

Hagfish and lamprey Hox genes reveal conservation of temporal colinearity in vertebrates

Juan Pascual-Anaya^{1*}, Iori Sato^{1†}, Fumiaki Sugahara^{1,2†}, Shinnosuke Higuchi¹, Jordi Paps³, Ren Yandong^{4,5}, Wataru Takagi¹, Adrián Ruiz-Villalba⁶, Kinya G. Ota⁷, Wen Wang⁴, Shigeru Kuratani¹

¹ *Evolutionary Morphology Laboratory, RIKEN, Kobe 650-0047, Japan*

² *Division of Biology, Hyogo College of Medicine, Nishinomiya 663-8501, Japan*

³ *School of Biological Sciences, University of Essex, Colchester CO4 3SQ, UK.*

⁴ *State Key Laboratory of Genetic Resources and Evolution, Kunming Institute of Zoology, Chinese Academy of Sciences, Kunming 650223, China.*

⁵ *University of Chinese Academy of Sciences, Beijing, 100049, China*

⁶ *Cell Therapy Program, Foundation of Applied Medical Research (FIMA), University of Navarra, Pamplona, Spain.*

⁷ *Laboratory of Aquatic Zoology, Marine Research Station, Institute of Cellular and Organismic Biology, Academia Sinica, Yilan 26242, Taiwan*

[†]Iori Sato and Fumiaki Sugahara contributed equally to this work.

*Correspondence to:

Juan Pascual-Anaya

e-mail: jpascualanaya@gmail.com

Evolutionary Morphology Laboratory, RIKEN, 2-2-3 Minatojima-minami, Chuo-ku, Kobe, Hyogo 650-0047, Japan

Tel.: +81-78-306-3589

Fax: +81-78-306-3370

Hox genes exert fundamental roles for proper regional specification along the main rostro-caudal axis of animal embryos. Hox genes are generally expressed in restricted spatial domains according to their position in the cluster (spatial colinearity), a feature that is conserved across bilaterians. In jawed vertebrates (gnathostomes), the position in the cluster also determines the onset of expression of Hox genes, a feature known as whole-cluster temporal colinearity (WTC), while in invertebrates this phenomenon is displayed as a subcluster-level temporal colinearity (STC). However, little is known about the expression profile of Hox genes in jawless vertebrates (cyclostomes), and therefore the evolutionary origin of WTC, as seen in gnathostomes, remains a mystery. Here we show that Hox genes in cyclostomes are expressed according to WTC during development. We have investigated the Hox repertoire and Hox gene expression profiles in three different species —a hagfish, a lamprey and a shark— encompassing the two major groups of vertebrates and found that these are expressed following a whole-cluster, temporally staggered pattern, indicating that WTC has been conserved during the last 500 million years despite drastically different genome evolution and morphological outputs between jawless and jawed vertebrates.

Hox genes are fundamental developmental genes with crucial roles for the early specification of embryonic structures along the main anterior-posterior axis of bilaterian animals¹. Hox genes are usually placed in the same genomic regions forming cluster(s). Hox clusters are thought to be the result of several tandem duplication events of an ancestral proto-Hox gene², and while most invertebrates generally have a single Hox cluster, vertebrate genomes present multiple clusters³. It is widely accepted that the genome of vertebrates has evolved through two rounds (2R) of whole genome duplication (WGD) events (but see ref. 4 for an alternative scenario), generating up to four paralogous loci for each single region of a pre-duplicative

genome⁵⁻⁸. Extant vertebrates are divided into two major groups: agnathans, represented by the monophyletic group of cyclostomes (hagfish and lampreys); and gnathostomes, encompassing all jawed-vertebrates in two major groups: cartilaginous fishes (e.g., sharks, rays and chimaeras) and bony vertebrates (e.g., teleosts, coelacanth, amphibians, reptiles, mammals). Tetrapod genomes, including mammals, contain four Hox clusters, named from HoxA to HoxD, as the result of these 2R-WGD (Fig. 1a). Although the 2R-WGD events are generally accepted, the timing of these events with respect to the divergence of cyclostomes and gnathostomes is still a matter of intense debate^{4,9-11}. Despite extended research on vertebrate genomes, this has mostly focused on representative species of gnathostomes, while cyclostomes have remained poorly understood. A recent study of the genome of the Artic lamprey, *Lethenteron camtschaticum* (or Japanese lamprey, *Lethenteron japonicum*), suggested that lampreys had probably undergone a third round of WGD event (3R-WGD)¹⁰. Whether this event is an independent, lineage-specific event remains a mystery, since the Hox complement of the hagfish is unknown.

The position of Hox genes in the cluster determines their expression patterns. Spatial colinearity refers to the property by which the anterior limit of expression of a given Hox gene is generally more rostral than its upstream (more 5') counterpart. Spatial colinearity is widely conserved among bilaterians studied so far, even in cases where the Hox cluster is completely atomized¹². Temporal colinearity refers to the phenomenon describing the temporal order of expression of Hox genes according to their position in the cluster, i.e., genes in the 3' part are expressed earlier, and was first described in the HoxD cluster of the mouse^{13,14}. Indeed, this so-called whole-cluster temporal colinearity (WTC)¹⁵ phenomenon had been described only in jawed vertebrates. The recent analysis of the scallop genome and the reanalysis of Hox gene expression in a wide range of invertebrates has revealed that Hox genes of these species follow what is called a subcluster-level temporal colinearity (STC),

i.e., that the cluster is divided into small, contiguous groups of Hox genes, each of these displaying temporal colinearity¹⁵. This situation leads to the uncertainty of what was the ancestral condition before deuterostomes and protostomes split. Moreover, temporal colinearity has not been described in any cyclostome species so far. In *L. camtschaticum*, Hox genes known to be in the same cluster were not expressed following temporal colinearity¹⁶, and both the Hox gene repertoire and expression of the other major group of cyclostomes, the hagfish, is mostly unknown^{3,17}. Therefore, the evolutionary origin of WTC as observed in gnathostomes remains obscure.

Here, we provide a comprehensive analysis of different transcriptomics and genomics resources for the Japanese inshore hagfish, *Eptatretus burgeri*. The hagfish Hox repertoire consists of at least 40 Hox genes, including six Hox4 genes that might suggest the presence of at least 6 Hox clusters, suggesting that the 3R-WGD described for the lamprey could be shared in cyclostomes. Finally, we have comprehensively compared the developmental expression levels of Hox genes during development of four different chordate species, including the hagfish and the lamprey, and conclude that temporal colinearity likely originated in the last common ancestor of chordates, and it was certainly well established at least in the last common ancestor of extant vertebrates.

Results and discussion

To gain insights into vertebrate Hox evolution (Fig. 1a), especially with regards to the evolution of temporal colinearity, we decided to comprehensively analyse the Hox repertoire and expression of Hox genes during development of both the lamprey and the hagfish. First, we screened both the developmental transcriptome and the genome of *E. burgeri*. For the developmental transcriptome, RNA-seq data was generated from three different whole

hagfish embryos at Bashford Dean stages 28/30, 35 and 40/45 (refs. 18 and 19; Fig. 1c-e) and from the head region of a hatched juvenile. In total, we found 40 bona fide Hox genes in the developmental transcriptome of the hagfish, including the 5' and 3' untranslated regions for most cases (Fig. 1b).

To determine the genomic organization of hagfish Hox genes we then screened a BAC library built from blood genomic DNA. We found 25 BAC clones spanning only 15 out of the 40 Hox genes (Supplementary Fig. 1, 2). Recently, it has been described that the lamprey genome goes through somatic rearrangements, differentially eliminating stretches of germ line-specific sequences, which might include protein-coding genes²⁰. Considering that the hagfish, which is known to go through a chromosome elimination process in somatic tissues during development²¹, might be losing Hox genes in somatic tissues, we decided to generate a draft genome using genomic DNA obtained from the testis (germ line) of a single individual. In our preliminary assembly, we found evidence for at least six Hox clusters containing all 40 Hox genes found in the transcriptome and three microRNAs, together with conserved syntenic non-Hox genes (Fig. 1b). The hagfish Hox repertoire and genomic organization are overall very similar to the one described in the *L. camtschaticum* genome¹⁰ (number of genes —43 in the lamprey— and putative clusters —six in the lamprey—), raising the possibility that the 3R-WGD event suggested to have occurred in the lamprey lineage¹⁰ took place before the split of lampreys and hagfish lineages. Surprisingly, we found a hagfish Hox13 gene (*Hox13VI*) enclosed by two conserved syntenic genes: *Lunapark* (*Lnp*) and *Even-skipped* (*Evx*). This suggests that a translocation event took place in the hagfish lineage, likely together with a severe disintegration of a cluster involving large Hox gene losses.

Phylogenetic analysis and best BLAST hits show that the hagfish genome contains representative Hox genes of most of the vertebrate paralogy groups (PG) between PG1 and PG14 (Fig. 2; Supplementary Figs. 3-7). Interestingly, the hagfish genome does not contain

any member of the PG12 (Fig. 1b, Supplementary Figs. 6, 7), a feature shared with the lamprey^{3,10,11} (Fig. 1a). Phylogenetic analysis of the posterior Hox genes suggests that a shared cyclostome loss of the PG12 is the most plausible scenario. We were, however, unable to clarify one-to-one orthology relationships between gnathostome HoxA-D paralogs and lamprey and hagfish Hox genes. Therefore, we named the hagfish Hox genes with a different nomenclature from the one used for the lamprey and gnathostomes counterparts, using roman numbers: I-VI (Fig. 1b).

The obscure orthology relationship between jawed and jawless vertebrate genes has been broadly described for both Hox and non-Hox gene families²². It is unclear whether the 2R-WGD events that took place during early vertebrate evolution are shared or not among cyclostomes and gnathostomes^{6,23,24}. The lack of one-to-one orthology relationships between genes from both groups can be taken as evidence for independent WGD events. However, despite their obscure phylogenetic relationship, Hox clusters of cyclostomes and gnathostomes can still be the result of an ancestral 2R-WGD, if the duplicated regions containing the Hox clusters had not completed rediploidization before the split of cyclostomes and gnathostomes²⁵. Consequently, certain number of phylogenetic analyses would support a shared WGD between cyclostomes and gnathostomes as it seems to be the case⁹. These would correspond to those genes that had differentiated into different alleles before the split of the two lineages.

Once confirmed the presence of clusters, we wondered whether the hagfish Hox genes were expressed according to the spatial colinearity rule. Spatial colinearity has been observed in the vast majority of bilaterians studied so far, included the lamprey^{16,26}, and with only few exceptions²⁷. In both the lamprey and gnathostomes, nested expression of anterior Hox genes is coupled to the morphological segmentation of the hindbrain into discrete rhombomeres, and this is controlled by a highly conserved gene regulatory network, established at least in

149 the last common ancestor of vertebrates²⁸. The hagfish hindbrain is, as in the lamprey²⁹,
 150 transiently segmented into rhombomeres during stage 45 (ref. 19). We investigated the
 151 expression pattern of *E. burgeri* anterior Hox genes in three different developmental stages,
 152 from mid-pharyngula (stage 40 and 45) to late-pharyngula (stage 53; Fig. 3), with especial
 153 focus on their putative expression pattern in the hindbrain. We found that several hagfish
 154 Hox1-5 genes were expressed with staggered anterior boundaries in the hindbrain, an
 155 expression pattern reminiscent of that of the lamprey^{26,28} and gnathostomes³⁰ (Fig. 3y, z). We
 156 also found Hox2-5 genes expressed colinearly in the pharynx of a juvenile at stage 53
 157 (Supplementary Fig. 8). In the hindbrain, the most rostral expression domain detected was
 158 that of *Hox2IV*, at the border between rhombomeres 1 and 2 (r1/2), from stage 40 (Fig3d, l,
 159 t). *Hox2III* signal is not revealed until stage 45, and is similar to that of *Hox2IV*, with its
 160 rostral limit apparently around the lateral edge of the diamond shape of the 4th ventricle,
 161 which in gnathostomes marks the r1/2 border³¹ (Fig. 3k). The expression of Hox2 genes from
 162 r2 rearwards is conserved in all vertebrates (Fig. 3z). In gnathostomes and the lamprey, r4 is
 163 characterized by a strong expression of Hox1. We were able to find only a very faint
 164 expression of only one of the Hox1 genes in the hagfish, *HoxIV*, not in r4 but probably
 165 within r7 with an unclear rostral limit (Fig. 3c). We were not able to find any expression for
 166 *HoxII* and *HoxIII*, which could still be expressed in r4 at different stages. Hagfish *Hox3VI*
 167 was expressed up to r6 (Fig. 3f, n), while, strikingly, *Hox3II* was found to be expressed in
 168 two domains: r5, and from r7 onwards, i.e., with r6 being *Hox3II*-negative (Fig. 3e, m, u). We
 169 also found that *Hox4IV* is expressed, as other vertebrate Hox4 genes, from r7 (Fig. 3h, p, v).
 170 *Hox4I* is expressed later in development, at stage 45, with a very similar pattern to that of
 171 *Hox4IV*, but slightly posteriorly (Fig. 3o). We also found a very weak signal for *Hox4VI* at
 172 stage 45 (Fig. 3q). *Hox5III* is expressed the most posteriorly, apparently from the most
 173 anterior part of the spinal cord at stage 40, its rostral limit shifting anteriorly into the

hindbrain by stage 45 and 53 (Fig. 3i, r, w), when transcripts of *Hox5IV* are also detected (Fig. 3s).

The evolution of the expression domains of Hox3 genes in the hindbrain of different vertebrates is particularly interesting. Considering the global expression pattern of Hox3 paralogs in each group, we observe that while in the lamprey (*Hox3α*) and the shark *Scyliorhinus canicula* (*Hoxb3*), Hox3 genes are expressed from r4 (refs. 26, 30, 32), in the hagfish and osteichthyans Hox3 genes are expressed up to r5 (Fig. 3z). There are two possible evolutionary explanations for this difference, both involving parallel evolutionary events: either a caudal shift of Hox3 expression domains from r4 to r5 convergently happened in both the lamprey and osteichthyan lineages, or, on the other hand, a rostral shift from r5 to r4 occurred in the lamprey and chondrichthyans. In a different lamprey species, *Petromyzon marinus*, the *Pm1Hox3* gene, orthologous of *L. camtschaticum* *Hox3α*, was found to be expressed from r5 like in mammals²⁸. This could favour the hypothesis of a convergent expression shift in both the Arctic lamprey and the shark as lineage or species-specific changes.

Following the spatial colinearity rule, most posterior PG Hox genes are expressed in the most caudal regions of the embryo. One of the expression domains of Hox13 paralogue genes are the most posterior parts of the hindgut. Concordantly, hagfish *Hox13II* and *Hox13VI* were found around the cloacal region of a juvenile (stage 60; Supplementary Fig. 10), as in the lamprey and other vertebrates³³. Vertebrate Hox14 genes have also been reported to be expressed in the most posterior parts of the hindgut of the lamprey and the shark³³. However, we were not able to detect any signal for *Hox14I* transcripts in the cloaca of the hagfish larva (Supplementary Fig.9).

Overall, vertebrate Hox code is generally conserved in hagfish, particularly in the case of the hindbrain, suggesting that the GRN underlying vertebrate hindbrain

segmentation²⁸ is well conserved in the hagfish. More important than the similarities, elucidating what specific regulatory inputs account for lineage-specific differences in the hindbrain Hox code, such as the striped expression of hagfish *Hox3II* and the different rostral limits of expression of different Hox3 genes in different vertebrates, will be helpful to determine how the hindbrain GRN diversified during vertebrate evolution, and what are the functional and morphological implications of these differences.

To unravel the evolutionary origin of WTC in vertebrates, we further carried out a comprehensive analysis of the developmental expression profile of Hox genes using embryos from both jawed and jawless vertebrates. Together with the RNA-seq data generated for *E. burgeri*, we sequenced RNA-seq libraries covering early to late developmental stages of the lamprey *L. camtschaticum*³⁴ and the gnathostome catshark *Scyliorhinus torazame*³⁵ and quantified the expression profiles of Hox genes. As expected, the expression profiles of *S. torazame* Hox genes were consistent with temporal colinearity across all the clusters found in other jawed vertebrates, showing a clear tendency of anterior Hox genes (Hox1-3) expressed at earlier stages and posterior ones (Hox9-14) at later (Fig. 4; Supplementary Fig. 10). Despite previous reports¹⁶, lamprey Hox genes (for which we found an unreported Hox1 gene, *Hox1ζ*) also followed the rule of temporal colinearity (Fig. 4; Supplementary Fig. 10). Interestingly, the HOX-γ cluster has completely lost the temporal colinearity (Supplementary Fig. 10). HOX-γ is one of the most degenerated clusters in the lamprey with only 4 Hox genes¹⁰ (Fig. 1a), which might be a direct consequence of the lack of temporal colinearity. In the hagfish, although obtaining a pool of embryos from a full developmental series is unfeasible, a similar tendency was also observed: levels of posterior *Hox11.I*, *Hox11.V* and all Hox13 genes are higher at stage 40-45 (comparable to lamprey stage 25-26) than at 28-30 (lamprey stage 22-23), while generally all anterior and central Hox genes levels are higher at stage 28-30 than at later stages (Fig. 4).

The above observations imply that cyclostome Hox expression profiles, as in gnathostomes, are consistent with the WTC rule, suggesting at least a vertebrate origin. In order to determine whether WTC was present before the origin of vertebrates, we investigated the Hox expression profiles of a chordate outgroup. Wang and colleagues¹⁵ described the tunicate *Ciona intestinalis* Hox gene expressions as according to the STC. However, their statement was based on the reanalysis of data from whole mount *in situ* hybridization³⁶, which is not a quantitative technique. Cephalochordates are the closest lineage to vertebrates with an intact Hox cluster, and is thus very informative in this regard. Expression profiles of Hox genes in the cephalochordate amphioxus *Branchiostoma belcheri*³⁷ show that amphioxus *Hox1* to *Hox5* are expressed in an anterior Hox/early-posterior/late manner. However, *Hox6*, *Hox10* and *Hox14* genes violate this pattern, consistent with our previous report³⁸, and *Hox7-8*, *Hox11-13* and *Hox15* were not detected during the stages assayed, contributing to the dismantling of the colinearity (Fig. 4). In most invertebrate species where STC has been described, Hox1-2 or Hox1-3 was the most anterior subgroup showing temporal colinearity¹⁵. The fact that in amphioxus Hox1-5 are expressed in temporal order as a single group indicates that this expression pattern is reminiscent from a genuine WTC, which was subsequently broken from Hox6 in the cephalochordate lineage (Fig. 4). In addition, amphioxus Hox6-15 genes might still follow WTC at later stages than the ones assayed here³⁷. The putative presence of WTC in both the cephalochordate and vertebrate lineages implies that it was likely present in the last common ancestor of vertebrates.

Taken together, our results depict a scenario in which chordate Hox genes are expressed following WTC, and protostome Hox genes according to STC. This, importantly, can offer a mechanistic answer to explain the radically different bauplans displayed by chordates and protostome invertebrates. Deschamps and Duboule³⁹ have recently proposed

that temporal colinearity, as seen in mammals (WTC), is displayed only by animals that follow a developmental strategy of anterior to posterior elongation, adding new regions to the main body axis from a posterior growth zone. This temporal activation of Hox genes during the posterior elongation, or Hox clock, translates during development into the spatial colinearity observed along the main anterior-posterior body axis. The fact that the lamprey and the hagfish also develop according to this posterior elongation –a developmental mode thus very well conserved across vertebrates–, together with the presence of WTC and spatial colinearity in the main axial structures of these animals (this study and refs. 16, 26), supports Deschamps and Duboule’s hypothesis³⁹. This implies that this mechanism was present in the last common ancestor of vertebrates, although some lineage-specific differences might have occurred in the mechanism transmitting the Hox clock from the posterior progenitors into the resulting axial structures (for instance, there are differences in the expression of Hox10 genes between lampreys and amniotes in the tailbud and axial mesoderm, see ref. 16). Ultimately, the question of whether the Hox cluster of the last common bilaterian ancestor was expressed according to either whole-cluster or subcluster modes of temporal expression remains open. A more detailed investigation of the temporal expression of Hox genes in non-chordate deuterostome groups (namely, ambulacrarians –e.g., sea urchins, sea stars, acorn worms–) will be thus needed to ultimately resolve this question^{40,41}.

It has been proposed that gnathostome Hox clusters are relatively compacted, or ‘organized’, due to a consolidation process that was associated with the emergence of meta-*cis* regulation of the cluster, and probably facilitated by the 2R-WGD events that occurred during vertebrate evolution⁴². Hox clusters of cyclostomes are, on the other hand, more akin to ‘disorganized’ types of clusters, like the one of amphioxus⁴² –because of their extremely large sizes–, suggesting that this consolidation did not start in the last common ancestor of vertebrates, but rather was a progressive gnathostome-specific process⁴² (Fig. 4). Further

functional analyses of the regulatory mechanisms of cyclostomes' Hox clusters, with special focus on determining the presence or absence of global regulatory elements outside the clusters, will be needed to clarify whether the consolidation process was indeed a consequence of the acquisition of a global regulatory mode for the cluster, or if, on the other hand, this meta-*cis* regulation was already present in the last common ancestor of vertebrates, before the consolidation process started. Moreover, the timing of the vertebrate 2R-WGD, i.e., whether or not these events are shared between gnathostomes and cyclostomes, is one of the most important questions that remain open about the origin of the vertebrate genome architecture, and solving it will be also helpful to decipher whether the vertebrate genome duplications facilitated the consolidation process.

METHODS

Animal sampling, experiments and aquarium maintenance

E. burgeri embryos (staged according to refs. 18 and 19) used in this study were obtained from adult hagfish individuals captured in the Japan Sea off Shimane prefecture as previously described⁴³, during August of a given year. Eggs were laid in a cage deposited in the natural environment in the sea in October of the same year. Deposited eggs were then incubated in laboratory aquariums with artificial sea water at 16 °C under controlled conditions, until developing embryos are apparent around February or March of the following year. Hagfish embryos used for RNA-seq were from adults captured in 2010, and were assayed in February 2011 (total developing embryos 8 out of >150 eggs). Hagfish embryos used for *in situ* hybridization were from adults captured in 2016 (stage 40), 2013 (stage 45) and 2014 (stage 60), and embryos fixed in March 2017, 2014 and 2015, respectively. Sections of stage 53 were from an embryo previously reported⁴⁴. Lamprey (*L. camtschaticum*) and cloudy catshark (*S. torazame*) embryos were obtained as previously described in refs. 45 and 44,

respectively. Lamprey and catshark embryos were staged according to refs. 34 and 35, respectively. The sampling and experiments were conducted according to the institutional and national guidelines for animal ethics, approved by the RIKEN Animal Experiments Committee (approval ID: H14-25-24).

RNA-seq data and transcriptome assemblies

Total RNA samples from three whole embryos of *E. burgeri* (Fig. 1c-e) and the head region of a hatched juvenile were used to prepare RNA-seq libraries and sequenced individually on different HiSeq and MiSeq platforms (one embryo at stage 28/30: Illumina TruSeq RNA Sample Prep Kit, non-strand-specific library, sequenced with a HiSeq1000 platform; one embryo at stage 35 and one at 40/45: one strand-specific library each using TruSeq RNA Sample Prep Kit modified with the dUTP method⁴⁶ and sequenced in a HiSeq2000, and a further non-strand-specific library with Illumina TruSeq RNA Library Prep Kit and sequenced in a MiSeq platform for the former; one juvenile's head: TruSeq RNA Sample Prep Kit, non-strand-specific library, HiSeq1500). Total RNA samples from separate pools of embryos of *L. camtschaticum* at stages 15/16, 18, 20, 22, 24, 26 and 28 (20-30 embryos per stage), and of *S. torazame* at stages 15/16 (6 embryos), 18 (9), 20 (10), 22 (9), 25 (5), 27 (5) and 28 (2) were used to prepare strand-specific libraries (Illumina TruSeq Stranded RNA Library Prep Kit). Lamprey and shark libraries were sequenced on a HiSeq1500 platform. Reads coming from mitochondrial DNA were filtered out using mirabait (bundled with MIRA). Then, reads were preprocessed with MIRA⁴⁷ v.4.9.5_2, using the option 'parameters = -GE:ppo=yes' in the manifest file. In the case of the hagfish, the resulting reads were then assembled with Trinity v2.1.1⁴⁸ following 3 different strategies: (1) assembly of all reads together; (2) idem, but including a digital normalization step within Trinity (--normalize_reads), and (3), assembly of RNA-seq data from each embryo separately and

324 further integrated with CD-HIT-EST⁴⁹ v4.6.4 with parameters ‘-c 0.98’. A fourth
 325 assembly was done with SOAPdenovoTrans v1.03⁵⁰ using all reads simultaneously and
 326 multiple k-mers (19, 21, 23, 25, 27, 29, 31: with ‘SOAPdenovo-Trans-31mer’
 327 command; and 41, 51, 61, 71, 81, 91: with ‘SOAPdenovo-Trans-127mer’), with a final
 328 integration with CD-HIT-EST. In the case of the lamprey and shark, reads were assembled
 329 according to 3 different pipelines: (1) assembly with Trinity v2.1.1 of reads coming from each
 330 pool of embryos independently, taking into account the strand-specific information (--
 331 SS_lib_type RF), and integrated with CD-HIT-EST v4.6.4; (2) idem, but not taking into
 332 account the strand-specific information; and (3) assembly of all reads together. In the case of
 333 the lamprey, a fourth assembly strategy was carried out by means of integrating a genome-
 334 guided assembly (option --genome_guided_bam of Trinity, mapping the RNA-seq reads
 335 to *L. camtschaticum* 1.0 genome¹⁰ with the splice-aware mapper HISAT2⁵¹) and the above *de*
 336 *novo* assembly #3, using the PASA v2.0.2⁵² pipeline
 337 (http://pasapipeline.github.io/#A_ComprehensiveTranscriptome). Finally, completeness
 338 assessments of all versions were done using CEGMA v2.5 and BUSCO v1.1b1 programs, as
 339 previously described⁵³ (Supplementary Tables 1-3). The most complete versions of *E.*
 340 *burgeri* and *L. camtschaticum* were selected for further analysis. In the case of *S. torazame*,
 341 although strategy #3 gave as a result a more complete transcriptome in general, it contained
 342 more fragmented Hox genes than version #1, and therefore we selected the latter. All *E.*
 343 *burgeri*, *L. camtschaticum* and *S. torazame* RNA-seq data have been deposited in NCBI
 344 GenBank under the BioProject number PRJNA371391. Amphioxus *B. belcheri* transcriptome
 345 was assembled using previously published RNA-seq data, from the NCBI’s SRA database,
 346 under BioProject numbers PRJNA310680³⁷ and PRJNA214454⁵⁴. *B. belcheri* RNA-seq reads
 347 from the former BioProject were subjected to adaptor trimming with cutadapt v1.10⁵⁵. All *B.*
 348 *belcheri* RNA-seq data was then assembled following the same above-mentioned pipeline for

the lamprey transcriptome (strategy #3), using the previously published *B. belcheri* genome⁵⁴ for the PASA pipeline.

All Trinity commands were executed using the `--group_pairs_distance 999` parameter value⁵⁶.

BAC library, clone screening and PacBio sequencing and assembly

Blood was drawn from the caudal sub-cutaneous sinus of one adult specimen of *E. burgeri* using a heparin-rinsed disposable syringe. The whole blood sample was immediately frozen in liquid nitrogen, and used for DNA extraction. A BAC library consisting of 129,024 clones, with an average insert size of 100 Kbp (~4.4X of the *E. burgeri* genome size), was constructed using the pCCBAC1 vector⁵⁷ [CopyControlTM BAC Cloning Kit (*Hind*III) (EPICENTRE)] and pooled into 96-well and 384-well plates according to the Matrix Pool and Superpool Strategy⁵⁸ by Amplicon Express (Pullman, WA, USA). The BAC library was screened for Hox-containing clones by means of PCR with specific primers⁵⁸. Positive BAC clones were extracted with the QIAGEN Large-Construct kit, and sequenced in different pools using SMRT technology on a PacBio RS platform using XL-C2 chemistry, or on a RSII platform using P4-C2 chemistry. pCCBAC1 vector sequence were masked using a script from ref. 59 with minor modifications. BAC clones were assembled using masked subreads with MIRA⁴⁷ v4.9.5_2. The sequence of the BAC clones used in this study have been deposited in GenBank (accession numbers MF182102-MF182109).

Genome sequencing and assembly

Germ line DNA for whole genome shotgun (WGS) sequencing, derived from the testis of a single male hagfish, *E. burgeri*, was sequenced on an Illumina HiSeq 2500TM platform. In total, we sequenced five pair-end (174-bp, 234-bp, 242-bp, 279-bp and 612-bp) and five

mate-pair (5-Kbp, 5~7-Kbp, 7~10-Kbp, 10~15-Kbp and 15~20-Kbp) libraries, generating >300X coverage of the estimated 2.906 Gb-long genome of the hagfish. All short-read data were corrected by SOAPec v2.01⁶⁰ using >40X data. Assembly of the hagfish genome was performed with ABySS v1.9.0⁶¹ with a k-mer size of 79, followed by a scaffolding step with SOAPdenovo v2.04-r241⁶⁰ software (parameter ‘-K 41 -d 1 -M 2 -F’). Gaps were finally filled with GapCloser v1.12-r6⁶⁰. The resulting assembly (size, ~2.59 Gb; N50, ~439 Kbp) was used for the screening of Hox clusters. Hox-containing scaffolds were then aligned against the BAC clones using MUMmer v3.23⁶² and visualized using mummerplot, bundled within the same software. Sequences of Hox-containing scaffolds, as well as those of *E. burgeri* *Lnp* and *Evx* (whose sequences are not complete in the genome) have been deposited in GenBank under accession numbers MF398213-MF398235. A publication with more detailed and in-depth analysis of the *E. burgeri* genome is now in preparation.

Identification of Hox genes

UniProt Knowledgebase (UniProtKB) database (<http://www.uniprot.org/>) was searched for entries containing the term “Hox” and restricted to Eumetazoans (name:hox, taxonomy:6072; UniProt release 2015_11). Resulting entries were downloaded and used as queries against the transcriptome assembly and genome of the hagfish by means of TBLASTN (NCBI BLAST v2.2.31+⁶³). The best BLAST hits were then used as queries against the whole UniProtKB database using BLASTX. Those transcripts whose reciprocal best hit was a Hox gene were kept and manually inspected for false positives. Lamprey Hox genes were downloaded from NCBI GenBank¹⁰ and blasted against our lamprey transcriptome assembly to identify Hox transcripts. *Hox4* η , *Hox7* ϵ , *Hox9* ζ , *Hox11* δ , *Hox13* α , *Hox13* ϵ , *Hox13* ζ and *Hox14* ϵ were not found in our transcriptome assembly. We found an extra, unreported Hox1 paralogous gene, which we named *Hox1* ζ (following the nomenclature from ref. 10). *Scyliorhinus canicula*

399 Hox genes³⁰ sequences were downloaded from NCBI GenBank and used as queries to
400 identify orthologous sequences in our *S. torazame* transcriptome by means of TBLASTN.
401 The *L. camtschaticum* *Hox1ζ* and *S. torazame* Hox gene sequences were deposited in
402 GenBank (accession numbers MF398236-MF398269).

403

404 **cDNA cloning and section *in situ* hybridization**

405 Selected Hox genes were cloned from cDNA prepared for a previous study⁴⁴ using specific
406 primers. *In situ* hybridization on paraffin wax-embedded sections of stage 45 and 60 hagfish
407 embryos was performed according to refs. 44, 45. Haematoxylin and eosin (H&E) staining on
408 paraffin sections of stage 60 embryo was carried out by standard protocol. H&E stained
409 sections were further stained with Alcian Blue⁶⁴.

410

411 **3D reconstruction of the hagfish embryos**

412 The 3D reconstruction images of hagfish embryos were reconstructed based on images taken
413 of 1 in every 10 histological sagittal sections at 6 μm, stained with standard haematoxylin and
414 eosin staining protocols for the stage 40 embryo, and 1 in 2 unstained sections at 8 μm for the
415 stage 45 embryo. Reconstructed images were acquired using Avizo software (Visualization
416 Sciences Group). Stage 53 reconstruction is from an embryo used previously¹⁹.

417

418 **Molecular phylogenetic analyses**

419 The Hox genes nucleotide sequences for different chordates and outgroups were mined from
420 NCBI GenBank, Ensembl (www.ensembl.org), EchinoBase
421 (<http://www.echinobase.org/Echinobase/>), or, in some instances, manually annotated (see
422 Supplementary Table 4 for accession numbers of genes used in the analyses). Hox genes
423 sequences of amphioxus species *B. lanceolatum* and *B. floridae* are from refs. 65-67. Five

424 datasets based on different gene content were assembled: 1) Anterior genes (Hox1-3), 2) Hox
425 4 genes, 3) Central genes (Hox4-8), 4) Posterior genes (Hox9-14), and 5) all Hox genes
426 together. The datasets were aligned using MAFFT v7.123b⁶⁸ using the “*auto*” option, regions
427 of ambiguous alignment trimmed with Gblocks v0.91b⁶⁹ using the less stringent options.
428 Alignments were visually inspected with BioEdit v7.2.6⁷⁰. Phylogenetic trees were inferred
429 with RAxML v8.2.10⁷¹ using a random starting tree, the evolutionary model LG + Gamma +
430 Invariants with empirical base frequencies, and 1000 rapid bootstrap replicates. Trees were
431 edited using FigTree v1.4.3 (<http://tree.bio.ed.ac.uk/software/figtree/>).

432

433 **Expression profiling of Hox genes.**

434 RNA-seq reads from individual embryos were used to quantify transcripts of the selected
435 transcriptomes of *E. burgeri* and *S. torazame* using Perl scripts
436 ‘align_and_estimate_abundance.pl’ and ‘align_and_estimate_abundance.pl’, bundled with
437 Trinity v.2.1.1, and using RSEM v1.2.28⁷² as quantification method
438 (<https://github.com/trinityrnaseq/trinityrnaseq/wiki/Trinity-Transcript-Quantification>). Hox
439 transcripts of *S. torazame* were directly quantified using RSEM with data from each
440 embryonic stage. TPM values from either genes (for Hox genes represented with a single
441 transcript in the assemblies) or isoforms (for Hox genes represented with several transcripts)
442 were then selected and a heat map analyses of the log(TPM+0.1) were conducted in R using
443 heatmap.2 (gplots package⁷³) scaling by gene (row Z-score), and implemented in RStudio
444 v1.0.136⁷⁴ [with R v3.3.0 (2016-05-03)⁷⁵]. *B. belcheri* Hox transcripts were quantified using
445 previously published DGE-seq data³⁷ with DGE-EM v1.0.0⁷⁶ software, and FPKM values
446 were analysed as above.

447

448 **Data availability.**

RNA-seq data generated in this study have been deposited in SRA, under the BioProject number PRJNA371391. Sequences generated and analysed in this study have been deposited in NCBI GenBank under accession numbers MF182102-MF182109 and MF398213-MF398269.

References

- 1 Pearson, J. C., Lemons, D. & McGinnis, W. Modulating Hox gene functions during animal body patterning. *Nat. Rev. Genet.* **6**, 893-904 (2005).
- 2 Garcia-Fernàndez, J. The genesis and evolution of homeobox gene clusters. *Nat. Rev. Genet.* **6**, 881-892 (2005).
- 3 Pascual-Anaya, J., D'Aniello, S., Kuratani, S. & Garcia-Fernandez, J. Evolution of Hox gene clusters in deuterostomes. *BMC Dev. Biol.* **13** (2013).
- 4 Smith, J. J. & Keinath, M. C. The sea lamprey meiotic map improves resolution of ancient vertebrate genome duplications. *Genome Res* **25**, 1081-1090 (2015).
- 5 Dehal, P. & Boore, J. L. Two rounds of whole genome duplication in the ancestral vertebrate. *PLoS Biol.* **3**, e314 (2005).
- 6 Putnam, N. H. *et al.* The amphioxus genome and the evolution of the chordate karyotype. *Nature* **453**, 1064-1071 (2008).
- 7 Van de Peer, Y., Maere, S. & Meyer, A. The evolutionary significance of ancient genome duplications. *Nat. Rev. Genet.* **10**, 725-732 (2009).
- 8 Van de Peer, Y., Maere, S. & Meyer, A. 2R or not 2R is not the question anymore. *Nat. Rev. Genet.* **11**, 166 (2010).

- 472 9 Kuraku, S., Meyer, A. & Kuratani, S. Timing of genome duplications relative to the
473 origin of the vertebrates: did cyclostomes diverge before or after? *Mol. Biol. Evol.*
474 **26**, 47-59 (2009).
- 475 10 Mehta, T. K. *et al.* Evidence for at least six Hox clusters in the Japanese lamprey
476 (*Lethenteron japonicum*). *Proc Natl Acad Sci U S A* **110**, 16044-16049 (2013).
- 477 11 Smith, J. J. *et al.* Sequencing of the sea lamprey (*Petromyzon marinus*) genome
478 provides insights into vertebrate evolution. *Nat. Genet.* **45**, 415-421 (2013).
- 479 12 Seo, H. C. *et al.* Hox cluster disintegration with persistent anteroposterior order
480 of expression in *Oikopleura dioica*. *Nature* **431**, 67-71 (2004).
- 481 13 Dollé, P., Izpisúa-Belmonte, J. C., Falkenstein, H., Renucci, A. & Duboule, D.
482 Coordinate expression of the murine *Hox-5* complex homoeobox-containing
483 genes during limb pattern formation. *Nature* **342**, 767-772 (1989).
- 484 14 Izpisúa-Belmonte, J. C., Falkenstein, H., Dollé, P., Renucci, A. & Duboule, D. Murine
485 genes related to the *Drosophila AbdB* homeotic genes are sequentially expressed
486 during development of the posterior part of the body. *EMBO J.* **10**, 2279-2289
487 (1991).
- 488 15 Wang, S. *et al.* Scallop genome provides insights into evolution of bilaterian
489 karyotype and development. *Nat. Ecol. Evol.* **1**, 0120 (2017).
- 490 16 Takio, Y. *et al.* Hox gene expression patterns in *Lethenteron japonicum*
491 embryos—Insights into the evolution of the vertebrate Hox code. *Dev. Biol.* **308**,
492 606-620 (2007).

- 493 17 Stadler, P. F. *et al.* Evidence for independent *Hox* gene duplications in the hagfish
494 lineage: a PCR-based gene inventory of *Eptatretus stoutii*. *Mol. Phylogenet. Evol.*
495 **32**, 686-694 (2004).
- 496 18 Dean, B. On the Embryology of *Bdellostoma stouti*: A General Account of Myxinoid
497 Development from the Egg and Segmentation to Hatching. 220-276 (G. Fischer,
498 1899).
- 499 19 Oisi, Y., Ota, K. G., Kuraku, S., Fujimoto, S. & Kuratani, S. Craniofacial development
500 of hagfishes and the evolution of vertebrates. *Nature* **493**, 175-180 (2013).
- 501 20 Smith, J. J., Baker, C., Eichler, E. E. & Amemiya, C. T. Genetic consequences of
502 programmed genome rearrangement. *Curr. Biol.* **22**, 1524-1529 (2012).
- 503 21 Kohno, S.-i., Kubota, S. & Nakai, Y. in *The Biology of Hagfishes*. 81-100 (Springer
504 Netherlands, 1998).
- 505 22 Kuraku, S. Impact of asymmetric gene repertoire between cyclostomes and
506 gnathostomes. *Semin. Cell Dev. Biol.* **24**, 119-127 (2013).
- 507 23 Ohno, S. *Evolution by Gene Duplication*. (Springer-Verlag, 1970).
- 508 24 Kasahara, M. The 2R hypothesis: an update. *Curr. Opin. Immunol.* **19**, 547-552
509 (2007).
- 510 25 Martin, K. J. & Holland, P. W. Enigmatic orthology relationships between *Hox*
511 clusters of the African butterfly fish and other teleosts following ancient whole-
512 genome duplication. *Mol. Biol. Evol.* **31**, 2592-2611 (2014).
- 513 26 Takio, Y. *et al.* Evolutionary biology: lamprey *Hox* genes and the evolution of
514 jaws. *Nature* **429**, 1 p following 262 (2004).

- 515 27 Schiemann, S. M. *et al.* Clustered brachiopod Hox genes are not expressed
516 collinearly and are associated with lophotrochozoan novelties. *Proc. Natl. Acad.*
517 *Sci. U. S. A.* **114**, E1913-E1922 (2017).
- 518 28 Parker, H. J., Bronner, M. E. & Krumlauf, R. A *Hox* regulatory network of hindbrain
519 segmentation is conserved to the base of vertebrates. *Nature* **514**, 490-493
520 (2014).
- 521 29 Kuratani, S., Horigome, N., Ueki, T., Aizawa, S. & Hirano, S. Stereotyped axonal
522 bundle formation and neuromeric patterns in embryos of a cyclostome,
523 *Lampetra japonica*. *J. Comp. Neurol.* **391**, 99-114 (1998).
- 524 30 Oulion, S. *et al.* Evolution of repeated structures along the body axis of jawed
525 vertebrates, insights from the *Scyliorhinus canicula* Hox code. *Evol. Dev.* **13**, 247-
526 259 (2011).
- 527 31 Wingate, R. J. The rhombic lip and early cerebellar development. *Curr. Opin.*
528 *Neurobiol.* **11**, 82-88 (2001).
- 529 32 Murakami, Y. *et al.* Segmental development of reticulospinal and branchiomotor
530 neurons in lamprey: insights into the evolution of the vertebrate hindbrain.
531 *Development* **131**, 983-995 (2004).
- 532 33 Kuraku, S. *et al.* Noncanonical role of Hox14 revealed by its expression patterns
533 in lamprey and shark. *Proc. Natl. Acad. Sci. U. S. A.*, 0710947105 (2008).
- 534 34 Tahara, Y. Normal stages of development in the lamprey, *Lampetra reissneri*
535 (Dybowski). *Zool. Sci.* **5**, 109-118 (1988).

536 35 Ballard, W. W., Mellinger, J. & Lechenault, H. A series of normal stages for
537 development of *Scyliorhinus canicula*, the lesser spotted dogfish (*Chondrichthyes*:
538 *Scyliorhinidae*). *J. Exp. Zool.* **267**, 318-336 (1993).

539 36 Ikuta, T., Yoshida, N., Satoh, N. & Saiga, H. *Ciona intestinalis* Hox gene cluster: Its
540 dispersed structure and residual colinear expression in development. *Proc. Natl.*
541 *Acad. Sci. U. S. A.* **101**, 15118-15123 (2004).

542 37 Yang, K. Y. *et al.* Transcriptome analysis of different developmental stages of
543 amphioxus reveals dynamic changes of distinct classes of genes during
544 development. *Sci. Rep.* **6**, 23195 (2016).

545 38 Pascual-Anaya, J. *et al.* Broken colinearity of the amphioxus Hox cluster. *EvoDevo*
546 **3** (2012).

547 39 Deschamps, J. & Duboule, D. Embryonic timing, axial stem cells, chromatin
548 dynamics, and the Hox clock. *Genes Dev.* **31**, 1406-1416 (2017).

549 40 Aronowicz, J. & Lowe, C. J. Hox gene expression in the hemichordate *Saccoglossus*
550 *kowalevskii* and the evolution of deuterostome nervous systems. *Integr. Comp.*
551 *Biol.* **46**, 890-901 (2006).

552 41 Gonzalez, P., Uhlinger, K. R. & Lowe, C. J. The Adult Body Plan of Indirect
553 Developing Hemichordates Develops by Adding a Hox-Patterned Trunk to an
554 Anterior Larval Territory. *Curr. Biol.* **27**, 87-95 (2017).

555 42 Duboule, D. The rise and fall of Hox gene clusters. *Development* **134**, 2549-2560
556 (2007).

557 43 Oisi, Y., Kakitani, O., Kuratani, S. & Ota, K. G. in *In Situ Hybridization Methods* (ed
558 Giselbert Hauptmann) 249-262 (Springer New York, 2015).

559 44 Sugahara, F. *et al.* Evidence from cyclostomes for complex regionalization of the
560 ancestral vertebrate brain. *Nature* **531**, 97-100 (2016).

561 45 Sugahara, F., Murakami, Y. & Kuratani, S. in *In Situ Hybridization Methods* (ed
562 Giselbert Hauptmann) 263-278 (Springer New York, 2015).

563 46 Sultan, M. *et al.* A simple strand-specific RNA-Seq library preparation protocol
564 combining the Illumina TruSeq RNA and the dUTP methods. *Biochem. Biophys.*
565 *Res. Commun.* **422**, 643-646 (2012).

566 47 Chevreux, B., Wetter, T. & Suhai, S. in *Computer Science and Biology: Proceedings*
567 *of the German Conference on Bioinformatics (GCB)* 45-56 (1999).

568 48 Grabherr, M. G. *et al.* Full-length transcriptome assembly from RNA-Seq data
569 without a reference genome. *Nat. Biotechnol.* **29**, 644-652 (2011).

570 49 Li, W. & Godzik, A. Cd-hit: a fast program for clustering and comparing large sets
571 of protein or nucleotide sequences. *Bioinformatics* **22**, 1658-1659 (2006).

572 50 Xie, Y. *et al.* SOAPdenovo-Trans: de novo transcriptome assembly with short
573 RNA-Seq reads. *Bioinformatics* **30**, 1660-1666 (2014).

574 51 Kim, D., Langmead, B. & Salzberg, S. L. HISAT: a fast spliced aligner with low
575 memory requirements. *Nat. Methods* **12**, 357-360 (2015).

576 52 Haas, B. J. *et al.* Improving the *Arabidopsis* genome annotation using maximal
577 transcript alignment assemblies. *Nucleic Acids Res.* **31**, 5654-5666 (2003).

578 53 Hara, Y. *et al.* Optimizing and benchmarking de novo transcriptome sequencing:
579 from library preparation to assembly evaluation. *BMC Genomics* **16**, 977 (2015).

580 54 Huang, S. *et al.* Decelerated genome evolution in modern vertebrates revealed by
581 analysis of multiple lancelet genomes. *Nat. Commun.* **5**, 5896 (2014).

582 55 Martin, M. Cutadapt removes adapter sequences from high-throughput
583 sequencing reads. *EMBnet.journal* **17**, 10-12 (2011).

584 56 Macmanes, M. D. On the optimal trimming of high-throughput mRNA sequence
585 data. *Front. Genet.* **5**, 13 (2014).

586 57 Wild, J., Hradecna, Z. & Szybalski, W. Conditionally amplifiable BACs: switching
587 from single-copy to high-copy vectors and genomic clones. *Genome Res.* **12**,
588 1434-1444 (2002).

589 58 Bouzidi, M. F. *et al.* A sunflower BAC library suitable for PCR screening and
590 physical mapping of targeted genomic regions. *Theor. Appl. Genet.* **113**, 81-89
591 (2006).

592 59 Huddleston, J. *et al.* Reconstructing complex regions of genomes using long-read
593 sequencing technology. *Genome Res.* **24**, 688-696 (2014).

594 60 Luo, R. *et al.* SOAPdenovo2: an empirically improved memory-efficient short-
595 read de novo assembler. *GigaScience* **1**, 18 (2012).

596 61 Simpson, J. T. *et al.* ABySS: a parallel assembler for short read sequence data.
597 *Genome Res.* **19**, 1117-1123 (2009).

598 62 Kurtz, S. *et al.* Versatile and open software for comparing large genomes. *Genome*
599 *Biol.* **5**, R12 (2004).

600 63 Camacho, C. *et al.* BLAST+: architecture and applications. *BMC Bioinformatics* **10**,
601 421 (2009).

602 64 Nowicki, J. L., Takimoto, R. & Burke, A. C. The lateral somitic frontier: dorso-
603 ventral aspects of antero-posterior regionalization in avian embryos. *Mech. Dev.*
604 **120**, 227-240 (2003).

605 65 Amemiya, C. T. *et al.* The amphioxus *Hox* cluster: characterization, comparative
606 genomics, and evolution. *J. Exp. Zool. B Mol. Dev. Evol.* **310**, 465-477 (2008).

607 66 Holland, L. Z. *et al.* The amphioxus genome illuminates vertebrate origins and
608 cephalochordate biology. *Genome Res.* **18**, 1100-1111 (2008).

609 67 Pascual-Anaya, J., D'Aniello, S. & Garcia-Fernàndez, J. Unexpectedly large number
610 of conserved noncoding regions within the ancestral chordate *Hox* cluster. *Dev.*
611 *Genes Evol.* **218**, 591-597 (2008).

612 68 Katoh, K. & Standley, D. M. MAFFT multiple sequence alignment software version
613 7: improvements in performance and usability. *Mol. Biol. Evol.* **30**, 772-780
614 (2013).

615 69 Castresana, J. Selection of conserved blocks from multiple alignments for their
616 use in phylogenetic analysis. *Mol. Biol. Evol.* **17**, 540-552 (2000).

617 70 Hall, T. A. BioEdit: a user-friendly biological sequence alignment editor and
618 analysis program for Windows 95/98/NT. *Nucleic Acids Symp. Ser.* **41**, 95-98
619 (1999).

620 71 Stamatakis, A. RAxML version 8: a tool for phylogenetic analysis and post-
621 analysis of large phylogenies. *Bioinformatics* **30**, 1312-1313 (2014).

622 72 Li, B. & Dewey, C. N. RSEM: accurate transcript quantification from RNA-Seq data
623 with or without a reference genome. *BMC Bioinformatics* **12**, 323 (2011).

624 73 gplots: Various R programming tools for plotting data. R package version 3.0.1
625 (2016).

626 74 RStudio: Integrated Development for R. (R. RStudio, Inc., Boston, MA, 2016).

627 75 R: A language and environment for statistical computing (Vienna, Austria, 2016).

Nicolae, M. & Măndoiu, I. in *Bioinformatics Research and Applications: 7th International Symposium, ISBRA 2011, Changsha, China, May 27-29, 2011. Proceedings* (eds Jianer Chen, Jianxin Wang, & Alexander Zelikovsky) 392-403 (Springer Berlin Heidelberg, 2011).

Acknowledgments

We thank Y. Oisi and S. Fujimoto for providing preliminary hagfish Hox sequences; J.M. Martín-Durán, I. Maeso, M. Irimia and C. Böhmer for fruitful discussions; O. Kakitani for hagfish sampling; K. Shirato for shark sampling; S. Shibuya and K. Yamamoto for maintenance of aquarium tanks; S. Kuraku, K. Itomi, C. Tanegashima, K. Tatsumi and O. Nishimura from the Phyloinformatics Unit, RIKEN CLST, for RNA-seq data production; J. Huddleston and E. Eichler for providing the code to mask BAC vector sequences from PacBio reads; I. Mandoiu for his help using DGE-EM software; and B. Chevreux for his help with the MIRA assembler. This work was supported by the Chinese Academy of Sciences program XDB13000000 to W.W., and by a Grant-in-Aid for Scientific Research (A) 15H02416 (Japan Society for the Promotion of Science), a Grant-in-Aid for Scientific Research on Innovative Areas (Research in a Proposed Research Area) 17H06384 (Ministry of Education, Culture, Sports, Science and Technology of Japan), and a Naito Grant for the Promotion of Focused Research (The Naito Foundation) to S.K.

Author contributions

J.P.-A. conceived the project, designed the experiments, and wrote the paper. J.P.-A, F.S., S.H. and W.T. obtained the hagfish embryos. J.P.-A, I.S., F.S., S.H., W.T. and A.R.-V performed experiments. K.O. built the BAC library. R.Y. and W.W. sequenced and

assembled the *E. burgeri* draft genome. J.P. performed the phylogenetic analyses. All authors analysed and discussed the data. All authors approved the final version of the manuscript.

Competing interests

The authors declare no competing financial interests.

Figure legends

Figure 1. Hox cluster evolution in chordates. a, Phylogenetic tree of chordates, showing the two major groups of vertebrates –cyclostomes (hagfish and lamprey) and gnathostomes (jawed vertebrates, e.g., mouse and shark) – together with cephalochordates (amphioxus), displaying their known Hox repertoires. Numbers on the nodes indicate the putative number of Hox clusters in each last common ancestor. **b**, *E. burgeri* Hox genes and clusters found in this study, drawn to scale. All Hox genes are transcribed in the same orientation, from left to right. Orientation of transcription of non-Hox syntenic genes are indicated by arrowheads. Solid horizontal lines correspond to single scaffolds. Double diagonal lines separate two contiguous scaffolds, based on BAC sequences connecting them (Supplementary Fig. 2). *Hox3II* and *Hox3VI* genes have corresponding exons 1 and 2 in two different scaffolds, which have been put together based on both BAC and transcriptomics evidences. e1, exon 1; e2, exon 2. Asterisk over *miR-10III* indicates that this microRNA is within a 5'UTR intron of *Hox4III*; hash symbol, *Hox6III* exon2 is not assembled in place, but in a separate small contig containing its sequence. **c-e**, *E. burgeri* embryos used for the transcriptomics analysis, at Dean stages 28-30 (**c**), 35 (**d**) and 40-45 (**e**). fb, forebrain; hb, hindbrain; mb, midbrain; ov, otic vesicle; ph, pharynx; som, somites. Scale bars, 1 mm.

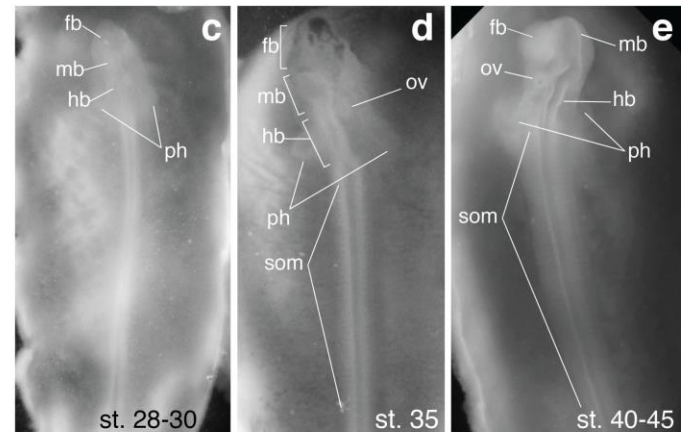
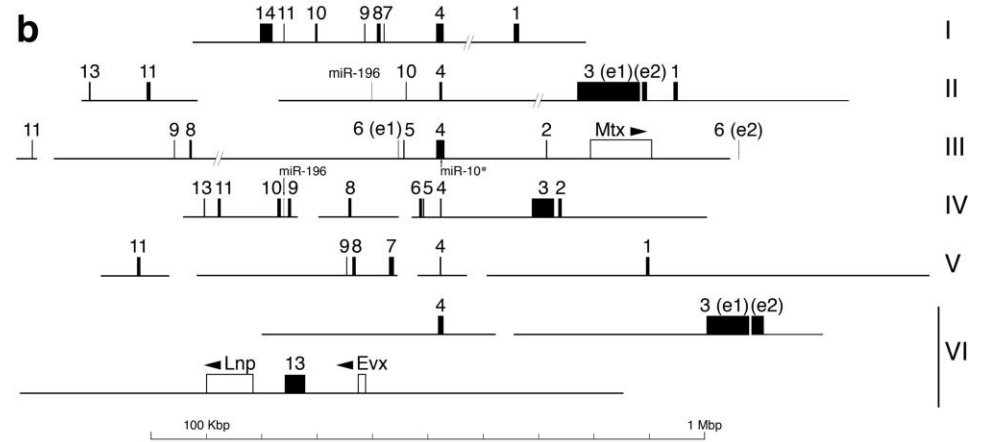
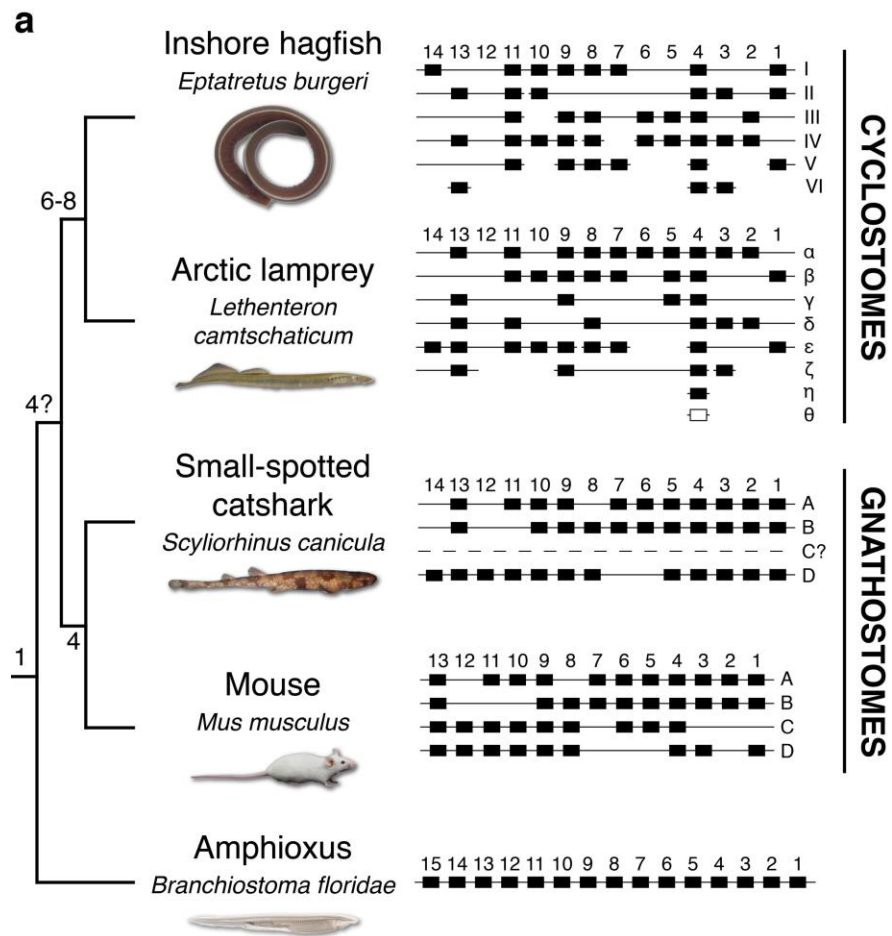
Figure 2. Molecular phylogenetic tree of vertebrate Hox genes. 1000-replicate Maximum Likelihood tree of representative Hox genes of all paralogy groups in vertebrates. The branches have been color-coded by paralogy group (Hox1-14). Red and blue branches denote *E. burgeri* and *L. camtschaticum* Hox genes, respectively. Black branches correspond to invertebrate Hox counterparts (amphioxus – *Branchiostoma floridae* and *Branchiostoma lanceolatum* – and sea urchin – *Strongylocentrotus purpuratus* –). Note that no hagfish or lamprey sequence have been found within the Hox12 group (denoted with square brackets). The same tree, with bootstrap values and branch tip names can be found in Supplementary Fig. 7.

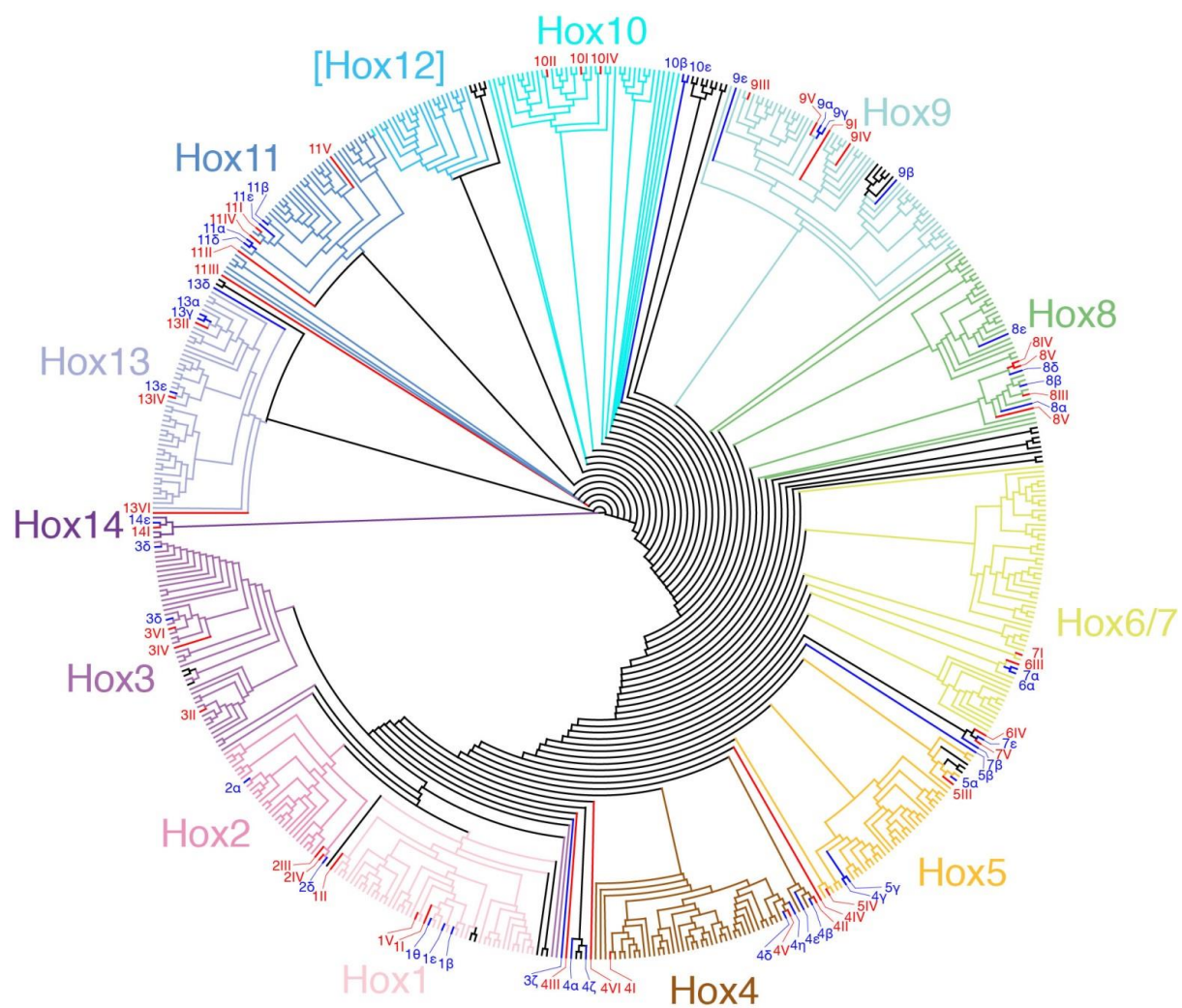
Figure 3. Spatial colinearity of hagfish Hox genes in the hindbrain of *E. burgeri* embryos. **a, b**, Embryos at stage Bashford Dean 40 (**a**) and 45 (**b**) used for *in situ* hybridizations on sections. The inset square brackets mark the head regions, used for sagittal sectioning. **a', a'', b', b''**, 3D Avizo reconstructions of the heads of the embryos shown in **a** and **b**, respectively, showing the main internal anatomy of the brain and main head structures. The central nervous systems are in purple; ectoderm is in light blue; endoderm is in yellow; otic vesicle in green; and notochord is in light red. These embryos are the source of the sections shown in **c-j** (stage 40) and **k-s** (stage 45). **c-w**, Spatial colinearity displayed by expression patterns of *E. burgeri* *Hox1IV* (**c**), *Hox2III* (**k**), *Hox2IV* (**d, l, t**), *Hox3II* (**e, m, u**), *Hox3VI* (**f, n**), *Hox4I* (**g, o**), *Hox4IV* (**h, p, v**), *Hox4VI* (**q**), *Hox5III* (**i, r, w**) and *Hox5IV* (**j, s**), revealed by *in situ* hybridization on sections of an embryos shown in **a, b** and **x**. **x, x'**, 3D Avizo reconstructions of the head of an embryo at stage Bashford Dean 53. The central nervous systems are in purple; ectoderm is in light blue; endoderm is in yellow; otic vesicle in green; and notochord is in light red. **y**, Expression patterns shown in (**c-j**) in the hindbrain are aligned according to rhombomere (r) segmentation, showing nested expression patterns of

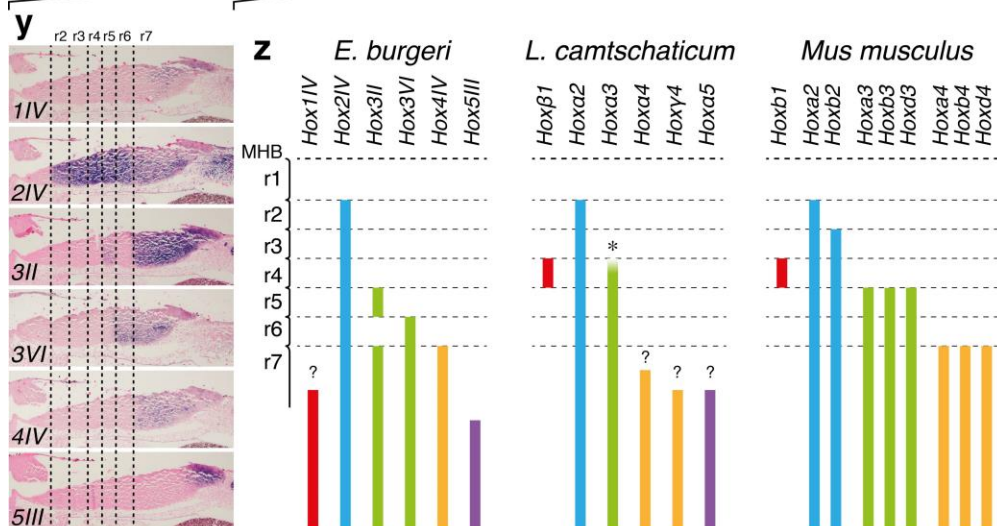
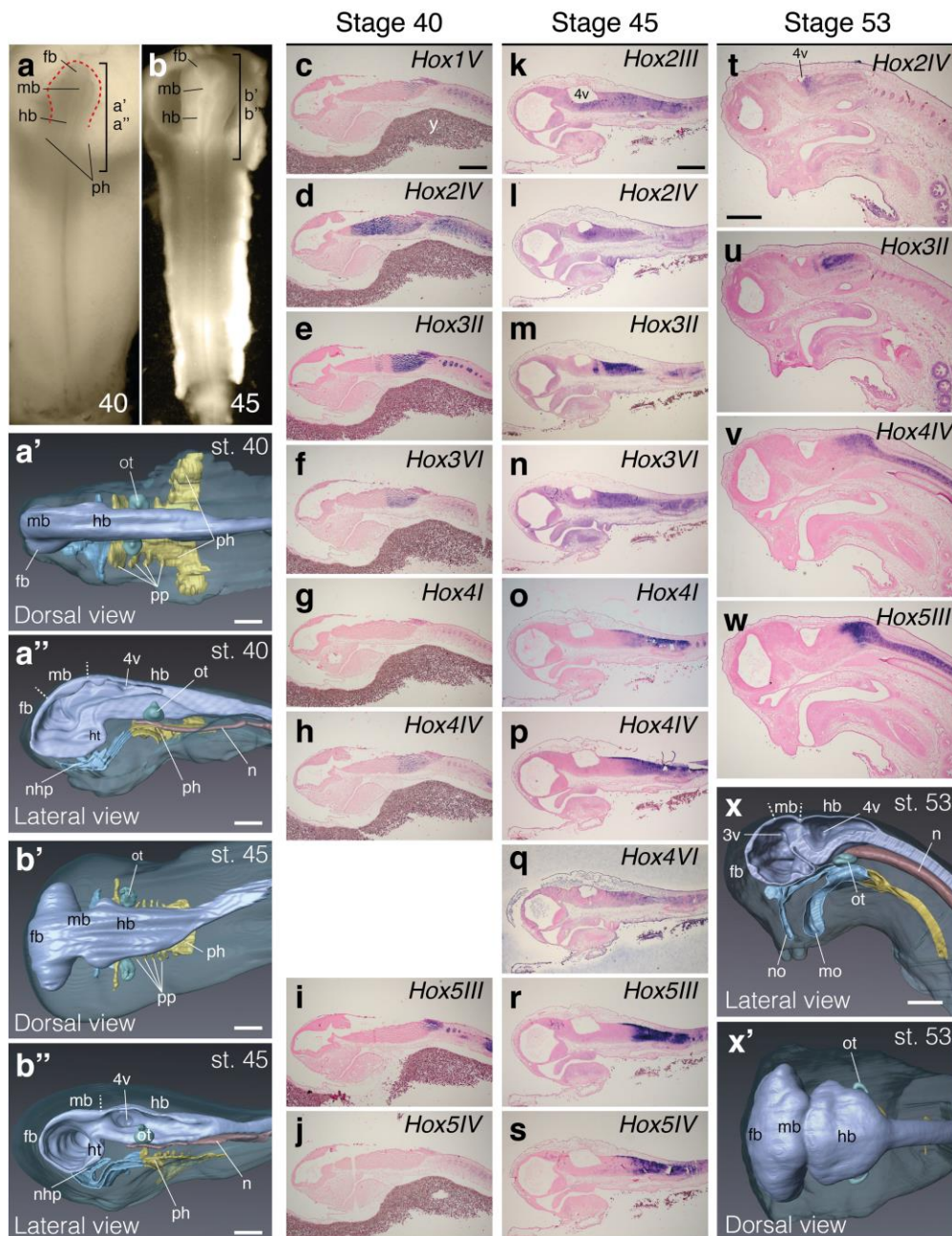
701 *Hox1-5* paralogs in the hagfish hindbrain. **z**, Schematic diagrams summarizing the expression
 702 patterns of Hox genes in the hindbrain of the lamprey, the hagfish and mouse, with nested
 703 anterior limits coinciding with rhombomere borders, and showing the overall conservation
 704 among the groups. f, forebrain; h, hindbrain; ht, hypothalamus; m, midbrain; mo, mouth; n,
 705 notochord; no, nasal opening; nhp, naso-hypophyseal plate; ot, otic vesicle; ph, pharynx; pp,
 706 pharyngeal pouches; y, yolk; 3v, 4v, third and fourth ventricles, respectively. Scale bars, 0.5
 707 mm. Asterisk indicates a different expression has been found in a separate species, *P.*
 708 *marinus*, in which *Pm1Hox3* rostral limit is on the r4/r5 border²⁸.

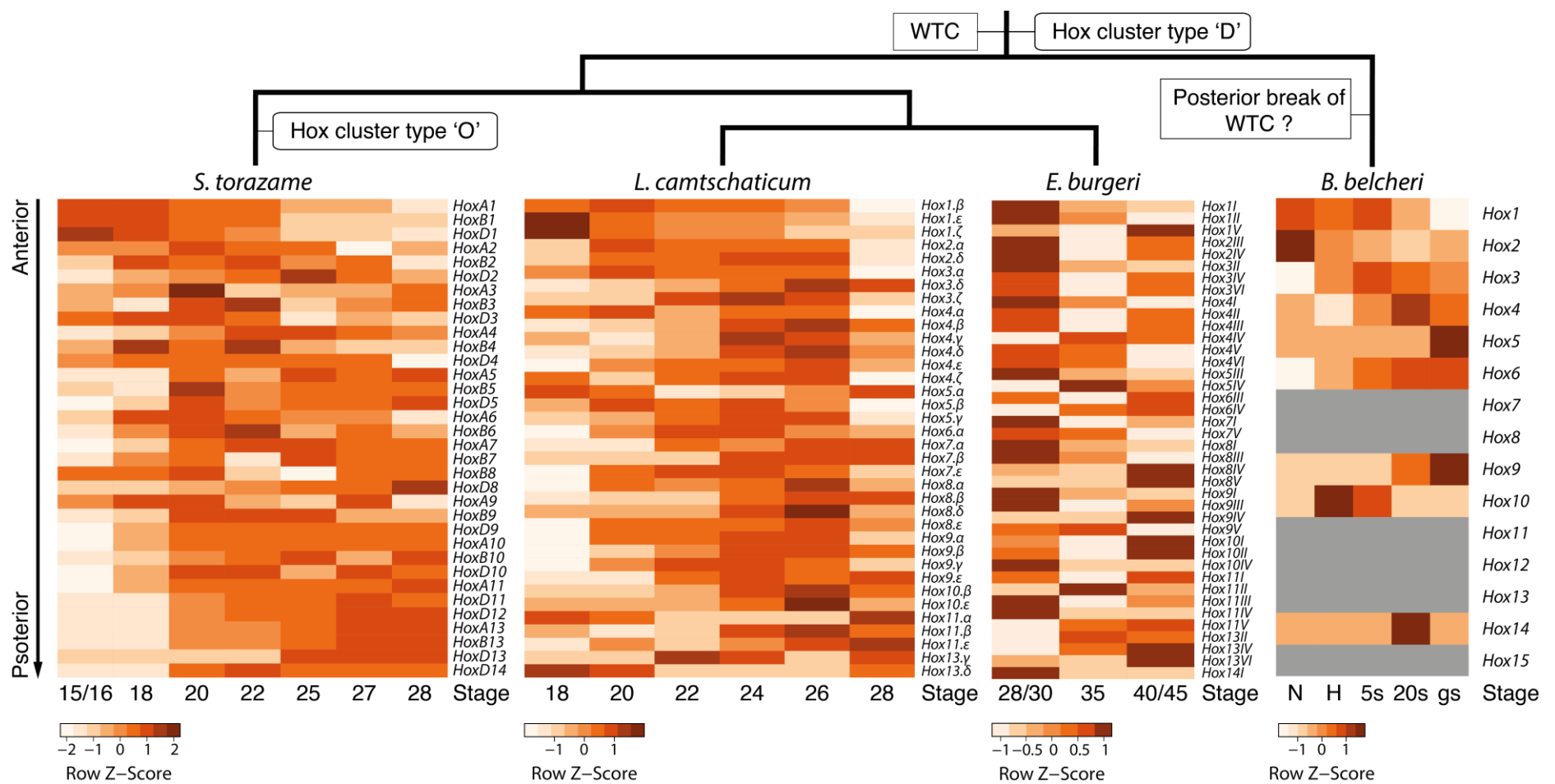
709

710 **Figure 4. Developmental expression profiling of Hox genes in chordates.** Heatmaps of
 711 Hox genes expression in *S. torazame* (gnathostome), *L. camtschaticum* and *E. burgeri*
 712 (agnathans), and *B. belcheri* (invertebrate chordate), coloured according to Z-score (standard
 713 deviations from mean expression level). Anterior Hox genes (top rows of heatmaps) tend to
 714 be expressed at higher levels at early stages of development than posterior genes (bottom
 715 rows of heatmaps) in both *S. torazame* and *L. camtschaticum*. On top, a phylogenetic tree
 716 with chordate relationships of the species studied here indicate the putative events that took
 717 place during evolution: in *B. belcheri*, temporal colinearity is appreciated between Hox1-5
 718 genes, indicating WTC was likely present in the last common ancestor of chordates, and a
 719 secondary escape of the posterior half of the cluster from it occurred independently in the
 720 amphioxus lineage. The large sizes of both amphioxus and agnathan Hox clusters implies that
 721 the common ancestor of vertebrates had a so-called ‘disorganized’ (D) cluster type, while the
 722 consolidation towards an ‘organized’ (O) type occurred in the gnathostome lineage⁴², after
 723 the split between jawed and jawless vertebrates. In *B. belcheri*, grey rows indicate genes with
 724 a FPKM value of 0 in all stages. N, amphioxus neurula stage; H, hatching stage; 5s, 5-somite
 725 stage; 20s, 20-somite stage; gs, 1- or 2-gill slit larvae.









729