following WGD.

In the common house spider *Parasteatoda tepidariorum*, which is believed to have undergone a single round of WGD¹⁰, paralogues of microRNAs were found to exhibit arm switching, a phenomenon whereby dominant microRNA arm usage is swapped among different tissues, developmental stages or species^{26,27}. We investigated microRNA arm switching in the sRNA transcriptomes generated here and compared this to their orthologues in various arthropods including fruitfly (*Drosophila melanogaster*), mosquito (*Aedes aegypti*), butterfly (*Heliconius melpomene*), beetle (*Tribolium castaneum*), water flea (*Daphnia pulex*), and tick (*Ixodes scapulari*)^{28,29}. By comparing dominant arm usage across different species, we found that many microRNAs, such as *miR-2788*, *miR-281* and *miR-iab-8* have

undergone microRNA arm switching (Fig. 4c, Supplementary Data 2). Moreover, we also observed microRNA arm switching in cases of microRNAs throughout different developmental time or tissues (Supplementary Data 2). These findings are congruent with the spider microRNA study^{10,27}.

WGD at the population level. An additional question that remains poorly explored is the evolutionary consequences of WGDs on gene duplicates at the population level. To address this question, individuals of both *C. rotundicauda* and *T. tridentatus* were collected from different locations across Asia and subjected to genome sequencing (Fig. 5a, Supplementary Tables 9 and 10). We first mapped sequencing reads to the mitochondrial genome and constructed evolutionary trees from mitochondrial data to identify population structure. Distinct subpopulations were identified within different regions in Asia, for example, populations from Hong Kong formed a distinct group from other locations in Asia, which may be due to the strong ocean currents surrounding Hong Kong preventing gene flow between these locations (Fig. 5b, Supplementary Fig. 26).

Next, we sought to identify how dynamic mutations observed at paralogous genes are in different individuals. We focussed on homeobox genes, calling single-nucleotide polymorphisms (SNPs) at the homeodomains of all annotated homeobox genes and found confident cases of both non-synonymous substitutions as well as pseudogenisation in the homeodomain of certain populations (Fig. 5c; Supplementary Data 7). In *T. tridentatus*, non-synonymous substitutions at the homeodomain of *Six3/6-like* and *Oncut-E* genes were revealed in individuals from Malaysian populations (Fig. 5c, Supplementary Fig. 27a). Similarly, for *C. roundiculata*, non-synonymous substitutions at the homeodomain of the *En-D* gene were also revealed in individuals from populations in Thailand (Supplementary Fig. 27b). This is the first evidence demonstrating that different gene duplicates are under different rates of mutation and selection at the individual level after WGD in invertebrates.

Importantly, unique pseudogenisation was discovered in the paralogue of *Unpg* in many individuals in the *C. rotundicauda* population located in Hong Kong (Fig. 5d). In 9 out of the 10 individuals sequenced from Hong Kong, an alternative form with a deletion in *Unpg-A1* was identified (Fig. 5d). Given that homeodomains are standardised as transcription factors with a sequence length of ~60–63 amino acids²⁰, the observed deletion suggests that in certain individuals, these genes are in the process of becoming pseudogenes. This is the first evidence demonstrating the ongoing and dynamic mutation rate of paralogues at the population level after WGD in invertebrates.

Conclusion. Despite its importance in evolution, the impacts of WGD remain understudied, particularly in invertebrates such as horseshoe crabs. This study provides evidence of the 3R WGD events in horseshoe crabs, shedding light on the evolutionary fates of genes and microRNAs at both the individual and population levels, as well as highlighting the genetic diversity of these amazing animals, with importance for understanding their evolution, genomics, and practical value for breeding programmes and conservation.

Methods

DNA, mRNA, and sRNA extraction and sequencing. Genomic DNA of the horseshoe crabs *C. rotundicauda* and *T. tridentatus* was isolated from the leg muscle of a single individual in each case, using the PureLink Genomic DNA Kit (Invitrogen). In addition, different tissues were dissected and homogenised in Trizol reagent (Invitrogen), and total RNA was isolated following the manufacturers' instructions. Blood samples of both species of horseshoe crab were drawn by syringe and directly transferred into Trizol reagent for RNA extraction. For egg, 1st, 2nd and 3rd instars of *T. tridentatus*, whole individuals were used for RNA extraction. Extracted gDNA was subject to quality control using gel electrophoresis. Qualified samples were sent to Novogene and Dovetail Genomics for library preparation and sequencing. In addition, a Chicago library was prepared by Dovetail Genomics using the method described by Putnam et al.³⁰. Briefly, ~500 ng of high molecular weight gDNA (mean fragment length = 55 kb) was reconstituted into artificial chromatin *in vitro* and fixed with formaldehyde. Fixed chromatin was digested with DpnII, the 5' overhangs filled in with biotinylated nucleotides, and free blunt ends were ligated. After ligation, crosslinks were reversed, and the DNA purified. Purified DNA was treated to remove biotin that was not internal to ligated fragments. The DNA was then sheared to ~350 bp mean fragment size and sequencing libraries were generated using NEBNext Ultra enzymes and Illumina-compatible adaptors. Biotin-containing fragments were isolated using streptavidin beads before PCR enrichment of each library. The libraries were sequenced on the Illumina HiSeq X platform. Dovetail HiC libraries were prepared as described previously³¹. Briefly, for each library, chromatin was fixed with formaldehyde in the nucleus and then extracted. Fixed chromatin was digested with DpnII, the 5' overhangs filled in with biotinylated nucleotides, and free blunt ends were ligated. After ligation, crosslinks were reversed and the DNA purified. Purified DNA was treated to remove biotin that was not internal to ligated fragments. The DNA was then sheared to ~350 bp mean fragment size and sequencing libraries were generated using NEBNext Ultra enzymes and Illumina-compatible adaptors. Biotin-containing fragments were isolated using streptavidin beads before PCR enrichment of each library. Details of the sequencing data can be found in Supplementary Tables 1 and 2.

Total RNA was subject to quality control using a Nanodrop spectrophotometer (Thermo Scientific), gel electrophoresis and analysis using the Agilent 2100 Bioanalyzer (Agilent RNA 6000 Nano Kit). High quality samples underwent library construction and sequencing at Novogene; polyA-selected RNA-sequencing libraries were prepared using TruSeq RNA Sample Prep Kit v2. Insert sizes and the concentration of final libraries were determined using an Agilent 2100 bioanalyzer instrument (Agilent DNA 1000 Reagents) and real-time quantitative PCR (TaqMan Probe), respectively. Small RNA (<200 nt) was isolated using the mirVana miRNA isolation kit (Ambion) according to the manufacturer's instructions. Small RNA was dissolved in the elution buffer provided in the mirVana miRNA isolation kit (Thermo Fisher Scientific) and submitted to Novogene for HiSeq small RNA library construction and 50 bp single-end sequencing. Detailed information for the sequencing data can be found in Supplementary Tables 7 and 8.

Genome, mRNA transcriptome, and sRNA assembly and annotation. To process the Illumina sequencing data, adaptors were trimmed and reads were filtered using the following parameters '-n 0.1 (i.e. removal if N accounted for 10% or more of reads) -1 4 -q 0.5 (i.e., removal if the quality value is lower than 4 and accounts for 50% or more of reads)'. FastQC was run for quality control³². If adaptor

contamination was identified, adaptor sequences were removed using minion³³. Adapter trimming and quality trimming was then performed with cutadapt v1.10³⁴. For each species, k-mers of the Illumina PE library of 500 bp insert size were counted using DSK version 2.1.0 with $k = 25$ ³⁵, and estimation of genome size, repeat content, and heterozygosity were analysed based on a k-mer-based statistical approach using the GenomeScope webtool³⁶. Kraken was used to estimate the percentage of reads that may result from contamination from bacteria³⁷. Chromium WGS reads were separately used to make a *de novo* assembly using Supernova (v 2.1.1), with the parameter '-maxreads = 23,154,066' for *C. rotundicauda*, and '-maxreads = 100,000,000' for *T. tridentatus*, respectively. The *de novo* assembly, shotgun reads, Chicago library reads, and Dovetail HiC library reads were used as input data for HiRise, a software pipeline designed for using proximity ligation data to scaffold genome assemblies³⁰. An iterative analysis was conducted. First, Shotgun and Chicago library sequences were aligned to the draft input assembly using a modified SNAP read mapper (<http://snap.cs.berkeley.edu>). The separation of Chicago read pairs mapped within draft scaffolds was analysed by HiRise to produce a likelihood model for genomic distance between read pairs, and the model was used to identify and break putative misjoins, to score prospective joins, and to make joins above a threshold. After aligning and scaffolding Chicago data, Dovetail HiC library sequences were aligned and scaffolded following the same method. After scaffolding, shotgun sequences were used to close gaps between contigs.

Raw sequencing reads of the transcriptomes were pre-processed with quality trimmed by trimmomatic (version 0.33, with parameters 'ILLUMINACLIP: TruSeq3-PE.fa:2:30:10 SLIDINGWINDOW:4:5 LEADING:5 TRAILING:5 MINLEN:25')³⁸. For the nuclear genomes, the genome sequences were cleaned and masked by Funannotate (v1.6.0, <https://github.com/nextgenusfs/funannotate>)³⁹, the softmasked assembly were used to run 'funannotate train' with parameters '-stranded RF-max_intronlen 350,000' to align RNA-seq data, ran Trinity, and then ran PASA⁴⁰. The PASA gene models were used to train Augustus in 'funannotate predict' step following manufacturers recommended options for eukaryotic genomes (<https://funannotate.readthedocs.io/en/latest/tutorials.html#non-fungal-genomes-higher-eukaryotes>). Briefly, the gene models were predicted by funannotate predict with parameters '-repeats2evm-protein_evidence uniprot_sprot.fasta-genemark_mode ET-busco_seed_species arthropoda-optimize_augustus-busco_db arthropoda-organism other-max_intronlen 350000', the gene models predicted by several prediction sources including GeneMark⁴¹, high-quality Augustus predictions (HiQ), PASA⁴⁰, Augustus⁴², GlimmerHMM⁴³ and snap⁴⁴ were passed to Evidence Modeler⁴⁰ (EVM Weights: {'GeneMark': 1, 'HiQ': 2, 'pasa': 6, 'proteins': 1, 'Augustus': 1, 'GlimmerHMM': 1, 'snap': 1, 'transcripts': 1}) and generated the final annotation files, and then used of PASA⁴⁰ to update the EVM consensus predictions, added UTR annotations and models for alternatively spliced isoforms. The protein-coding genes which cannot hit to nr db by DIAMOND blastp (version v0.9.22.123)⁴⁵ with e value 1e-5 were removed.

To process small RNA data, we removed small RNA sequencing raw reads with Phred quality score less than 20, and adaptor sequences were trimmed. Processed reads of length 18–27 bp were then mapped to their respective horseshoe crab genome and analysed using the mirDeep2 package⁴⁶. To identify conserved microRNAs, the predicted horseshoe crab microRNA hairpins were compared against metazoan microRNA precursor sequences from miRBase²⁸ using BLASTn (e value 0.01)⁴⁷. Predicted microRNAs were manually examined. Novel microRNAs were defined only when they fulfilled the unique features of microRNAs (MirGeneDB 2.0 <https://mirgenedb.org/information>)²⁹. In addition, the copy number of microRNA loci was examined by using microRNA hairpins confirmed above to BLAST against each horseshoe crab genome.

Annotation of repetitive elements. Repetitive elements were identified using an in-house pipeline. First, elements were identified using RepeatMasker ver. 4.0.8⁴⁸ with the *Arthropoda* RepBase⁴⁹ repeat library. Low-complexity repeats were ignored (-nolow) and a sensitive (-s) search was performed. Following this, a *de novo* repeat library was constructed using RepeatModeler ver. 1.0.11⁵⁰, including RECON ver. 1.08⁵¹ and RepeatScout ver. 1.0.5⁵². Novel repeats identified by RepeatModeler were analysed with a 'BLAST, Extract, Extend' process to characterise elements along their entire length⁵³. Consensus sequences and classification information for each repeat family were generated. The resulting *de novo* repeat library was utilised to identify repetitive elements using RepeatMasker. Repetitive element association with genomic features were determined using BedTools ver. 2.26.0⁵⁴. 'Genic' repetitive elements were defined as those overlapping loci annotated as genes \pm 2 kb and identified using the BedTools window function. All plots were generated using Rstudio ver. 1.2.1335 with R ver. 3.5.1⁵⁵ and ggplot2 ver. 3.2.1⁵⁶.

Annotation of gene families and phylogenetic analyses. Potential gene family sequences were first retrieved from the two genomes using tBLASTn⁴⁷. Identity of each putatively identified gene was then tested by comparison to sequences in the NCBI nr database using BLASTx. For homeobox gene retrieval, sequences were also analysed using the BLAST function in HomeoDB. For phylogenetic analyses of gene families, DNA sequences were translated into amino acid sequences and aligned to other members of the gene family; gapped sites were removed from

alignments and phylogenetic trees were constructed using MEGA. Homeobox genes in the *Limulus* genome assembly (GCF_000517525) was obtained from previous study⁸.

Synteny analyses. Synteny blocks were computed using SyMAP v4.2 (Synteny Mapping and Analysis Program) with default parameters except Min Dots from 2 to 7 (Minimum number of anchors required to define a syntenic block = 2–7) and ‘mask_all_but_genes = 1’ to mask non-genic sequence⁵⁷.

Population genomic analyses. After quality control using FastQC³², adaptors and low-quality bases were removed from the read ends using FASTP⁵⁸ with ‘-qualified_quality_phred 30-length_required 25’ and other default parameters, followed by a second round of quality control using FastQC. The trimmed reads were mapped to the unmasked mitochondrion genome (NC_012574 of *T. tridentatus* and NC_019623 of *C. rotundicauda*) using bwa (version 0.7.12-r1039) with default parameters. The mapped reads were sorted using SortSam of picard, and duplicated reads were removed using MarkDuplicates of picard. HaplotypeCaller from the Genome Analysis Toolkit GATK (version 4, <https://gatk.broadinstitute.org/hc/en-us>) was used to estimate the general variant calling file for each individual, and then combined by GenotypeGVCFs to a single variant calling file. Hard filtering of the SNP calls was carried out with Fisher strand bias (FS > 60.0), mapping quality MQ < 40.0, and thresholding by sequencing coverage based on minimum coverage (DP < 100) and maximum coverage (DP > 1500). The SNPs were annotated with SnpEff (version 4.3T, <http://snpeff.sourceforge.net/index.html>)⁵⁹.

Filtered SNPs were used to generate population tree. The model-based software programme STRUCTURE Version 2.3.4. 81 was used for population analysis. To determine most appropriate *k* value, burn-in Markov Chain Monte Carlo (MCMC) replication was set to 50,000 and data were collected over 100,000 MCMC replications in each run. Two independent runs were performed setting the number of population (*k*) from 2 to 10 using a model allowing for admixture and correlated allele frequencies. The basis of this kind of clustering method is the allocation of individual samples to *k* clusters. The *k* value was determined based on the rate of change in Ln(P(D)) between successive *k*, stability of grouping pattern across two run and sample information about the material in Supplementary Data 7, Supplementary Tables 9 and 10. Evolutionary divergence of within and between four different location horseshoe crab samples was performed using MEGA 7 (Molecular Evolutionary genetic analysis)⁶⁰ following maximum composite likelihood model with 1000 bootstrap iterations of all samples. Principal coordinate analysis (PCoA) and UPGMA phylogenetic analysis was conducted to further assess the population subdivisions. PCoA was performed based on distance matrix using DARwin V.6.0.21 and UPGMA tree was constructed based on the simple matching dissimilarity (DARwin).

Trimmed reads were mapped to the homeobox sequences using bwa (version 0.7.12-r1039) with default parameters. The mapped reads were sorted using SortSam of picard, and duplicated reads were removed using MarkDuplicates of picard. HaplotypeCaller from the Genome Analysis Toolkit GATK (version 4, <https://gatk.broadinstitute.org/hc/en-us>) was used to estimate the general variant calling file for each individual, and then combined by GenotypeGVCFs to a single variant calling file. Hard filtering of the SNP calls was carried out with Fisher strand bias (FS > 60.0), mapping quality (MQ < 40.0), QualByDepth (QD < 2.0), MappingQualityRankSumTest (MQRankSum < -12.5), ReadPosRankSumTest (ReadPosRankSum < -8.0) as <https://gatkforums.broadinstitute.org/gatk/discussion/2806/howto-apply-hard-filters-to-a-call-set>. The filtered out SNPs were then annotated with SnpEff (version 4.3 T, <http://snpeff.sourceforge.net/index.html>)⁵⁹. The missense mutation of the homeobox domain were manually checked with samtools tview.

MicroRNA arm switching detection. The expression levels of 5p and 3p arms of microRNAs in the horseshoe crabs were calculated based on the number of sequencing reads mapped to the respective arm region in the predicted microRNA hairpin using bowtie/mirDeep2. The expression of different arms of microRNAs from different species were mapped according to the previous method⁶¹ or referred to the data from MirGeneDB 2.0²⁹. The arm usage ratio (AUR) of each microRNA was calculated using the formula $AUR = 5p/(5p + 3p)$, where 5p and 3p refer to the read counts of predicted 5p and 3p arms, respectively. The AUR ranged from 0 to 1, with smaller values indicating the tendency of 3p preference and larger values indicating the tendency of 5p preference. 5p and 3p dominance was defined where $AUR > 0.7$ and < 0.3 , respectively. No arm preference was defined when AUR ranged from 0.3 to 0.7. The overall arm preference (OAP) of each horseshoe crab microRNA was defined by evaluating their arm dominance in multiple tissue samples. If more than 70% of all tissue samples showed one type of arm dominance, then this type of arm dominance was defined as the OAP of this microRNA. Otherwise, no OAP was defined.

Statistics and reproducibility. Sample size is outlined in Supplementary Data 7. All analyses are reproducible with access to genetic data (see ‘Data availability’).

Reporting summary. Further information on research design is available in the Nature Research Reporting Summary linked to this article.

Data availability

The final genome assemblies have been deposited on NCBI with accession numbers WCHO00000000 and WCHN00000000. The raw reads generated in this study have been deposited to the NCBI database under the BioProject accession no. PRJNA574021 and PRJNA574023, for *C. rotundicauda* and *T. tridentatus* respectively. The genome annotation files are deposited in Figshare <https://doi.org/10.6084/m9.figshare.13172414>⁶². All other data, if any, are available upon reasonable request.

Received: 22 April 2020; Accepted: 21 December 2020;

Published online: 19 January 2021

References

- Holland, P. W. H. More genes in vertebrates? *J. Struct. Func. Genom.* **3**, 75–84 (2003).
- Van de Peer, Y., Mizrahi, E. & Marchal, K. The evolutionary significance of polyploidy. *Nat. Rev. Genet.* **18**, 411–424 (2017).
- Sémon, M. & Wolfe, K. H. Reciprocal gene loss between Tetraodon and zebrafish after whole genome duplication in their ancestor. *Trends Genet.* **23**, 108–112 (2007).
- Jailon, O., Aury, J. M. & Wincker, P. “Changing by doubling”, the impact of Whole Genome Duplications in the evolution of eukaryotes. *Comptes Rendus. Biol.* **332**, 241–253 (2009).
- Van de Peer, Y., Maere, S. & Meyer, A. The evolutionary significance of ancient genome duplications. *Nat. Rev. Genet.* **10**, 725–732 (2009).
- Moriyama, Y. & Koshiba-Takeuchi, K. Significance of whole-genome duplications on the emergence of evolutionary novelties. *Brief. Funct. Genom.* **17**, 329–338 (2018).
- Nossa, C. et al. Joint assembly and genetic mapping of the Atlantic horseshoe crab genome reveals ancient whole genome duplication. *GigaScience* **3**, 9 (2014).
- Battelle, B.-A. et al. Opsin repertoire and expression patterns in horseshoe crabs: evidence from the genome of *Limulus polyphemus* (Arthropoda: Chelicerata). *Genome Biol. Evol.* **8**, 1571–1589 (2016).
- Kenny, N. J. et al. Ancestral whole-genome duplication in the marine chelicerate horseshoe crabs. *Heredity* **116**, 190–199 (2016).
- Schwager, E. E. et al. The house spider genome reveals an ancient whole-genome duplication during arachnid evolution. *BMC Biol.* **15**, 62 (2017).
- Rudkin, D. M. & Young, G. A. Horseshoe crabs—an ancient ancestry revealed. In *Biology and Conservation of Horseshoe Crabs* (eds. Tanacredi, J. T., Botton, M. L. & Smith, D.) 25–44 (Springer, US, 2013). <https://www.springer.com/gp/book/9780387899589>.
- John, A. B. et al. A review on fisheries and conservation status of Asian horseshoe crabs. *Biodivers. Conserv.* **29**, 3573–3598 (2018).
- Obst, M., Faurby, S., Bussarawit, S. & Funch, P. Molecular phylogeny of extant horseshoe crabs (Xiphosura, Limulidae) indicates Paleogene diversification of Asian species. *Mol. Phylogenet. Evol.* **62**, 21–26 (2012).
- Gong, L. et al. Chromosomal level reference genome of *Tachypleus tridentatus* provides insights into evolution and adaptation of horseshoe crabs. *Mol. Ecol. Resour.* **19**, 744–756 (2019).
- Liao, Y. Y. et al. Draft genomic and transcriptome resources for marine chelicerate *Tachypleus tridentatus*. *Sci. Data* **6**, 190029 (2019).
- Iwasaki, Y., Iwami, T. & Sekiguchi, K. Karyology. In *Biology of Horseshoe Crabs* (ed. Sekiguchi, K.) 309–314 (Science House, Inc., Tokyo, 1988).
- Chénais, B., Caruso, A., Hiard, S. & Casse, N. The impact of transposable elements on eukaryotic genomes: From genome size increase to genetic adaptation to stressful environments. *Gene* **509**, 7–15 (2012).
- Holland, P. W. H. The dawn of amphioxus molecular biology—a personal perspective. *Int. J. Dev. Biol.* **61**, 585–590 (2017).
- Putnam, N. H. et al. The amphioxus genome and the evolution of the chordate karyotype. *Nature* **453**, 1064–1071 (2008).
- Holland, P. W. H. Evolution of homeobox genes. *Rev. Dev. Biol.* **2**, 31–45 (2013).
- Brooke, N. M., Garcia-Fernández, J. & Holland, P. W. H. The ParaHox gene cluster is an evolutionary sister of the Hox gene cluster. *Nature* **392**, 920–922 (1998).
- Hui, J. H. et al. Features of the ancestral bilaterian inferred from *Platymereis dumerilii* ParaHox genes. *BMC Biol.* **7**, 43 (2009).
- Hui, J. H. et al. Extensive chordate and annelid macrosynteny reveals ancestral homeobox gene organization. *Mol. Biol. Evol.* **29**, 157–165 (2012).
- Ballesteros, J. A. & Sharma, P. P. A critical appraisal of the placement of Xiphosura (Chelicerata) with account of known sources of phylogenetic error. *Syst. Biol.* **68**, 896–917 (2019).

25. Lozano-Fernandez, J. et al. Increasing species sampling in chelicerate genomic-scale datasets provides support for monophyly of Acari and Arachnida. *Nat. Commun.* **10**, 2295 (2019).
26. Griffiths-Jones, S., Hui, J. H. L., Marco, A. & Ronshaugen, M. MicroRNA evolution by arm switching. *EMBO Rep.* **12**, 172–177 (2011).
27. Leite, D. J. et al. Pervasive microRNA duplication in chelicerates: insights from the embryonic microRNA repertoire of the spider *Parasteatoda tepidariorum*. *Genome Biol. Evol.* **8**, 2133–2144 (2016).
28. Kozomara, A. & Griffiths-Jones, S. miRBase: annotating high confidence microRNAs using deep sequencing data. *Nucleic Acids Res.* **42**, D68–D73 (2014).
29. Fromm, B. et al. MirGeneDB 2.0: the metazoan microRNA complement. *Nucleic Acids Res.* **48**, D132–D141 (2019).
30. Putnam, N. H. et al. Chromosome-scale shotgun assembly using an in vitro method for long-range linkage. *Genome Res.* **26**, 342–350 (2016).
31. Lieberman-Aiden, E. et al. Comprehensive mapping of long-range interactions reveals folding principles of the human genome. *Science* **326**, 289–293 (2009).
32. Andrews, S. FastQC: a quality control tool for high throughput sequence data. <http://www.bioinformatics.babraham.ac.uk/projects/fastqc> (2010).
33. Davis, M. P. A., vanDongen, S., Abreu-Goodger, C., Bartonicek, N. & Enright, A. J. Kraken: a set of tools for quality control and analysis of high-throughput sequence data. *Methods* **63**, 41–49 (2013).
34. Martin, M. Cutadapt removes adapter sequences from high-throughput sequencing reads. *EMBnet* **17**, 3 (2011).
35. Rizk, G., Lavenier, D. & Chikhi, R. DSK: K-mer counting with very low memory usage. *Bioinformatics* **29**, 652–653 (2013).
36. Vurture, G. W. et al. GenomeScope: fast reference-free genome profiling from short reads. *Bioinformatics* **33**, 2202–2204 (2017).
37. Wood, D. E. & Salzberg, S. L. Kraken: Ultrafast metagenomic sequence classification using exact alignments. *Genome Biol.* **15**, R46 (2014).
38. Bolger, A. M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for Illumina sequence data. *Bioinformatics* **30**, 2114–2120 (2014).
39. Palmer, J. & Stajich, J. Funannotate: eukaryotic genome annotation pipeline. <http://funannotate.readthedocs.io/> (2018).
40. Haas, B. J. et al. Automated eukaryotic gene structure annotation using EVIDENCEModeler and the program to assemble spliced alignments. *Genome Biol.* **9**, R7 (2008).
41. Lomsadze, A., Ter-Hovhannisyan, V., Chernoff, Y. O. & Borodovsky, M. Gene identification in novel eukaryotic genomes by self-training algorithm. *Nucleic Acids Res.* **33**, 6494–6506 (2005).
42. Stanke, M. et al. AUGUSTUS: ab initio prediction of alternative transcripts. *Nucleic Acids Res.* **34**, W435–W439 (2006).
43. Majoros, W. H., Pertea, M. & Salzberg, S. L. TigrScan and GlimmerHMM: two open source ab initio eukaryotic gene-finders. *Bioinformatics* **20**, 2878–2879 (2004).
44. Korf, I. Gene finding in novel genomes. *BMC Bioinform.* **5**, 59 (2004).
45. Buchfink, B., Xie, C. & Huson, D. H. Fast and sensitive protein alignment using DIAMOND. *Nat. Methods* **12**, 59–60 (2014).
46. Friedländer, M. R., Mackowiak, S. D., Li, N., Chen, W. & Rajewsky, N. miRDeep2 accurately identifies known and hundreds of novel microRNA genes in seven animal clades. *Nucleic Acids Res.* **40**, 37–52 (2012).
47. Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local alignment search tool. *J. Mol. Biol.* **215**, 403–410 (1990).
48. Smit, A. F. A., Hubley, R. R. & Green, P. R. RepeatMasker Open-4.0. <http://repeatmasker.org> (2013).
49. Jurka, J. et al. Repbase Update, a database of eukaryotic repetitive elements. *Genome Res.* **110**, 462–467 (2005).
50. Smit, A. F. A. & Hubley, R. RepeatModeler Open-1.0. <http://repeatmasker.org> (2015).
51. Bao, Z. & Eddy, S. R. Automated *de novo* identification of repeat sequence families in sequenced genomes. *Genome Res.* **12**, 1269–1276 (2002).
52. Price, A. L., Jones, N. C. & Pevzner, P. A. *De novo* identification of repeat families in large genomes. *Bioinformatics* **21**, 351–358 (2015).
53. Platt, R. N., Blanco-Berdugo, L. & Ray, D. A. Accurate transposable element annotation is vital when analyzing new genome assemblies. *Genome Biol. Evol.* **8**, 403–410 (2016).
54. Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics* **26**, 841–842 (2010).
55. Team, R. C. R. A language and environment for statistical computing. (ISBN 3-900051-07-0. R Foundation for Statistical Computing. Vienna, Austria, 2013. <http://www.R-project.org>, 2005).
56. Wickham, H. *ggplot2: Elegant Graphics for Data Analysis*. (Springer, 2016).
57. Soderlund, C., Bomhoff, M. & Nelson, W. M. SyMAP v3.4: a turnkey synteny system with application to plant genomes. *Nucleic Acids Res.* **39**, e68 (2011).
58. Chen, S., Zhou, Y., Chen, Y. & Gu, J. fastp: an ultra-fast all-in-one FASTQ preprocessor. *Bioinformatics* **34**, i884–i890 (2018).
59. Cingolani, P. et al. A program for annotating and predicting the effects of single nucleotide polymorphisms, SnpEff: SNPs in the genome of *Drosophila melanogaster* strain *w1118*. *Fly* **6**, 80–92 (2012).
60. Kumar, S., Stecher, G. & Tamura, K. MEGA7: molecular evolutionary genetics analysis version 7.0 for bigger datasets. *Mol. Biol. Evol.* **33**, 1870–1874 (2016).
61. Marco, A., Hui, J. H., Ronshaugen, M. & Griffiths-Jones, S. Functional shifts in insect microRNA evolution. *Genome Biol. Evol.* **2**, 686–696. <https://doi.org/10.1093/gbe/evq053> (2010).
62. Nong, W. Horseshoe crab genomes reveal the evolution of genes and microRNAs after three rounds of whole genome duplication. Figshare. <https://doi.org/10.6084/m9.figshare.13172414.v2> (2020).
63. Shingate, P. et al. Chromosome-level assembly of the horseshoe crab genome provides insights into its genome evolution. *Nat Commun.* **11**, 2322 (2020).

Acknowledgements

The authors thank Peter Holland for discussion, B. Desany for helping with the assembly of 10X Genomics data, and F. Cheung, R. Leung, W. Tong, W. Yiu, and H. Yu for collection of some of the RNA. This research was supported by the Hong Kong Research Grant Council GRF Grant 14103516 and 14100919; Environment and Conservation Fund Project 28/2017; Agriculture, Fisheries and Conservation Department of HKSAR Government, and the School of Life Sciences of The Chinese University of Hong Kong (JHLH). P.Y.Q. and J.W.Q. are supported by Hong Kong Branch of Southern Marine Science and Engineering Guangdong Laboratory (Guangzhou) (SMSEGL20SC01). A.H. is supported by a Biotechnology and Biological Sciences Research Council (BBSRC) David Phillips Fellowship (BB/N020146/1). T.B. is supported by a studentship from the Biotechnology and Biological Sciences Research Council-funded South West Biosciences Doctoral Training Partnership (BB/M009122/1).

Author contributions

J.H.L.H. conceived and supervised the study. W.N. carried out the genome assemblies and analyses, gene model predictions, microRNA mapping, and synteny analyses. Z.Q. carried out the microRNA annotation, final checking microRNA copies and arm switching analyses. Y.L. carried out the homeobox gene analyses, synteny analyses, and the SNP analyses in populations. T.B.O. carried out the homeobox gene identifications and tree construction. A.Y.P.W. carried out the gene and microRNA copies analyses and arm switching analyses. H.Y.Y. provided animal husbandry and logistics. H.T.L. carried out novel microRNA analyses. S.N. carried out the population structure analyses. T.B. and A.H. performed the T.E. analyses. T.S. involved in the final version of genome assembly, and J.C. involved in earlier version of genome assembly. T.F.C., H.S.K., S.M.N., G.P., P.Y.Q., J.W.Q., K.Y.Y., S.S.T., W.G.B., S.G.C., J.H.L.H. applied and obtained the funding. N.I., S.P., A.J., S.G.C. collected and provided samples in field. S.S.T., W.G.C., S.G.C., A.H., J.H.L.H. drafted the first version of the manuscript. All authors provided comments and approved the manuscript.

Competing interests

The authors declare no competing interests.

Additional information

Supplementary information The online version contains supplementary material available at <https://doi.org/10.1038/s42003-020-01637-2>.

Correspondence and requests for materials should be addressed to J.H.L.H.

Reprints and permission information is available at <http://www.nature.com/reprints>

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2021

¹School of Life Sciences, Simon F.S. Li Marine Science Laboratory, State Key Laboratory of Agrobiotechnology, The Chinese University of Hong Kong, Hong Kong, China. ²Centre for Ecology and Conservation, University of Exeter, Penryn, UK. ³Dovetail Genomics, Scotts Valley, CA, USA. ⁴State Key Laboratory of Agrobiotechnology, School of Life Sciences, The Chinese University of Hong Kong, Hong Kong, China. ⁵School of Life Sciences, The Chinese University of Hong Kong, Hong Kong, China. ⁶School of Biological Sciences, The University of Hong Kong, Hong Kong, China. ⁷Leibniz Institute of Natural Product Research and Infection Biology – Hans Knöll Institute, Jena, Germany. ⁸Department of Ocean Science and Hong Kong Branch of Southern Marine Science and Engineering Guangdong Laboratory (Guangzhou), Hong Kong University of Science and Technology, Hong Kong, China. ⁹Department of Biology, Hong Kong Baptist University, Hong Kong, China. ¹⁰Department of Computer Science and Engineering, The Chinese University of Hong Kong, Hong Kong, China. ¹¹Institute of Marine Biotechnology, Universiti Malaysia Terengganu, Terengganu, Malaysia. ¹²Department of Bioscience and Biotechnology, Fakir Mohan University, Balasore, India. ¹³Institute of Tropical Biodiversity and Sustainable Development, University Malaysia Terengganu, 20130 Kuala Nerus, Terengganu, Malaysia. ¹⁴Research Division, Association for Biodiversity Conservation and Research (ABC), Odisha 756003, India. ¹⁵Institute of Oceanography and Maritime Studies (INOCEM), Kulliyah of Science, International Islamic University, Kuantan, Malaysia. ¹⁶Department of Cell and Systems Biology, University of Toronto, Toronto, Canada. ¹⁷Department of Biology, Queen’s University, Toronto, Canada. ¹⁸Department of Chemistry, City University of Hong Kong, Hong Kong, China. ¹⁹These authors contributed equally: Wenyan Nong, Zhe Qu, Yiqian Li, Tom Barton-Owen, Annette Y. P. Wong. ✉email: jeromehui@cuhk.edu.hk