

Reconstruction of the Vertebrate Ancestral Genome Reveals Dynamic Genome Reorganization in Early Vertebrates

Yoichiro Nakatani^{1, #}, Hiroyuki Takeda², Yuji Kohara³, Shinichi Morishita^{1, 4, #}

¹ Department of Computational Biology, Graduate School of Frontier Sciences, The University of Tokyo, Kashiwa 277-0882, Japan

² Department of Biological Sciences, Graduate School of Science, The University of Tokyo, Tokyo 113-0033, Japan

³ Center for Genetic Resource Information, National Institute of Genetics, Mishima 411-8540, Japan

⁴ Bioinformatics Research and Development (BIRD), Japan Science and Technology Agency (JST), Tokyo 102-8666, Japan

[#] Correspondence should be addressed to Y.N. (nakatani@cb.k.u-tokyo.ac.jp) and S.M. (moris@cb.k.u-tokyo.ac.jp).

Supplementary information

1	Introduction	1
2	Identification of ohnologs and orthologs	1
2.1	Identification of vertebrate ohnologs	1
2.2	Identification of teleost ohnologs and human–medaka orthologs	3
2.3	Identification of human–chicken orthologs	4
3	Identification of conserved vertebrate linkage (CVL) blocks	4
3.1	Identification of doubly conserved synteny (DCS) blocks	4
3.2	Combining fragmented DCS blocks into CVL blocks	5
3.3	Identification of the boundaries of CVL blocks	5
4	Refinement of CVL blocks	6
5	Ancestral vertebrate and gnathostome proto-chromosomes	7
6	Ancestral osteichthyan and amniote karyotypes	8
7	Effect of changing key parameter values on reconstruction of the vertebrate ancestral genome	8
8	Phylogenetic tree of vertebrates	9
9	Supplementary figures	10
10	Supplementary tables	59
11	References	62

1 Introduction

How is it possible to reconstruct ancestral vertebrate genomes without knowing the genomes of outgroup species? The fundamental principle of ancestral genome reconstruction is that the process requires three genomes, usually two descendent genomes and one outgroup genome. The three genomes do not have to be from different species. For example, extant teleost fishes have two copies of the ancestral teleost genome because of the WGD that occurred in the ancestral teleost. By comparing these two copies of the ancestral genome with the human genome, the ancestral teleost genome has been reconstructed in previous studies (Postlethwait et al. 2000; Naruse et al. 2004; Jaillon et al. 2004; Woods et al. 2005; Kohn et al. 2006). Similarly, the human genome contains four copies of the ancestral vertebrate genome produced by two rounds of WGD, which was sufficient to reconstruct the ancestral vertebrate genome at the second round of WGD. It might be difficult to reconstruct the ancestral vertebrate proto-chromosomes if numerous rearrangements occurred during the two rounds of WGD; however, in reality, only a few major interchromosomal rearrangements were observed between the two WGDs, allowing us to reconstruct the ancestral vertebrate proto-karyotype with only a few ambiguities. Alternative scenarios remained unresolved for three of the reconstructed ancestral vertebrate proto-chromosomes; nevertheless, our reconstruction was sufficient to reveal dynamic genome reorganization in early vertebrates.

2 Identification of ohnologs and orthologs

2.1 Identification of vertebrate ohnologs

Genes may be duplicated at many points in the course of evolution. The two rounds of whole-genome duplication (2R WGD) occurred after the divergence of urochordates and vertebrates, and before the divergence of ray-finned and lobe-finned fishes. Our goal was to identify ohnologs produced by 2R WGD while excluding other duplicated genes by inspecting duplication points in the phylogenetic tree. Supplementary Fig. S1A illustrates an ideal case in which the WGD events generated four and eight ohnologs in the human and medaka genomes, respectively, and all the ohnologs and their counterpart *Ciona* genes were preserved throughout evolution. The evolutionary distance between a pair of proteins is estimated by the inverse of the BLASTP raw score between the pair (Dehal et al. 2005). Assuming that recently duplicated genes are closer than pairs of older duplicates, the following properties hold in Supplementary Fig. S1A.

- The four human genes best match the same *Ciona* gene among all *Ciona* genes. Conversely, the four human genes to which the *Ciona* gene is most similar are the same four human genes. The four human genes are more similar to each other than to the *Ciona* gene because the *Ciona* gene diverged before the 2R WGD events.
- Each human gene has two orthologs in ray-finned fishes, such as medaka, that diverged after the 2R WGD events. Therefore, the human gene is closer to the ray-finned fish orthologs than to other human genes.

- None of the four human genes was copied by gene duplication after the divergence from ray-finned fishes.

These properties allow us to identify one group of four human ohnologs that corresponds to one *Ciona* gene in the ideal example in Supplementary Fig. S1A.

However, in practice, all ohnologs and their corresponding *Ciona* genes are rarely conserved due to the loss of genes, and Supplementary Fig. S1B illustrates such a difficult example. Here we show phenomena that disrupt the conservation in the order of evolution and we explain how ohnologs can be detected even in the presence of these problems.

- **Problem 1.** One serious problem is that the *Ciona* gene model is still incomplete; for example, *Ciona-2* is unknown while the similar *Ciona-1* gene has been sequenced, as indicated in Supplementary Fig. S1B. It then appears that the human genes numbered 1–8 best match *Ciona-1* and are therefore assigned to a single group, even though human-1–4 and human-5–8 should have been categorized into two separate groups. To resolve this problem, we further divide the single group into subgroups so that the distance between any two genes in a subgroup is bounded by the distance between *Ciona-1* and its best-matching human gene. This rule partitions the single group into two subgroups, correctly separating human-1–4 from human-5–8. One may be concerned that the divergence of gene evolution rates among *Ciona* genes, as well as human and medaka genes, could affect the analysis and produce false-positive groups of ohnologs. However, in the subsequent steps presented later, we will reduce the effect as possible by looking at neighboring ohnologs in synteny blocks because it is unlikely to observe a long series of erroneous ohnologs.
- **Problem 2.** Another serious problem is caused by the loss of genes. Supplementary Fig. S1B illustrates that several medaka genes were lost, making it difficult to determine that human-1 and 2 were duplicated before the split of humans and medaka. Specifically, this fact cannot be confirmed by checking that human-1 (or 2) is closer to a medaka ortholog than to human-2 because all medaka orthologs of human-1 and 2 were lost (see Supplementary Fig. S1C). This is a worst case scenario, and in most instances, we were able to recognize ohnologs by utilizing remaining medaka genes. For example, although medaka-7 was lost, we could determine that human-3 and -4 must have diverged before the split of medaka and humans because human-4 is closer to medaka-8 than to human-3. Similarly, human-8 is closer to medaka-14 than to human-7, indicating that they are ohnologs.
- **Problem 3.** The last major problem is that lineage-specific gene duplications create highly similar copies of a gene, for example, human-5 and -6 in Supplementary Fig. S1C. These duplicates are not treated as ohnologs because the distance between them is smaller than the distances to their medaka orthologs.

Based on the above findings, we implemented the steps of locating ohnolog candidates by generalizing the procedure proposed by Dehal and Boore (Dehal et al. 2005) to handle the above problems.

- **Step 1: Identification of vertebrate gene families.** Vertebrate protein sequences of human, mouse, dog,

chicken, *Tetraodon*, and *Takifugu* (version2) were obtained from Ensembl, and *Ciona* (version1) from JGI. After running an all against all BLASTP search, we used BLASTP hits with an E-value $<1e-10$ to identify vertebrate gene families that were expected to share common ancestral genes by categorizing vertebrate genes into distinct groups such that all genes in a group had the same best matching *Ciona* gene in terms of BLASTP raw scores. Therefore, each group had one representative *Ciona* gene at this point. However, to cope with Problem 1, each group was further divided, ensuring that the distance between any two genes in different subgroups was larger than the distance between the representative *Ciona* gene and its best-matching human gene in the group.

- **Step 2: Extraction of ohnologs created by 2R WGD from vertebrate gene families.** Individual vertebrate gene families should include ohnologs created by 2R WGD before the split of ray-finned and lobe-finned fishes, but they may also contain duplicated genes copied recently in particular lineages. To eliminate the latter class of lineage-specific duplicates, any two human genes were not considered ohnologous if they were closer to each other than to any medaka gene, which implied that they were copied after the medaka–human divergence. In addition to such non-ohnologous genes, we found that some groups contained numerous human genes generated by both gene duplications and the 2R WGD. Because it was difficult to distinguish ohnologs created solely by the 2R WGD process from others created the combined processes, we used only groups that included at most four duplicated human genes before the divergence of ray-finned and lobe-finned fishes, and treated human genes with no duplication after the divergence of ray-finned and lobe-finned fishes as ohnologs.

Considering the incompleteness of the *Ciona* genome, we also conducted the analysis described above using the sea urchin as an outgroup instead of *Ciona*, and combined the two ohnolog candidate sets. The collected ohnolog candidates are displayed along the individual human chromosomes in Supplementary Fig. S2. Our analysis essentially followed the method of Dehal and Boore (Dehal et al. 2005), and these two methods produced almost consistent paralogous regions across human chromosomes. We also compared several other methods developed in previous studies (Dehal et al. 2005; Dehal et al. 2006; Blomme et al. 2006) and found that phylogenetic tree construction was efficient in reducing erroneous identification of ohnologs; however, it also reduced the number of identified ohnologs among paralogous chromosomal regions. Since we developed a statistical reconstruction method that was robust against erroneously identified ohnologs, we omitted phylogenetic tree analysis to obtain a better reconstruction result as a whole.

2.2 Identification of teleost ohnologs and human–medaka orthologs

Following the procedure outlined in the previous section for detecting ohnologous human genes created by the 2R WGD, we attempted to identify ohnologous medaka genes created in the teleost WGD event. For this purpose, we used human genes as the outgroup, replacing *Ciona* genes, and we utilized pufferfish genes to test whether a pair of medaka genes were ohnologs or the result of gene duplication after the split of medaka and pufferfish. In Step 1, considering the substantial collection of human genes, which contrasted sharply

with the partial collection of *Ciona* genes, we categorized teleost genes into one group if they were best-matched with the same human gene. From a group of medaka genes, the gene with the highest similarity score to the outgroup human gene was selected as a representative. The representative gene qualified as a medaka ortholog for the outgroup human gene if the medaka genes in a group had at most one duplication event before the divergence of medaka and pufferfish. These medaka orthologs are also juxtaposed along the human chromosomes at the bottom of Supplementary Fig. S2.

2.3 Identification of human–chicken orthologs

Any human and chicken genes were orthologous if they were reciprocal best matches and neither of them had teleost genes with a higher similarity score. The chicken ortholog genes are also plotted in Supplementary Fig. S2.

3 Identification of conserved vertebrate linkage (CVL) blocks

Chromosomal segments in the vertebrate ancestor are distributed throughout the human genome due to intensive interchromosomal rearrangements (see Fig. 1). However, few interchromosomal rearrangements took place in the teleost lineages after the teleost WGD (Postlethwait et al. 2000; Naruse et al. 2004; Jaillon et al. 2004). Thus, although chromosomal segments in the vertebrate ancestor may have been broken into smaller segments and distributed over the human genome, their counterparts are likely to be highly preserved in the two (or three due to some chromosomal fissions) medaka chromosomes that were derived from the same chromosome in the ancestral teleost karyotype. Therefore, for detecting conserved vertebrate linkage (CVL) blocks, we attempted to identify blocks of human genes that had medaka orthologs on medaka chromosomes originating from the same chromosome in the ancestral teleost karyotype.

3.1 Identification of doubly conserved synteny (DCS) blocks

To this end, in the initial step, we utilized doubly conserved synteny (DCS) (Kellis et al. 2004; Jaillon et al. 2004) and identified the correspondence between human genes on one human chromosome and medaka orthologs on two duplicated medaka chromosomes by comparing the medaka and human genomes. Supplementary Fig. S3 summarizes the result. Some small DCS regions were added by manual inspection, but the ancestral teleost karyotype was unchanged. DCS is useful in providing an overview of correspondence between human and medaka chromosomes. For example, in Supplementary Fig. 2, human chromosome 1 has three major distinct DCS correspondences with the medaka genome, namely, medaka chromosomes 17-4, 5-7, and 11-16 (medaka chromosomes 11 and 22 also has a DCS correspondence as indicated in Supplementary Fig.S2, which is estimated to be a result of a translocation from 16 to 22 after the teleost WGD event). DCS does not immediately indicate CVL blocks; Supplementary Fig. 2 shows that each DCS region is not consecutive on human chromosome 1, but is partitioned into smaller blocks, and these

blocks should be recognized as CVL blocks.

3.2 Combining fragmented DCS blocks into CVL blocks

In principle, DCS regions should be divided if interrupted by inversions or translocations in the human lineage. However, DCS blocks are frequently fragmented into smaller blocks by translocations after the teleost WGD event in the medaka lineage. Here, we illustrate this problem, and afterwards we will present a solution to it. Supplementary Fig. S4A shows an ideal case for CVL block construction. In the green region of HSA6, most of the human genes are syntenic to the reconstructed teleost ancestor chromosome TEL-a because their medaka orthologs are found in its daughter chromosomes OLA22 and OLA24. This region constitutes one CVL block because it is not interrupted. However, interchromosomal translocations in the medaka lineage create some problems, and Supplementary Fig. S4B illustrates the effect of translocation in the medaka lineage. The four yellow genes were originally located in chromosome TEL-a of the teleost ancestor, but were translocated from OLA22 to OLA16 after WGD. Because of this translocation, the yellow human genes in the green region of HSA6 have medaka orthologs in OLA16 and were erroneously assigned to TEL-b. In this case, teleost ancestor synteny in the green region is interrupted by yellow genes even though the green region should be treated as one CVL block because it has no major rearrangements in the human lineage. Therefore, in the construction of CVL blocks, we need to let a small number of genes break into other CVL blocks to avoid partitioning them into smaller parts. However, allowing too many such synteny-interrupting genes may result in the merging of two unrelated CVL blocks.

To decide how many synteny-interrupting genes should be permitted, we investigated the size distribution of contiguous human genes that had orthologs on the same medaka chromosome (Supplementary Fig. S4C), observing that 4336 of 4356 contiguous synteny regions were smaller than 10 genes in size. A threshold smaller than 10 genes was likely to divide DCS regions into too many smaller CVL blocks, while a larger threshold often failed to identify small CVL blocks inside DCS regions. Therefore, we decided to divide a DCS region into two blocks if two proximate human genes in the DCS region were interrupted by at least 10 genes belonging to different DCS regions, e.g., in Supplementary Fig. S4B, the green region becomes one CVL block because the maximum size of contiguous synteny-interrupting genes is two. If the interrupting region with at least 10 genes is a DCS region, it is treated as a single CVL block as illustrated in Supplementary Fig. S5A. Stated another way, the minimum number of genes in a CVL block is set to ten. One may wonder how the reconstruction is affected if the minimum threshold is changed. We will discuss this issue later.

3.3 Identification of the boundaries of CVL blocks

After the division into CVL blocks, the boundary between neighboring CVL blocks was not always clear; rather, two neighboring blocks often overlapped and some genes were mixed in the boundary region (Supplementary Fig. S5A). Furthermore, many genes (depicted as white boxes) did not have orthologs in the

medaka genome probably because counterpart genomic regions were deleted, or massive mutations produced pseudogenes or genes of different functions. We considered putting unassigned genes surrounded by a pair of “assigned” genes on a CVL block into the same CVL block, but the region between the pair could contain some genes associated with teleost ancestor chromosomes other than that of the CVL block, as well as unassigned genes. To properly eliminate this noise, we applied the condition that the unassigned genes are put into the CVL block if the region involves only unassigned genes except for at most one gene assigned to a medaka ortholog that is not mapped to any teleost ancestor chromosome. The exception was tolerated because about 90% of the medaka genome is covered by mapped scaffolds, and small-scale interchromosomal translocations in the medaka lineage may have put such an unassigned gene into the CVL.

4 Refinement of CVL blocks

Here we discuss a serious issue that may arise during the construction of CVL blocks, and we present a solution to the problem. The major problems in the reconstruction of CVL blocks are genome rearrangements that took place after the 2R WGD in the ancestral vertebrate genome and before the divergence of ray-finned and lobe-finned fishes. If the genome of a cartilaginous fish were available, it would provide valuable information to resolve this problem. These rearrangements are still present in the human and fish genomes, making it difficult to reconstruct the ancestral vertebrate proto-karyotype. Supplementary Fig. S6A illustrates two fission events that took place prior to the osteichthyan ancestor. The four chromosomal fragments originating from these fissions survived as CVL blocks, and CVL blocks derived from the same vertebrate proto-chromosome are correctly assigned to one connected component in the CVL graph. Then, these blocks are combined into ancestral gnathostome chromosomes as described in the Methods section of the main text. For example, B4L have many common ohnologs with BL1 as well as B1R, but it shares none with B4R, indicating that a fission event break one proto-chromosome into B4L and B4R. In this case, no need exists to refine the CVL blocks.

However, chromosomal fusions that occurred before the osteichthyan ancestor are likely to have serious consequences, that is, undesirable CVL blocks with two fragments originating from multiple distinct proto-chromosomes, as shown in Supplementary Fig. S6B. Such unqualified CVL blocks could also be produced by numerous chromosomal rearrangements occurring independently in distinct ray-finned and lobe-finned fish lineages (Supplementary Fig. S6C). Creating such improper CVL blocks should be avoided, but they can only be detected after the reconstruction of the ancestral vertebrate karyotype by checking whether they have a significantly great amount of ohnologs in more than one vertebrate proto-chromosome. Supplementary Fig. S6D illustrates how CVL block A4+B1L can be identified as a fused CVL block. In this case, CVL block A4+B1L is paralogous to A1, A2, and A3, so it is assigned to the red connected component. Since the B1L region has few ohnologs in A1–A3, but several ohnologs in B2–B4L, the B1L region can be identified by mapping red and green ohnologs to CVL block A4+B1L. Specifically, each CVL block is checked for whether ohnologs from two connected components are distributed nonrandomly over the CVL

block by conducting a Mann–Whitney *U*-test. If the two-tailed probability is <0.01 , the CVL block is divided into two blocks. For example, in Supplementary Fig. S6D, genes between two red ohnologs that are not interrupted by green ohnologs are assigned to new CVL block, and vertebrate proto-chromosomes can be reconstructed correctly.

Supplementary Fig. S6E shows a more complicated case of fused CVL blocks, in which A4+B1L is paralogous to red and green CVL blocks, and one large connected component corresponding to two ancestral vertebrate proto-chromosomes is obtained. In this case, CVL block A4+B1L cannot be identified as a fused CVL block as in Supplementary Fig. S6D because A4 and B1L have ohnologs in the same connected component. To avoid making such erroneous connected components, we checked whether a connected component can be divided into two subcomponents that share significantly fewer ohnologs. If the probability of finding fewer ohnologs by chance (see Methods) was less than 0.01 (this parameter value did not affect the result if it is changed to 0.1 or 0.001), the CVL block was identified as a fused CVL block and its edges were removed from the CVL graph except for the edge to the most significantly paralogous CVL block. Subsequently, we conducted the CVL block refinement step in Supplementary Fig. S6D. In our analysis, for example, CVL block #63 was identified as a fused CVL block, and divided into new CVL blocks #63 and #113. This division is validated independently by chicken synteny since CVL blocks #63 and #113 have orthologs in different chicken chromosomes (see Supplementary Fig. S2, human chromosome 11). After dividing the CVL blocks, we constructed CVL graphs and reconstructed ancestral vertebrate proto-chromosomes.

In the initial step, 109 CVL blocks were generated, and subsequently, the refinement process produced a total of 118 CVL blocks. Striking examples can be seen in human chromosome 17, in which CVL blocks #82 and #84 were divided into two parts that were consistent with the break points of human–chicken synteny (Supplementary Fig. S2, human chromosome 17).

5 Ancestral vertebrate and gnathostome proto-chromosomes

Reconstructed ancestral vertebrate proto-chromosomes are listed in Supplementary Fig. S7. For each ancestral vertebrate proto-chromosome, the most significant five reconstruction candidates are shown at the top of the table; the left column indicates significance in terms of probability and the remaining columns indicate CVL blocks constituting individual sister chromosomes. In the Methods of the main text, we have described how to compute the significance of reconstructed gnathostome proto-chromosomes in order to select the optimum one.

Reconstruction using CVL blocks with a smaller number of genes is less reliable than using larger CVL blocks, which makes it difficult to rebuild smaller proto-chromosomes. In contrast, larger proto-chromosomes with more ohnologs are more reliable. One effective way to reconfirm the reconstruction is to examine how CVL blocks in one proto-chromosome are clustered in the ancestral teleost karyotype and the chicken genome. For example, the largest sister chromosome of proto-chromosome C consists of seven

CVL blocks and all of the blocks are in the same teleost ancestor chromosome k and chicken chromosome 1, supporting the ancestral linkage of the seven CVL blocks.

6 Ancestral osteichthyan and amniote karyotypes

In the reconstruction of osteichthyan and amniote ancestors, synteny regions in the chicken genome for CVL blocks must be identified. We joined the genes constituting a CVL block into synteny blocks by applying Bourque *et al.*'s "gene7" method (Bourque et al. 2005) with an additional restriction that synteny blocks should not be extended outside of the CVL block. Specifically, we joined two orthologs in a CVL block if they satisfy the conditions of "gene7."

"Two genes A and B are joined together if there are up to two intervening genes between them in every species, with certain constraints on flipping A and B: at two intervening genes, the relative orientations of A and B must be the same in all species or one of them can be flipped, and with less than two intervening genes, either or both can be flipped. Next, we discarded blocks supported by less than three genes." (Bourque et al. 2005)

Then, we applied the 2-of-3 rule to reconstruct the ancestral chromosomes. CVL blocks in the chromosomes of vertebrate, osteichthyan, and amniote ancestors, and synteny block size between CVL blocks and chicken chromosomes are listed in Supplementary Tables S1–S3.

7 Effect of changing key parameter values on reconstruction of the vertebrate ancestral genome

The following three parameters are essential to reconstruct CVL blocks and the vertebrate ancestral genome:

1. the threshold on the number of genes in a CVL block (see Subsection 3.2 in the supplementary document),
2. the significance threshold for testing if two CVL blocks are paralogous (see Methods in the main text), and
3. the significance threshold for the Mann-Whitney U-test to decide whether a CVL block is divided (see Section 4 in the supplementary document).

The default values of individual parameters are 10, 1E-4, and 1E-2, respectively. To see the effect of these parameters on our analysis, we reconstructed the vertebrate ancestral genome with parameter values that were lower or higher than the default values. Supplementary Table S4 presents the results when the first parameter is set to 7, 10 and 13, the second parameter to 5E-4, 1E-4 and 5E-5, and the third parameter to 5E-2, 1E-2 and 1E-3. Although the number of CVL blocks is somewhat affected by the change of parameter values, the numbers of vertebrate groups (or, proto-chromosomes) in vertebrate, gnathostome, osteichthyan, and

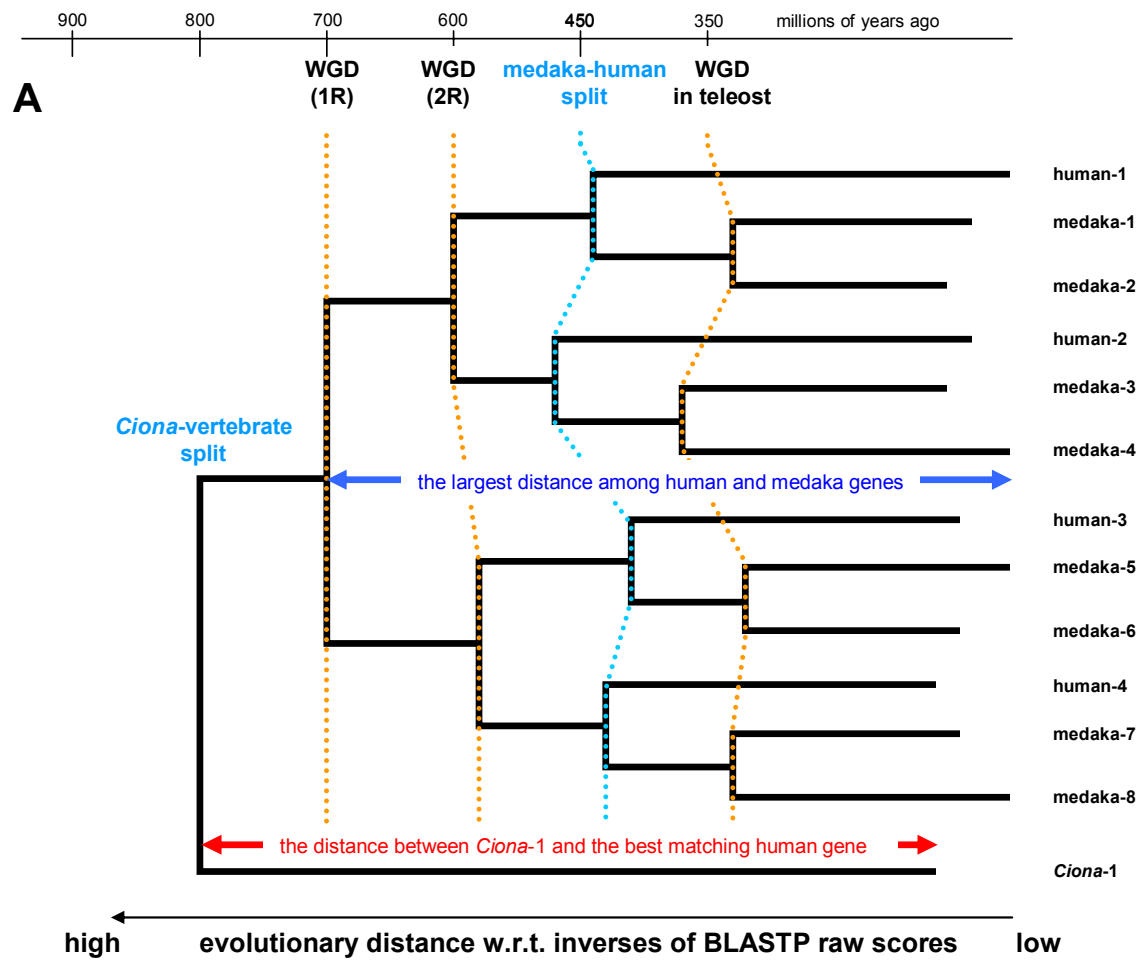
amniote ancestors are almost stable and are consistent with respective numbers in the default setting of parameters, thereby reconfirming our scenario of karyotype evolution. Minor changes are observed in Cases 3, 4, and 6. In Case 3, the first parameter is set to 13, which is higher than the default value 10. The setting is likely to merge fragmented DCS blocks into CVL blocks, producing five gnathostome proto-chromosomes for vertebrate proto-chromosome C. In Cases 4 and 6, the vertebrate proto-chromosomes C and F are fused into one.

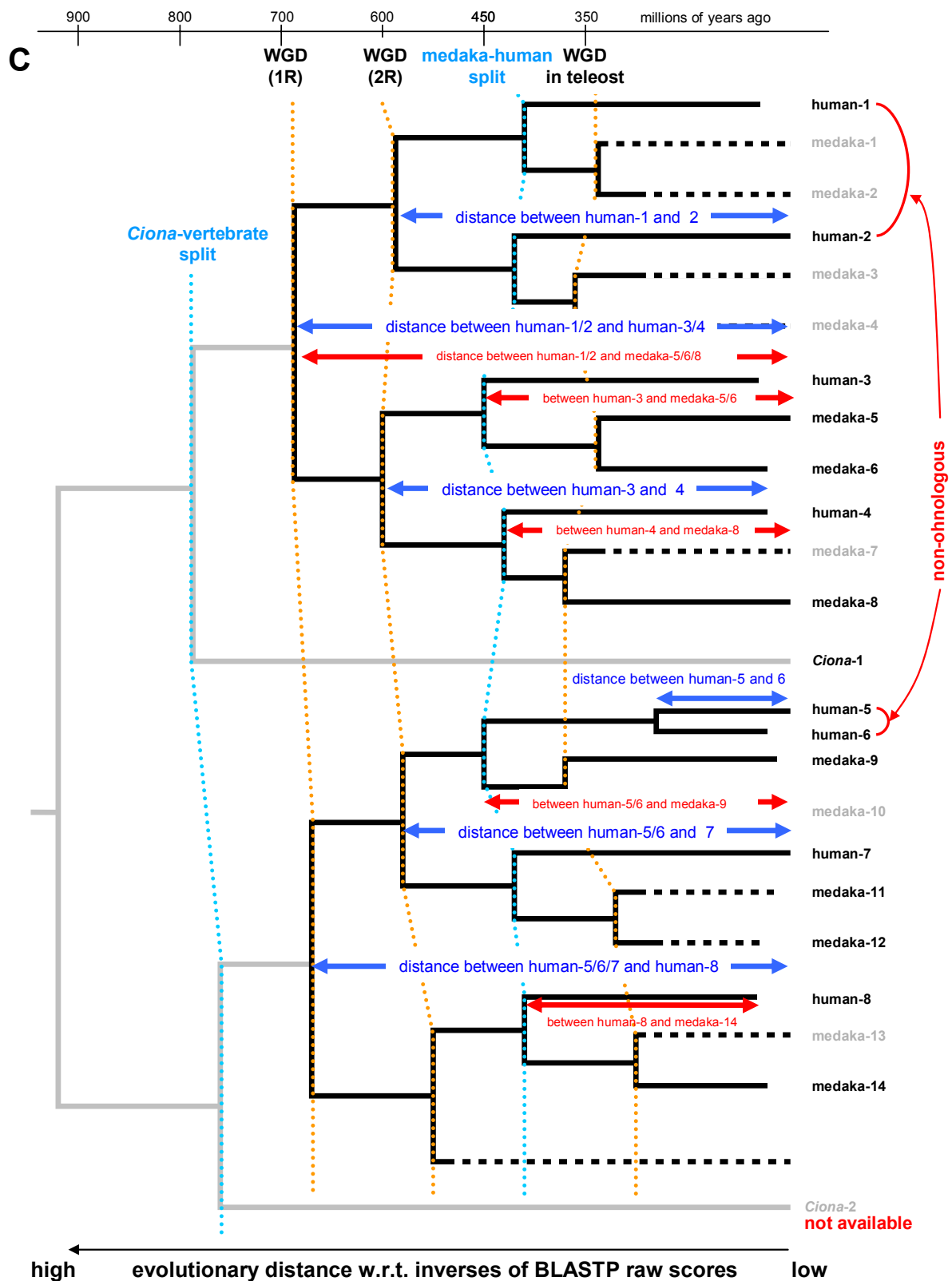
8 Phylogenetic tree of vertebrates

The phylogenetic tree in Fig. 6 is based on information presented in previous reports. We referred to Yamanoue *et al.* (Yamanoue et al. 2006) for the divergence times of torafugu, spotted green pufferfish, medaka, zebrafish, and the sarcopterygian–actinopterygian split, and to Inoue *et al.* (Inoue et al. 2005) for the divergence times of bichir, sturgeon, paddlefish, gar, and bowfin. We referred to Hedges and Poling (Hedges et al. 1999) for Reptilia, to Woodburne *et al.* (Woodburne et al. 2003) for Monotremata and Metatheria, to Springer *et al.* (Springer et al. 2003) for Mammalia, and to Blair and Hedges (Blair et al. 2005) for the rest, although some of the phylogenetic relationships and divergence times presented remain controversial (Meyer et al. 2003; Benton et al. 2007). The phylogenetic timing of the two rounds of whole-genome duplication in the vertebrate ancestor is cited from Stadler *et al.* (Stadler et al. 2004) and in the teleost ancestor, Hoegg *et al.* (Hoegg et al. 2004) and Crow *et al.* (Crow et al. 2006). The right column in Fig. 4 shows distributions of chromosome numbers up to $2n=100$; some species have more chromosomes, but were not included because of space limitations. Supplementary Fig. S8 shows the complete distributions. Chromosome number data were obtained from the Animal Genome Size Database [Gregory (2006): <http://www.genomesize.com>].

9 Supplementary figures

Supplementary Figure S1

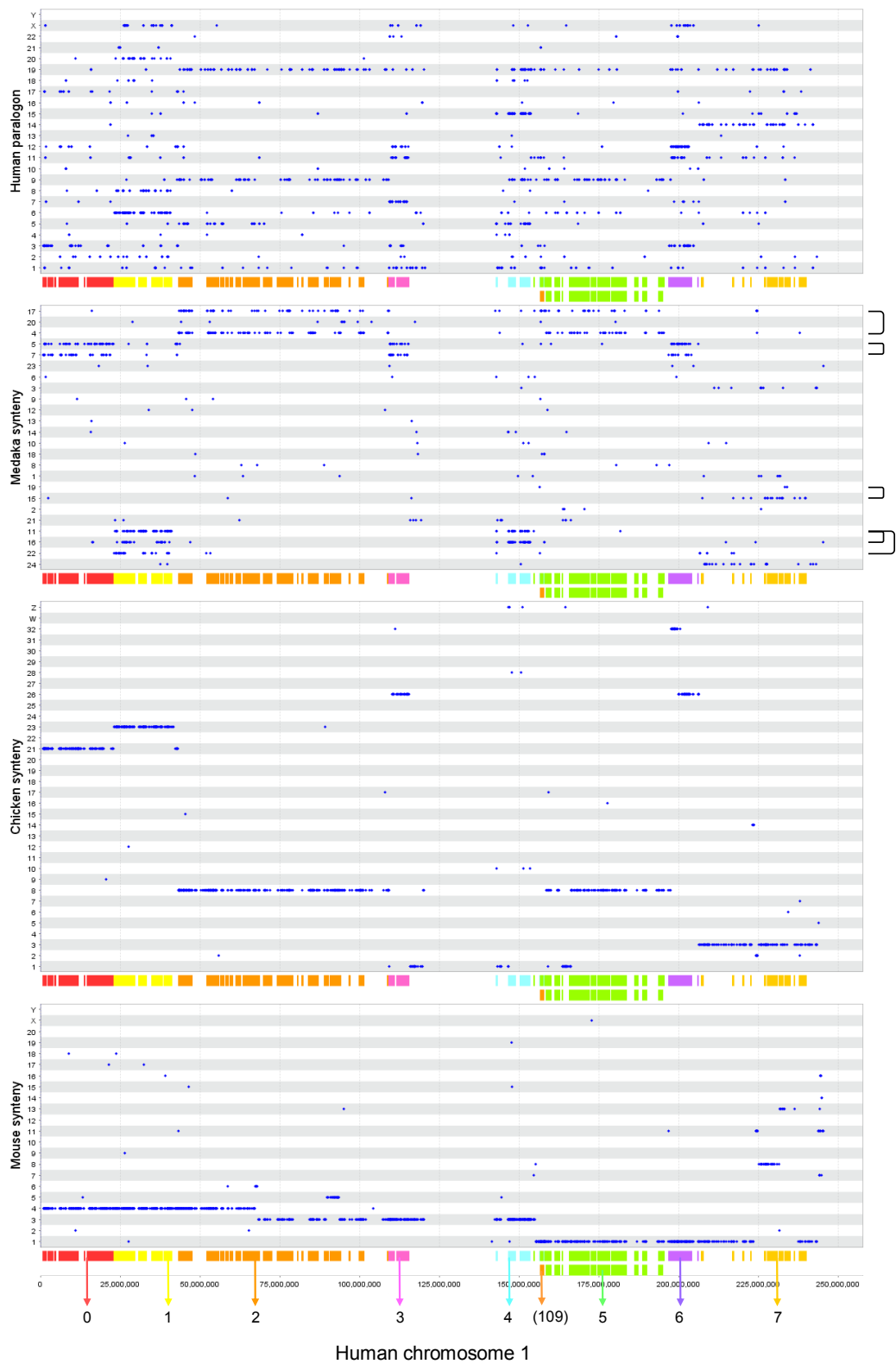


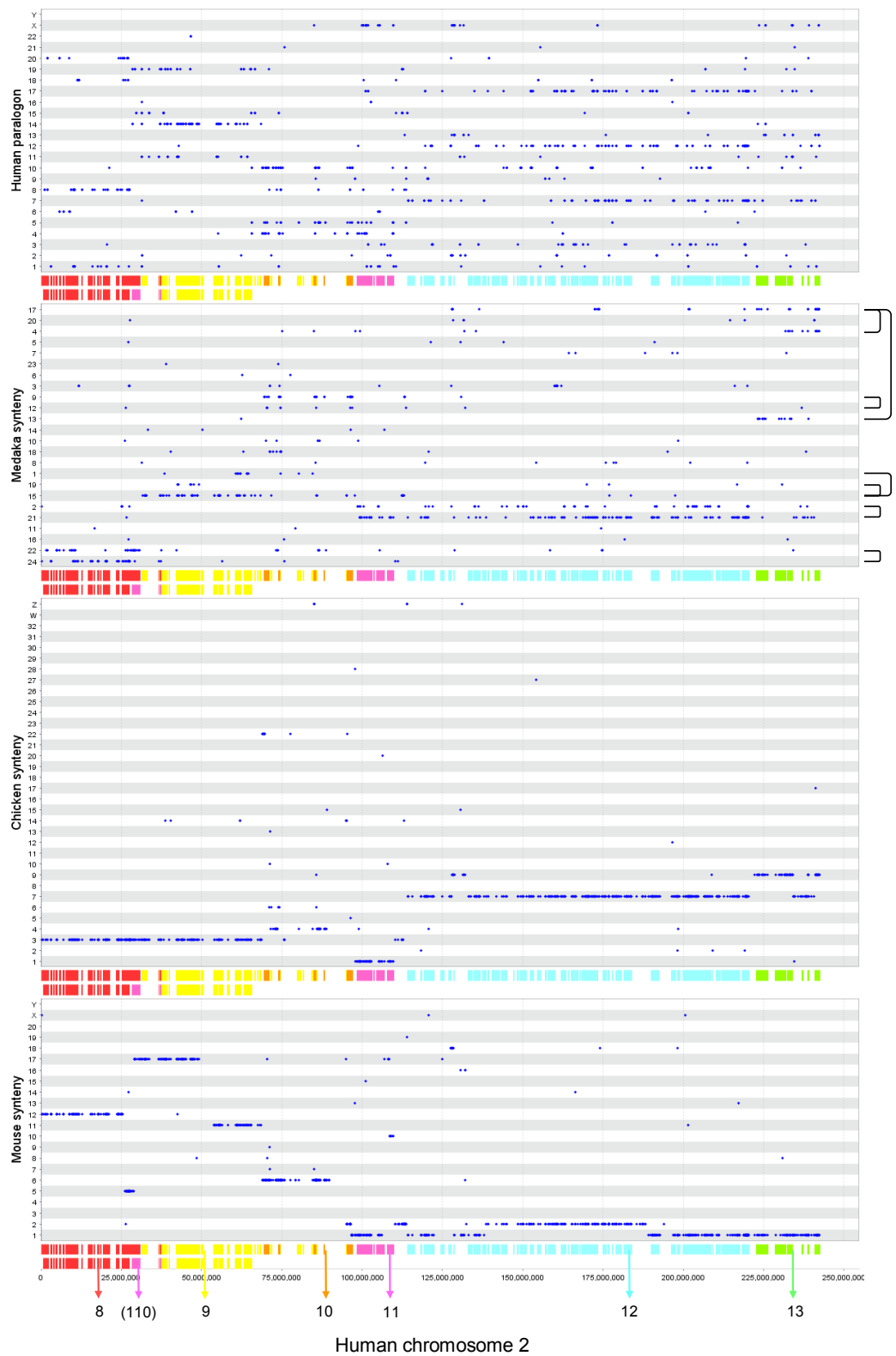


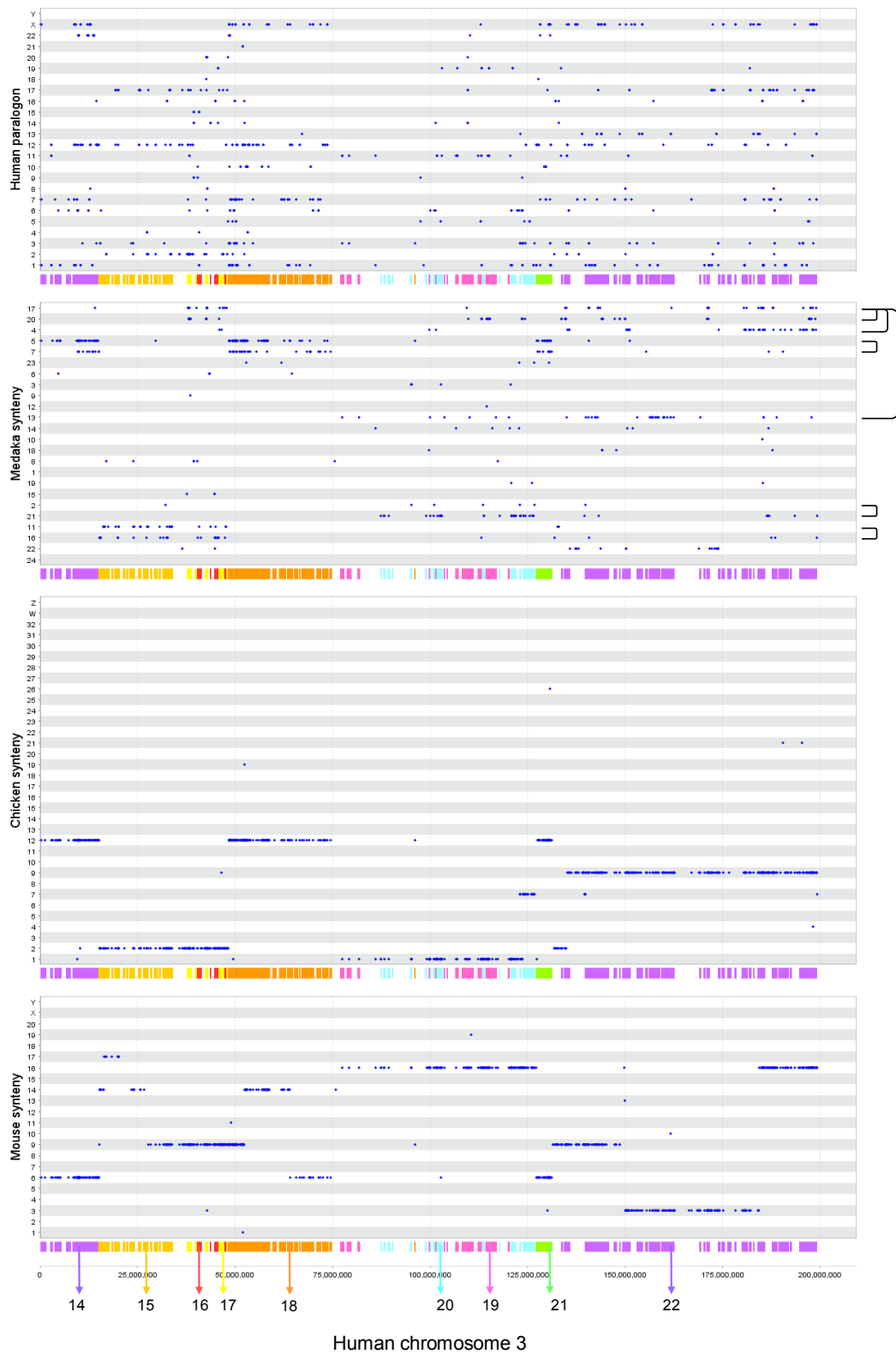
Phylogenetic tree of human, medaka, and *Ciona* genes. Evolutionary distances were calculated as inverses of BLASTP raw scores between pairs of two genes. **A.** The two rounds of whole genome duplication events in the vertebrate ancestor and the teleost whole genome duplication generated four human and eight medaka ohnologs. **B.** Several medaka genes were lost in the medaka lineage, one chromosome was lost before the medaka-human split, and *Ciona*-2 has not yet been sequenced. All human and medaka genes were most similar to *Ciona*-1 and were therefore temporarily categorized into one group. **C.** This temporary group was divided into two subgroups generated after the *Ciona*–vertebrate split such that the distance between any genes in individual subgroups is bounded by the distance between *Ciona*-1 and the best matching human gene.

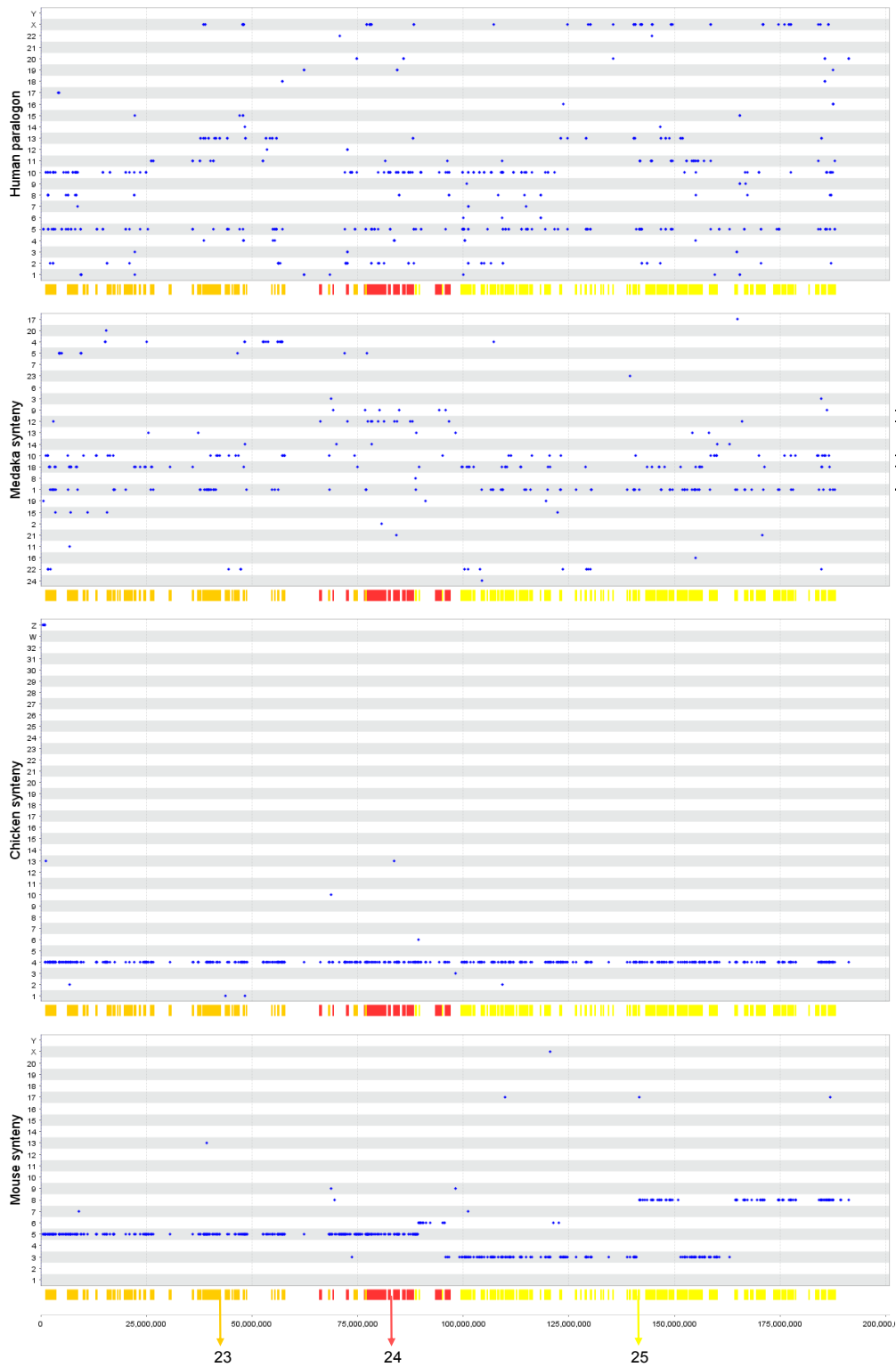
Supplementary Figure S2

Human ohnologs and orthologs of medaka, chicken, and mouse. Genes in one CVL block have the same color. CVL blocks in the upper row indicate the original 109 blocks (see “Identification of CVL blocks”), while those in the lower row, indicate the refined blocks (see “Refinement of CVL blocks”). These blocks are displayed along individual human chromosomes from the p to the q telomere and are numbered from 0 to 117. Those numbered from 109 to 117 were isolated in the CVL block refinement step, and these numbers are enclosed in parentheses to highlight this modification. Duplicated medaka chromosome pairs are connected by lines on the right side of the medaka ortholog plot.

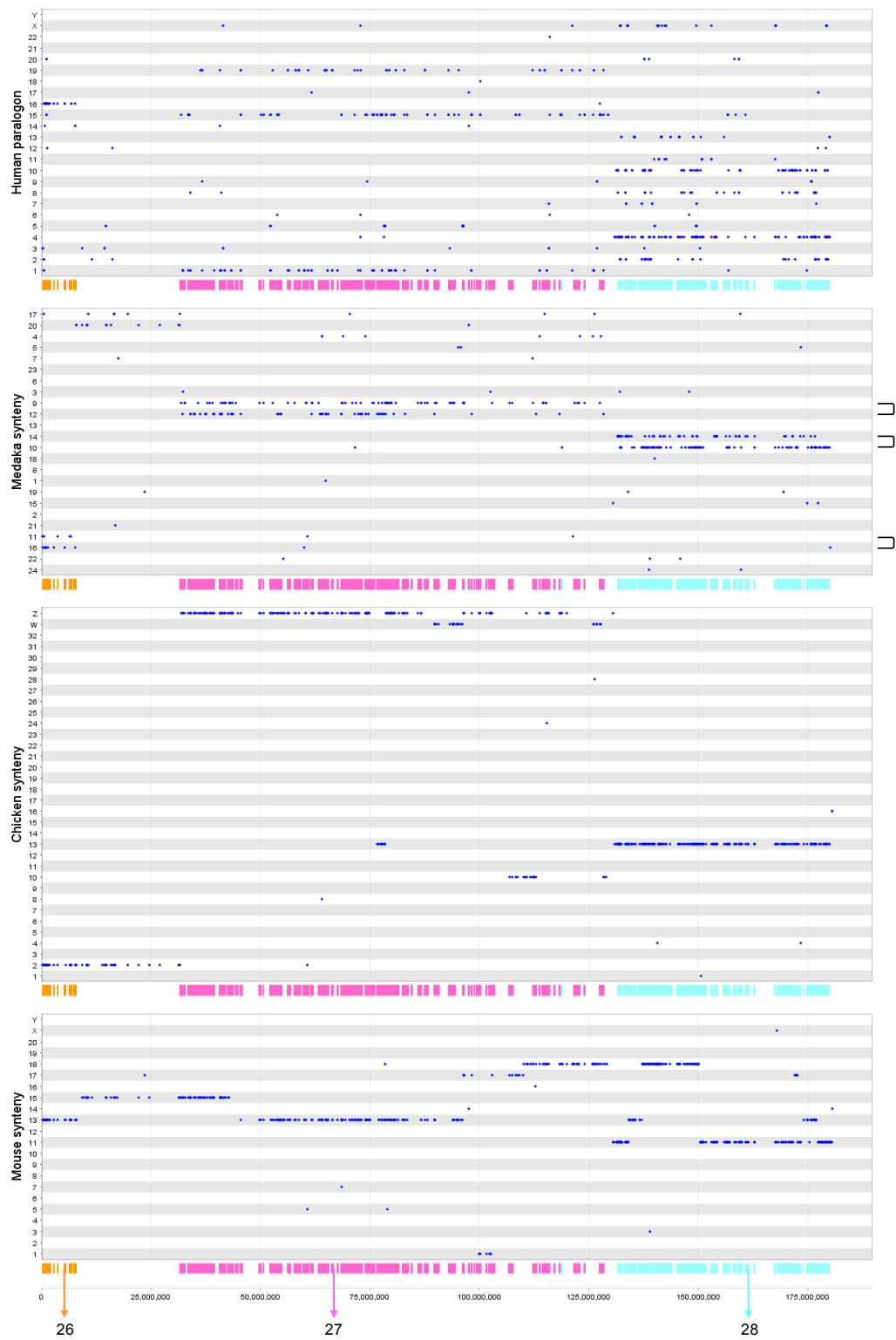




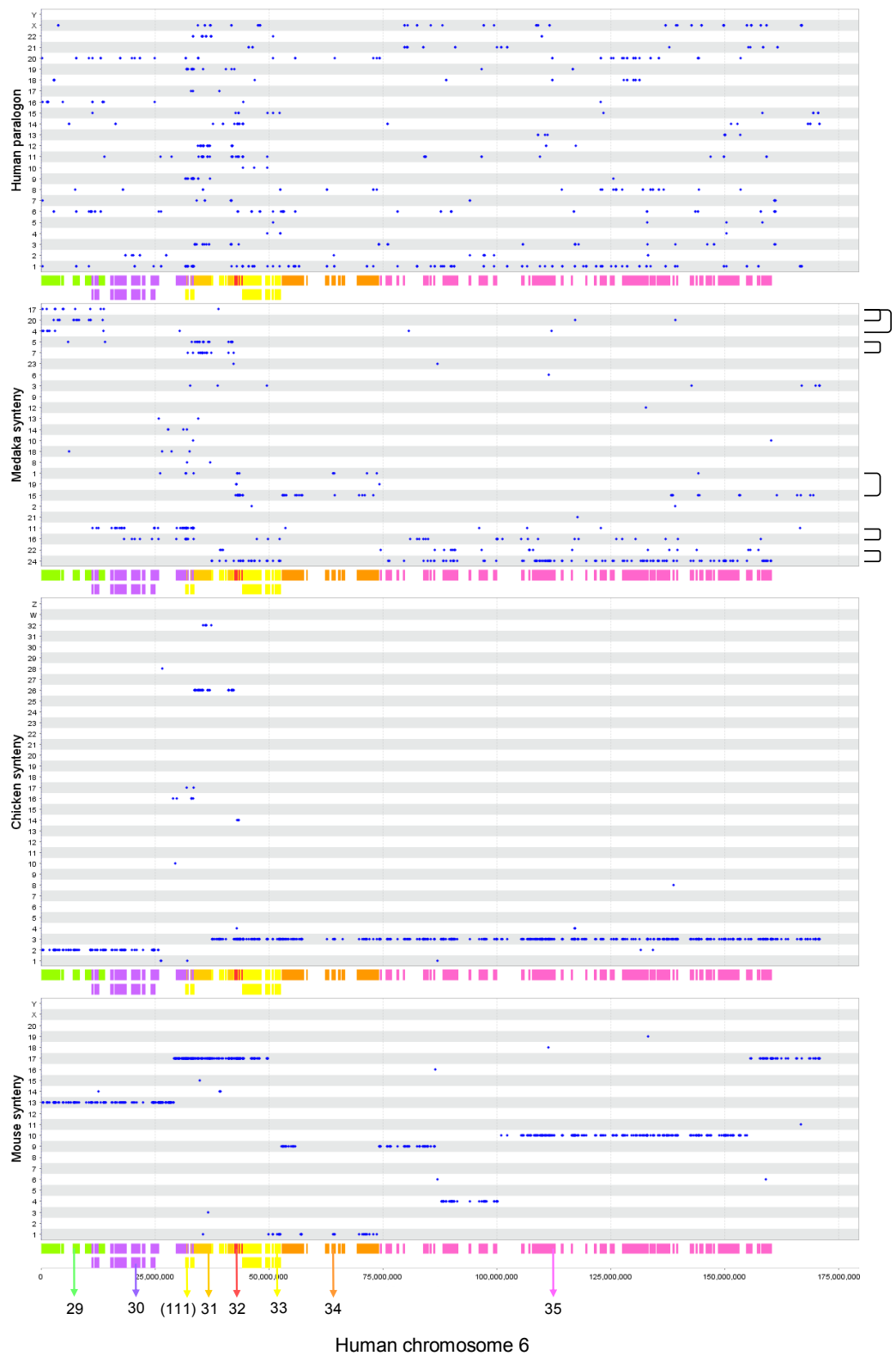


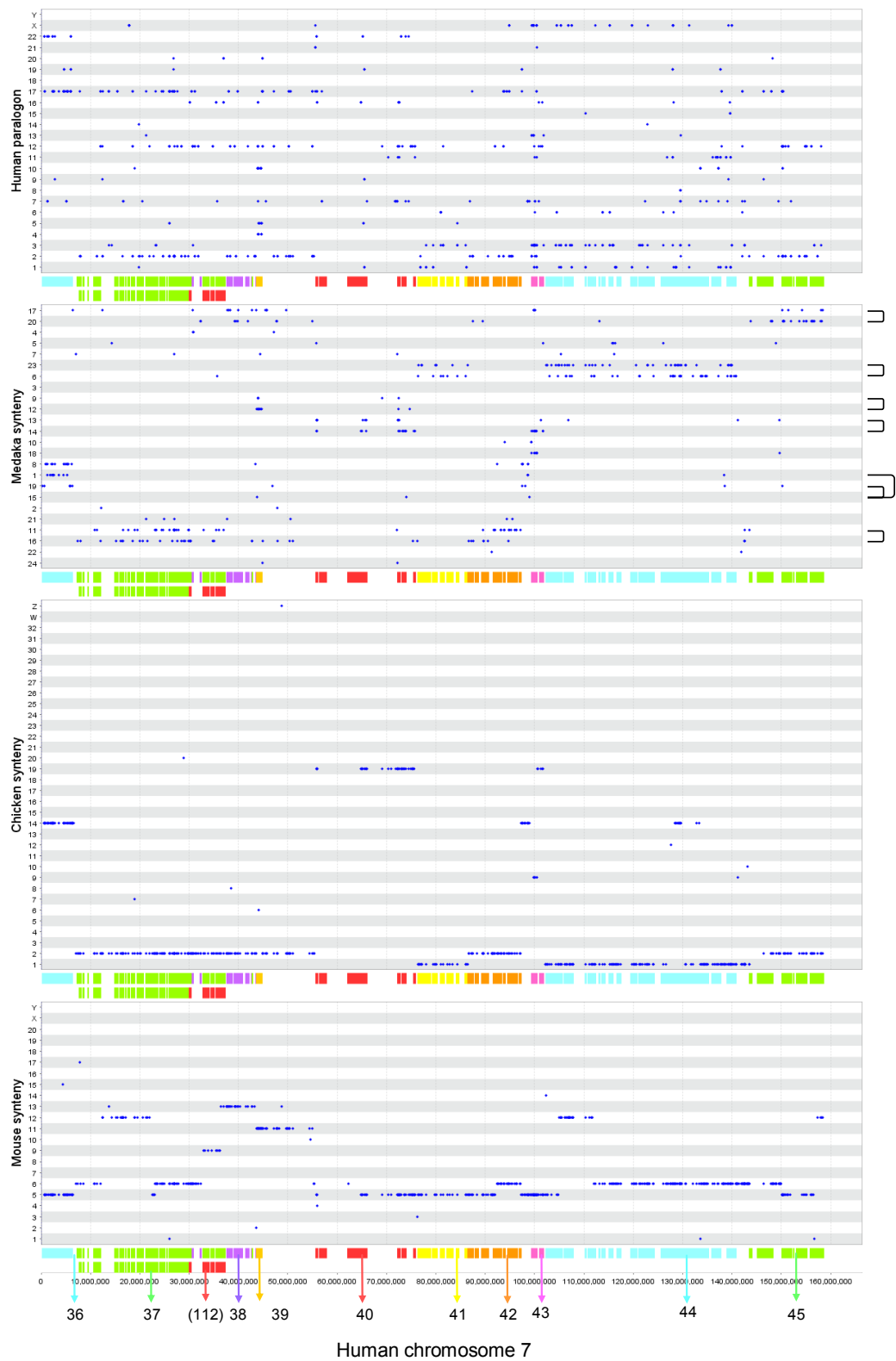


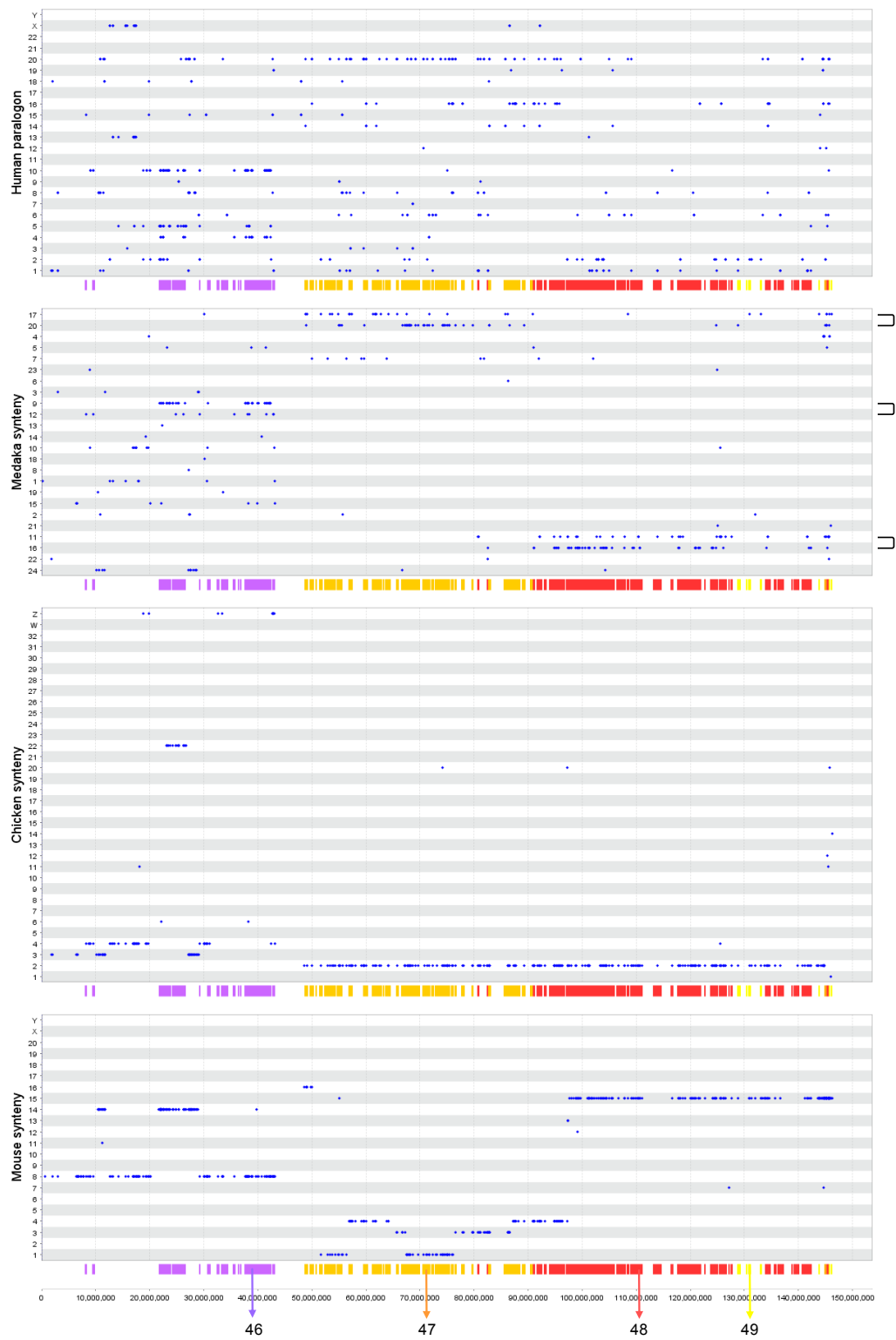
Human chromosome 4



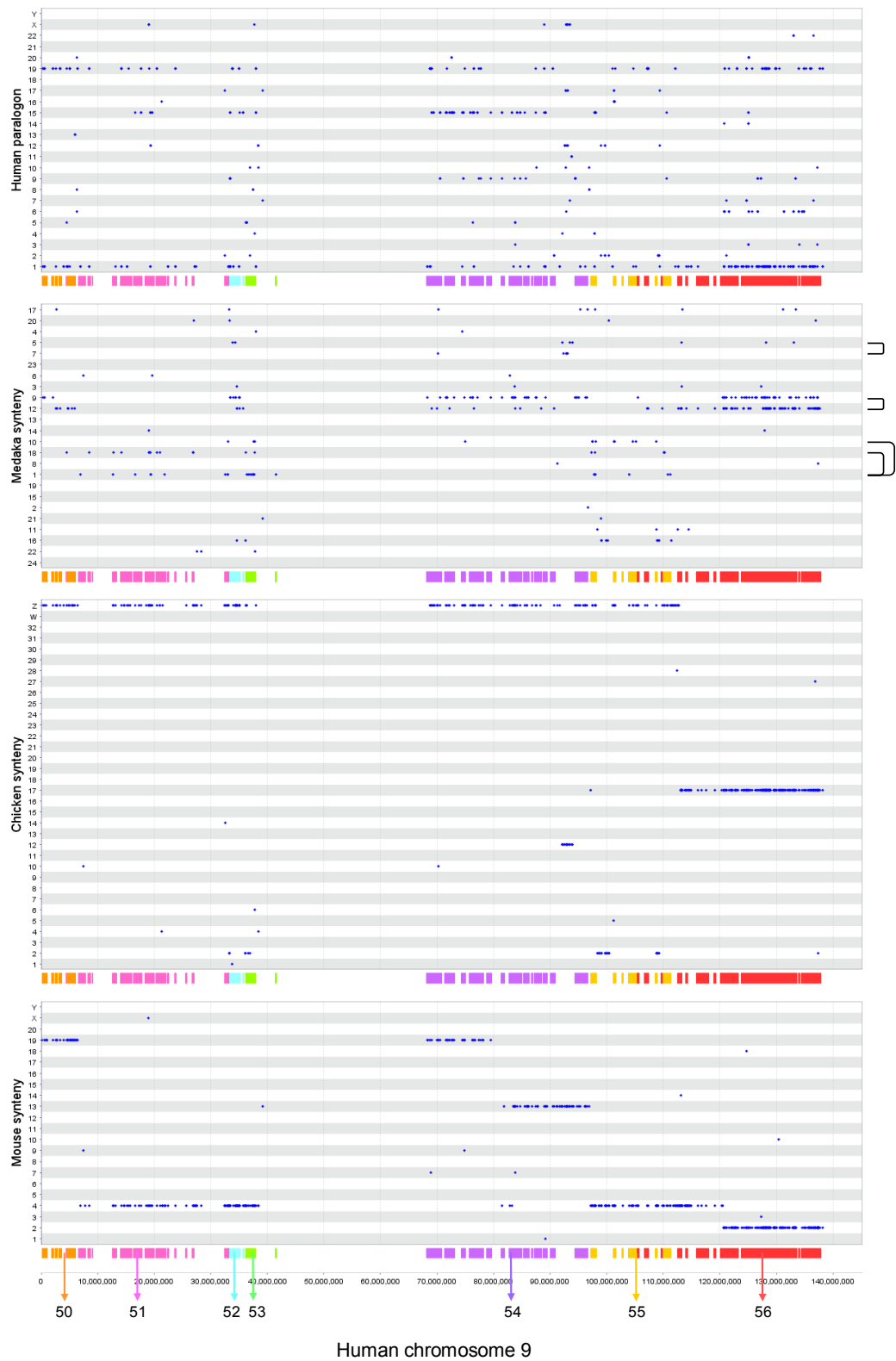
Human chromosome 5

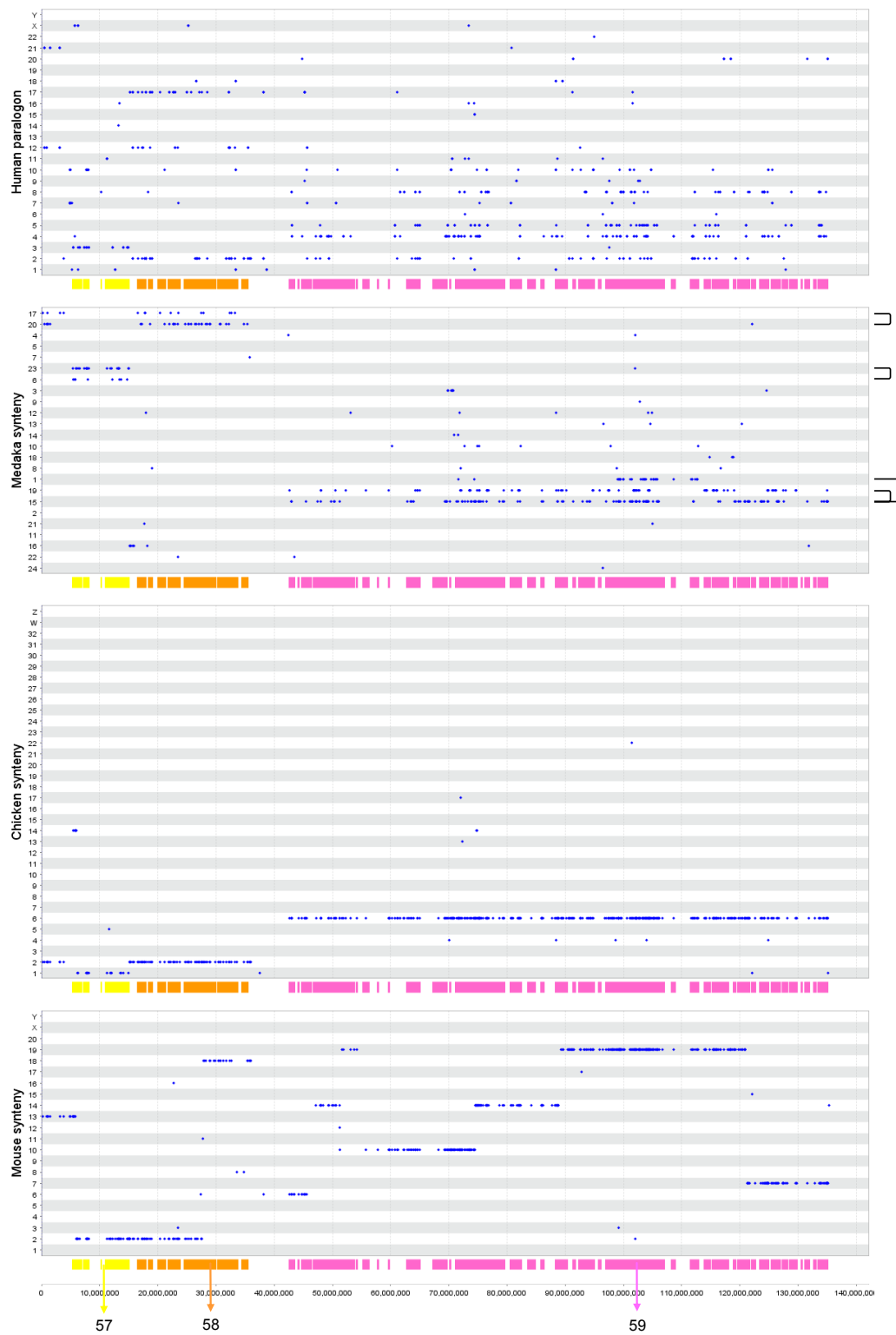




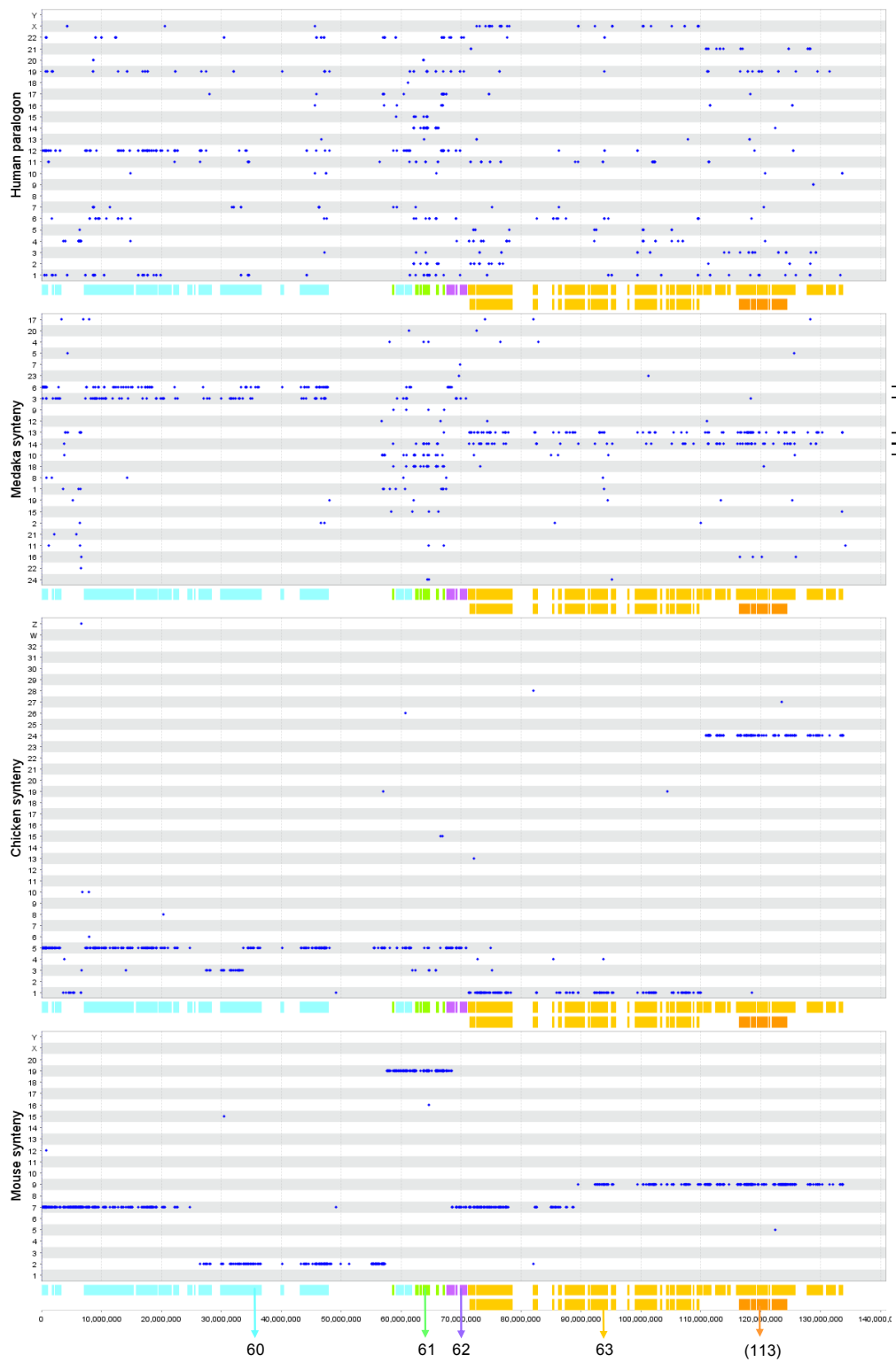


Human chromosome 8

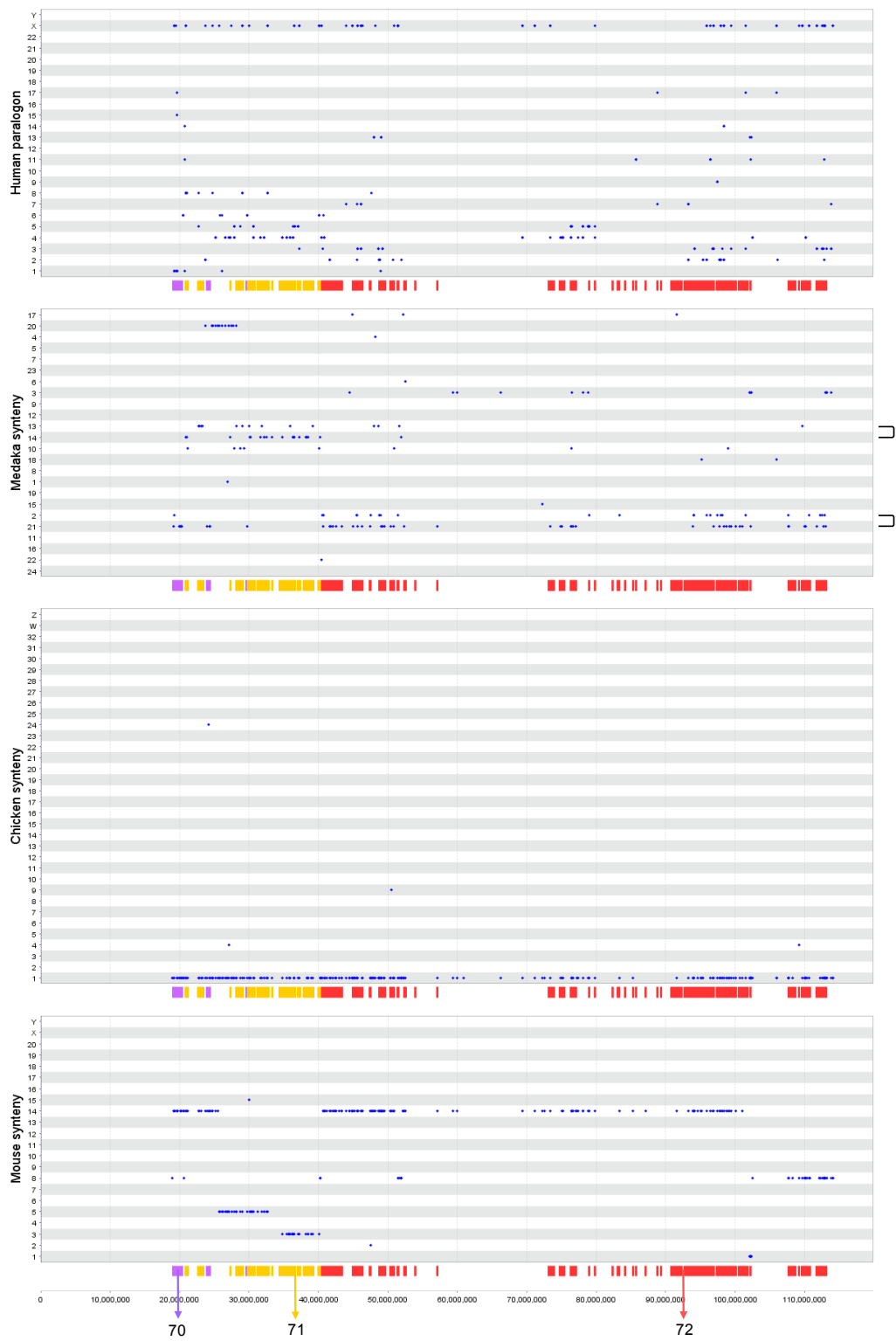




Human chromosome 10



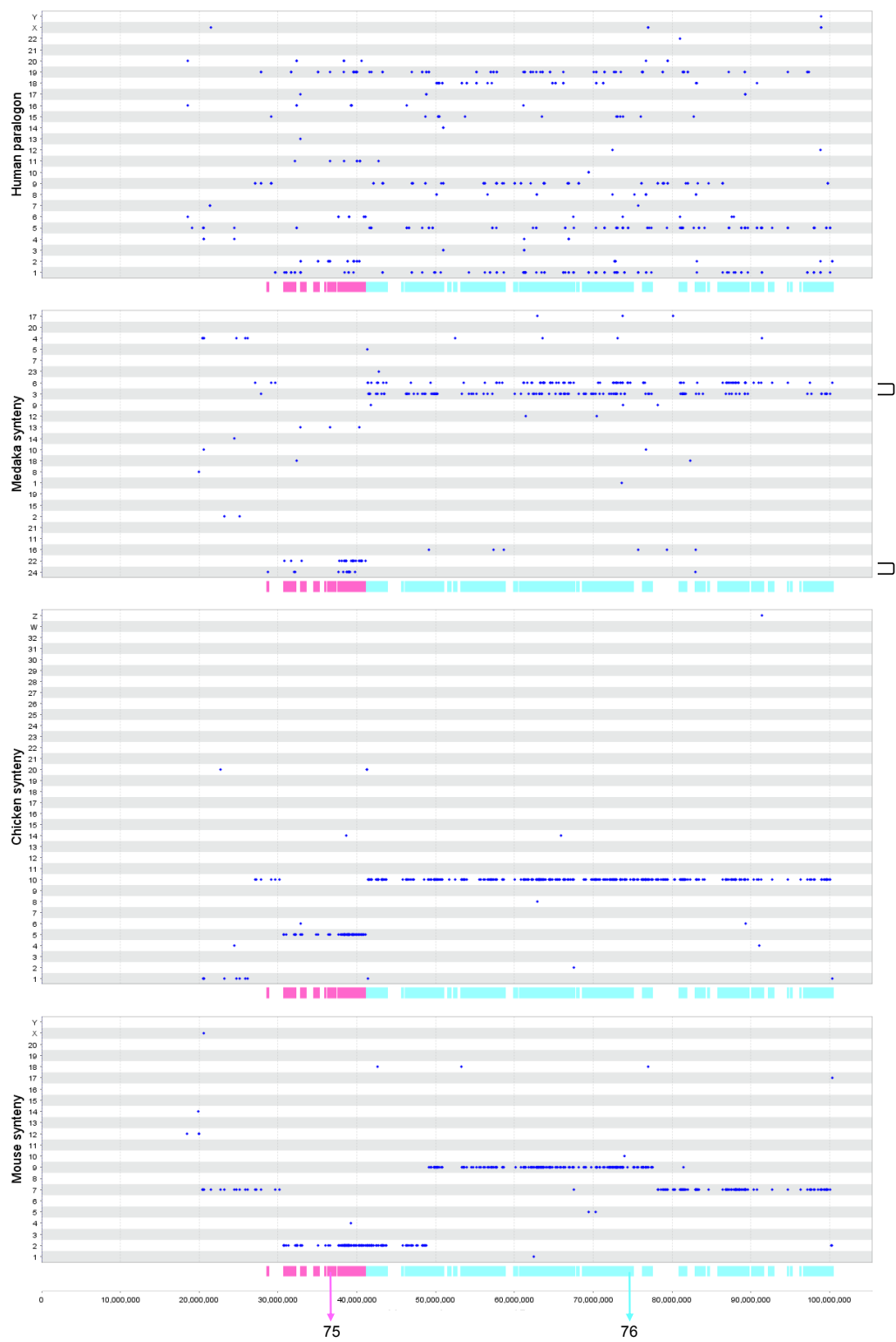
Human chromosome 11



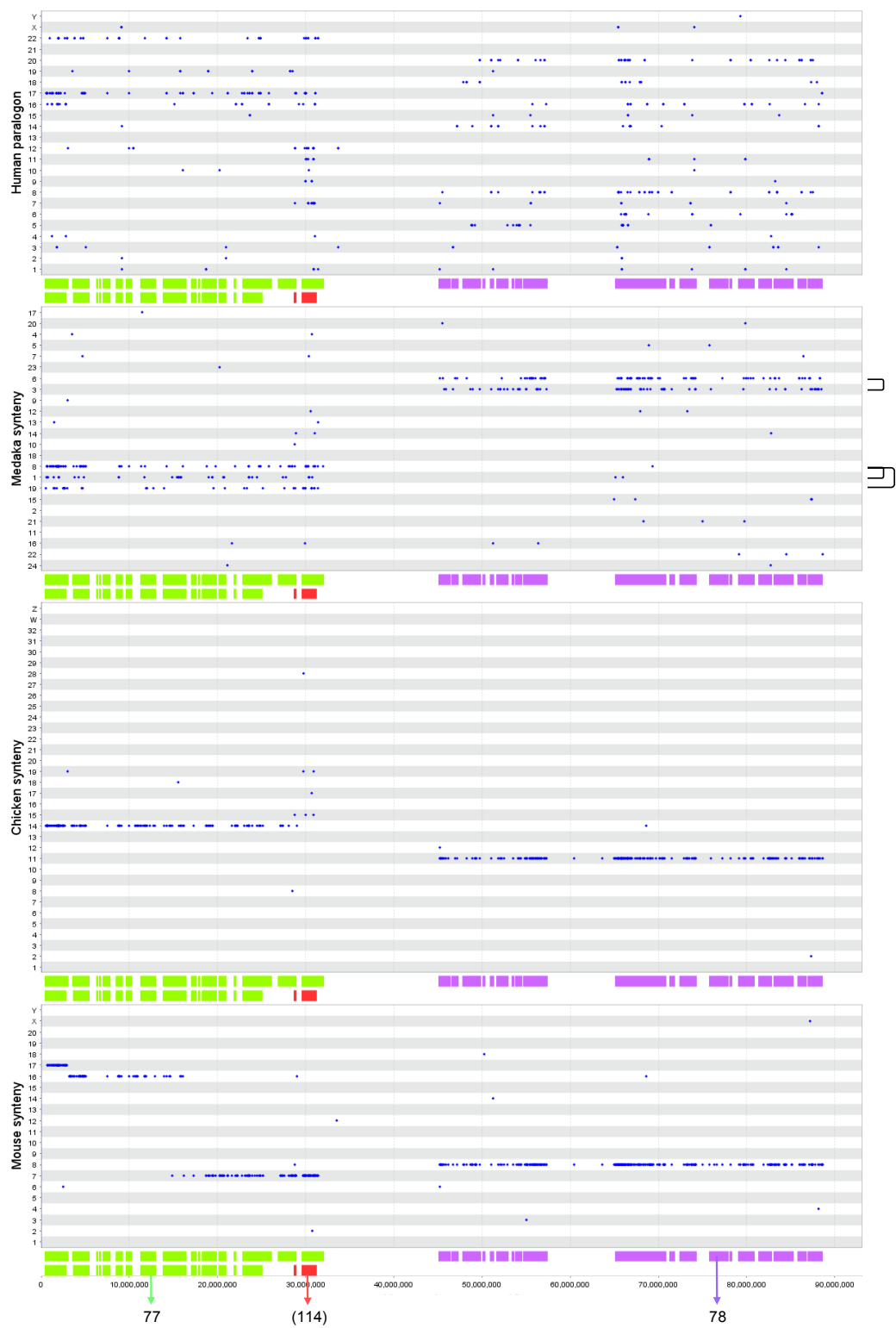
Human chromosome 13



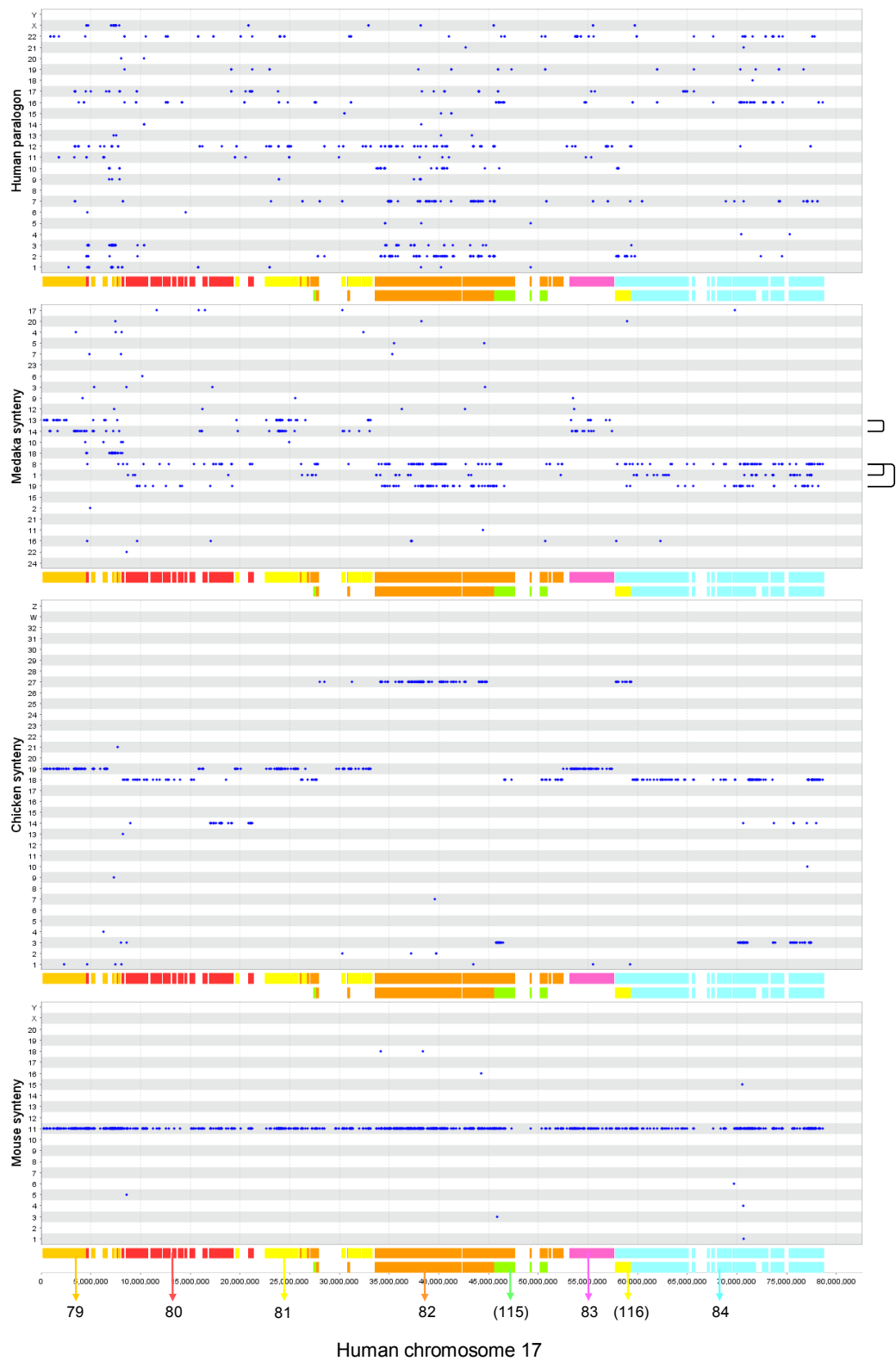
Human chromosome 14

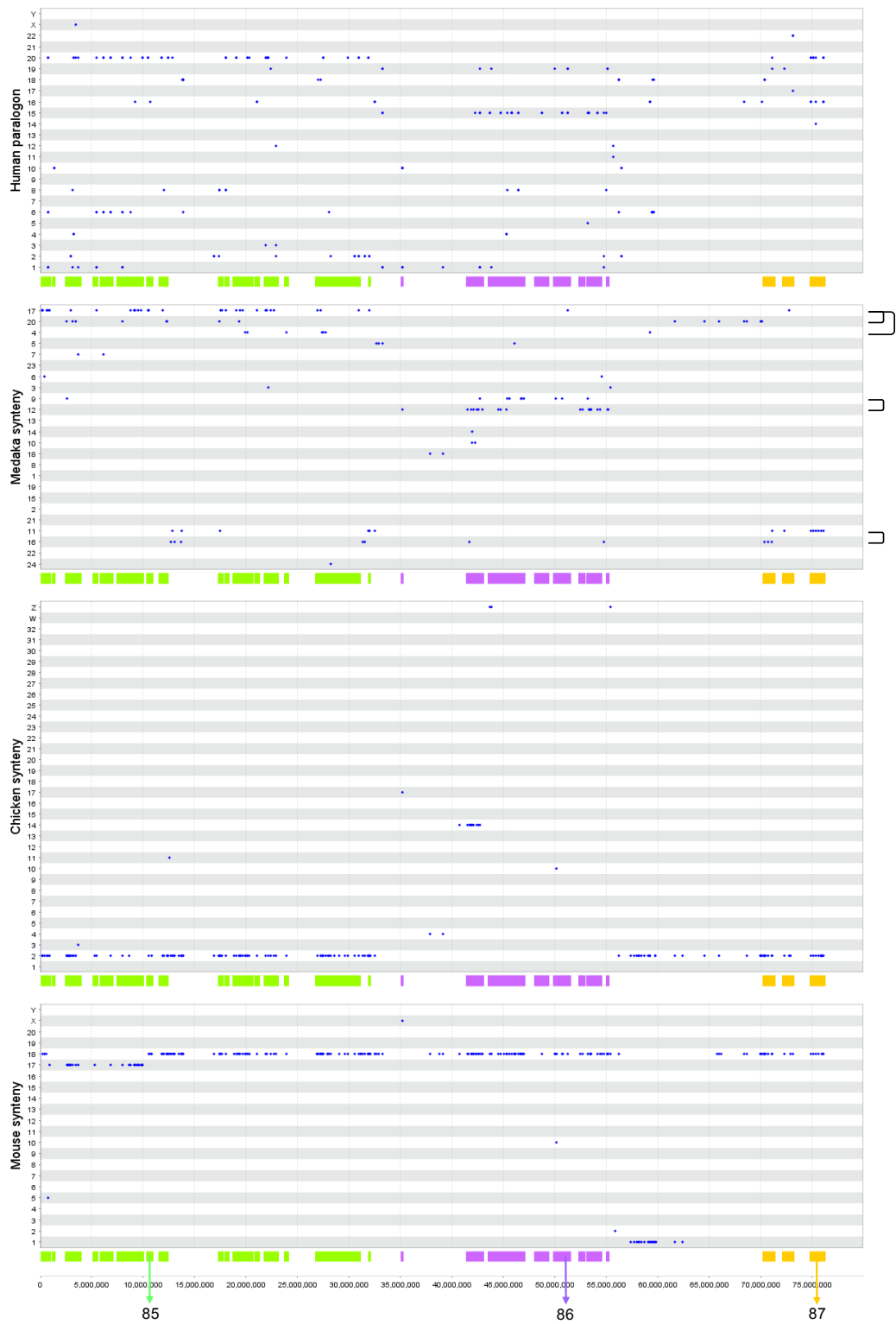


Human chromosome 15

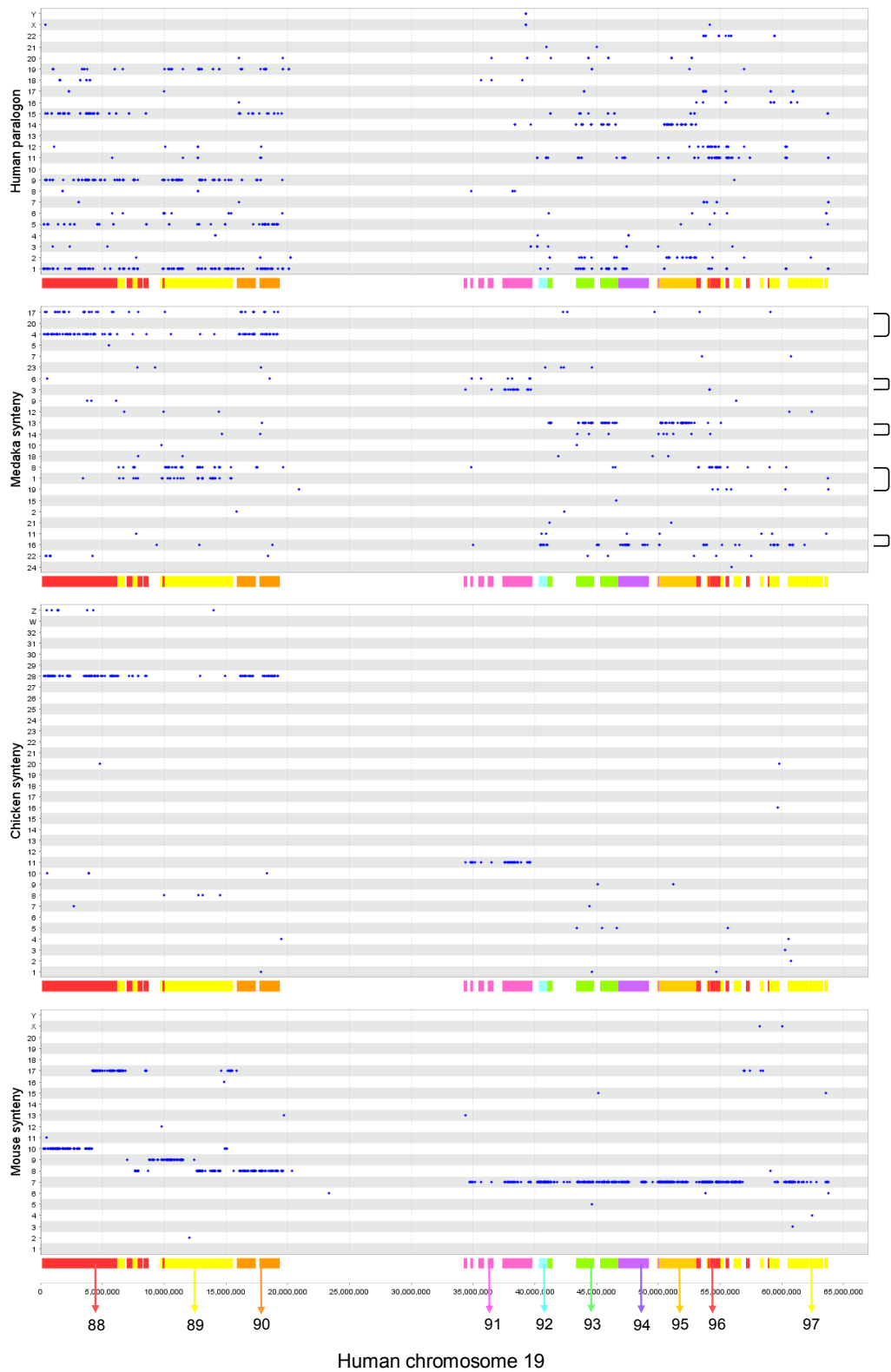


Human chromosome 16

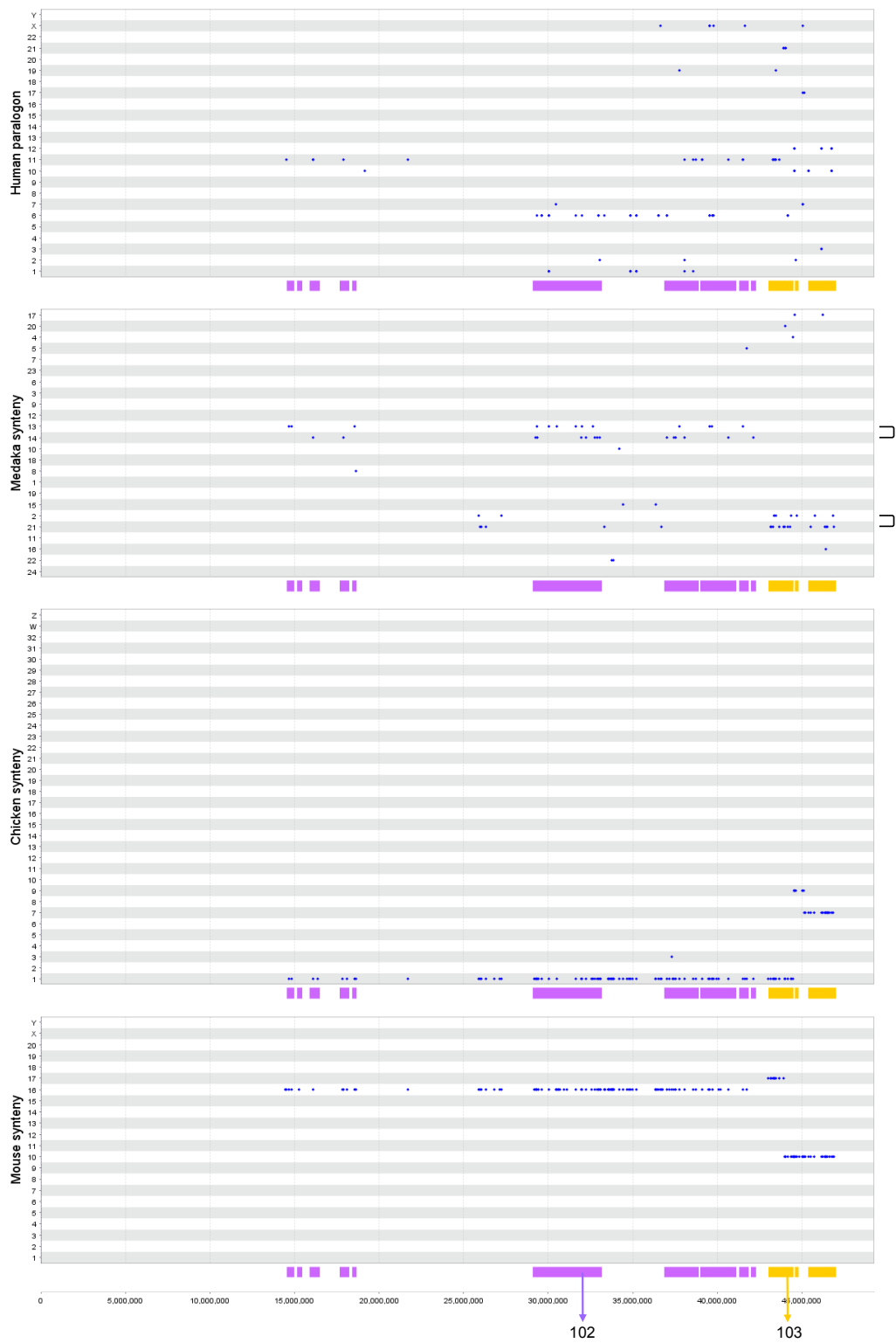




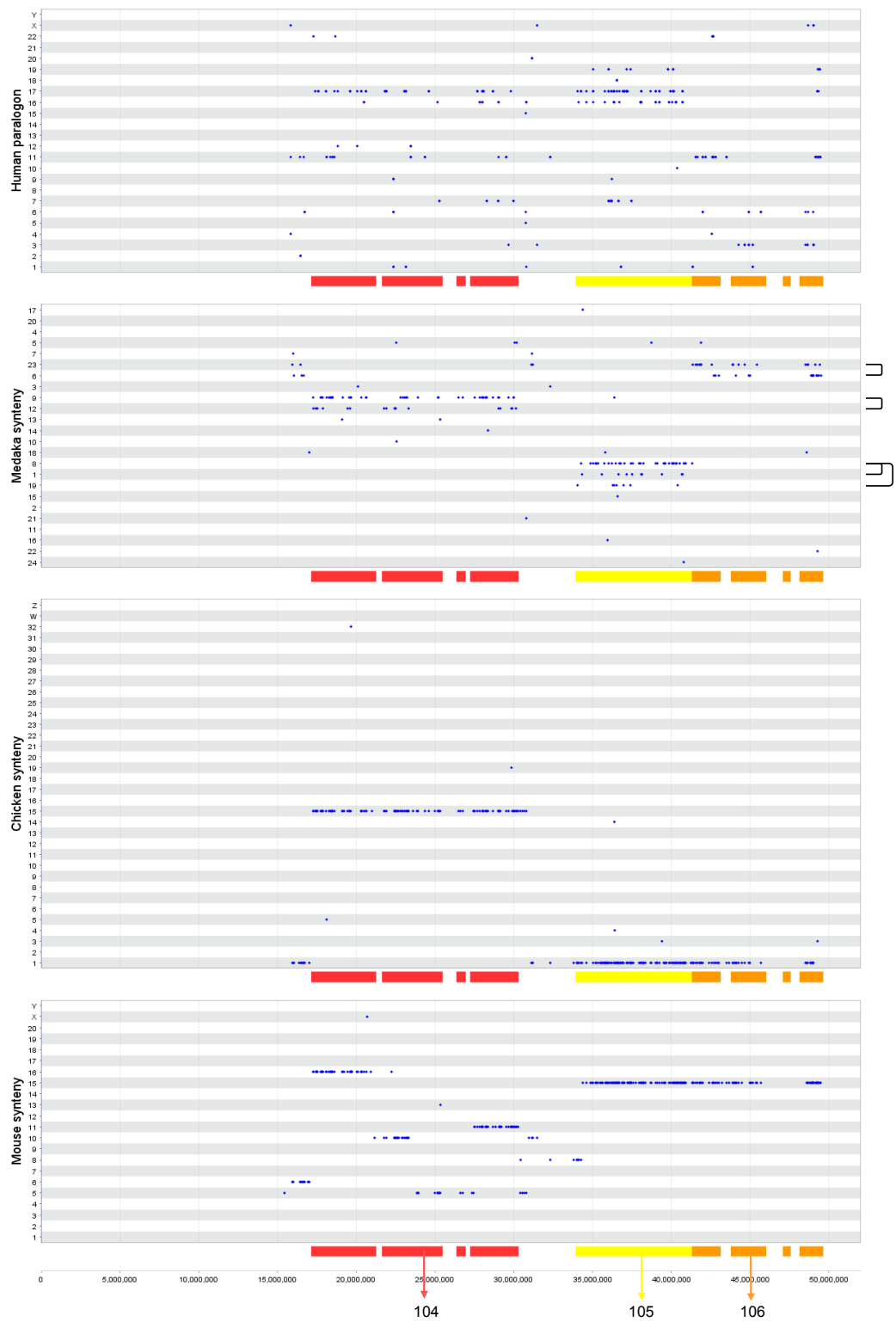
Human chromosome 18



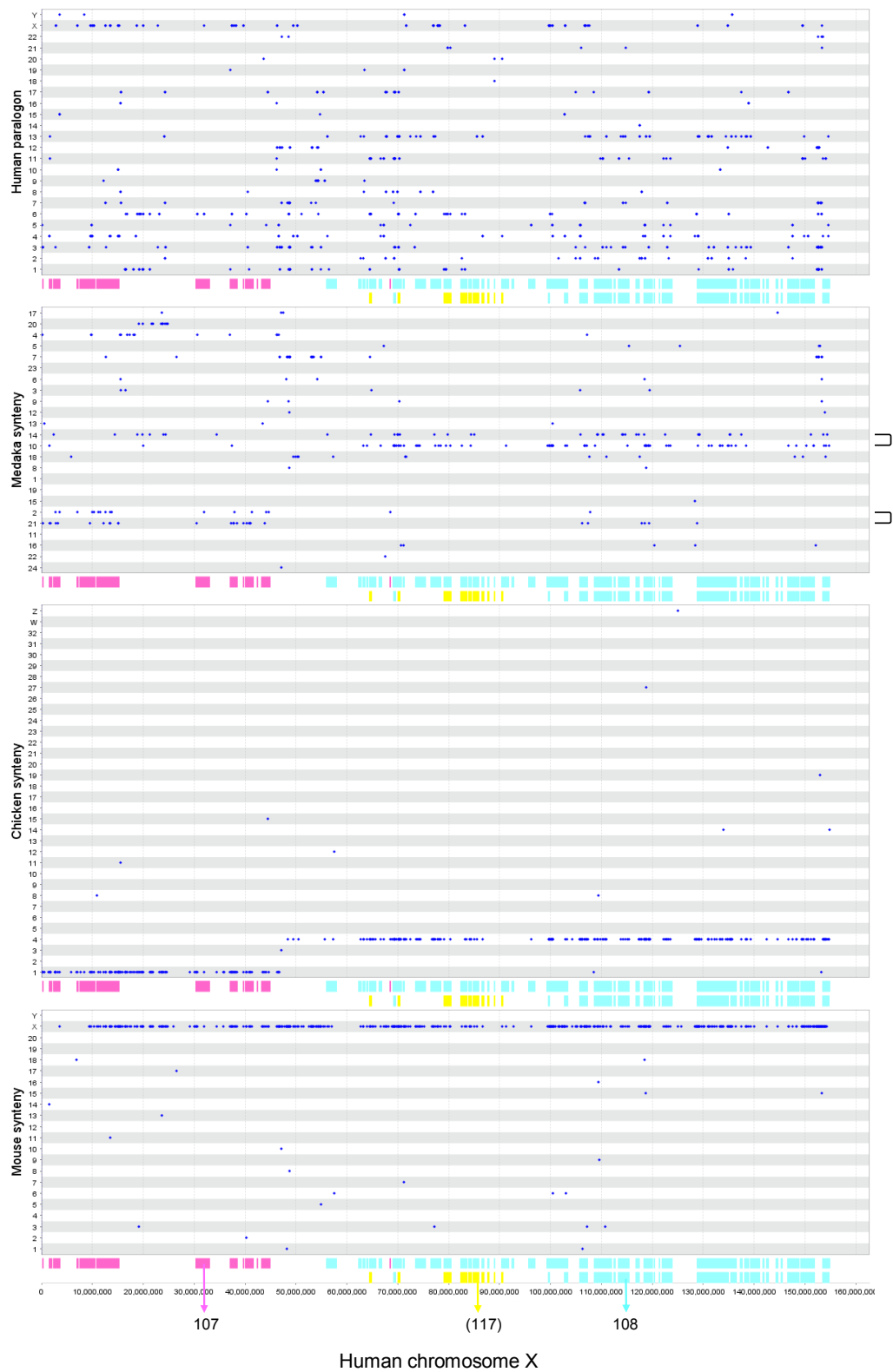


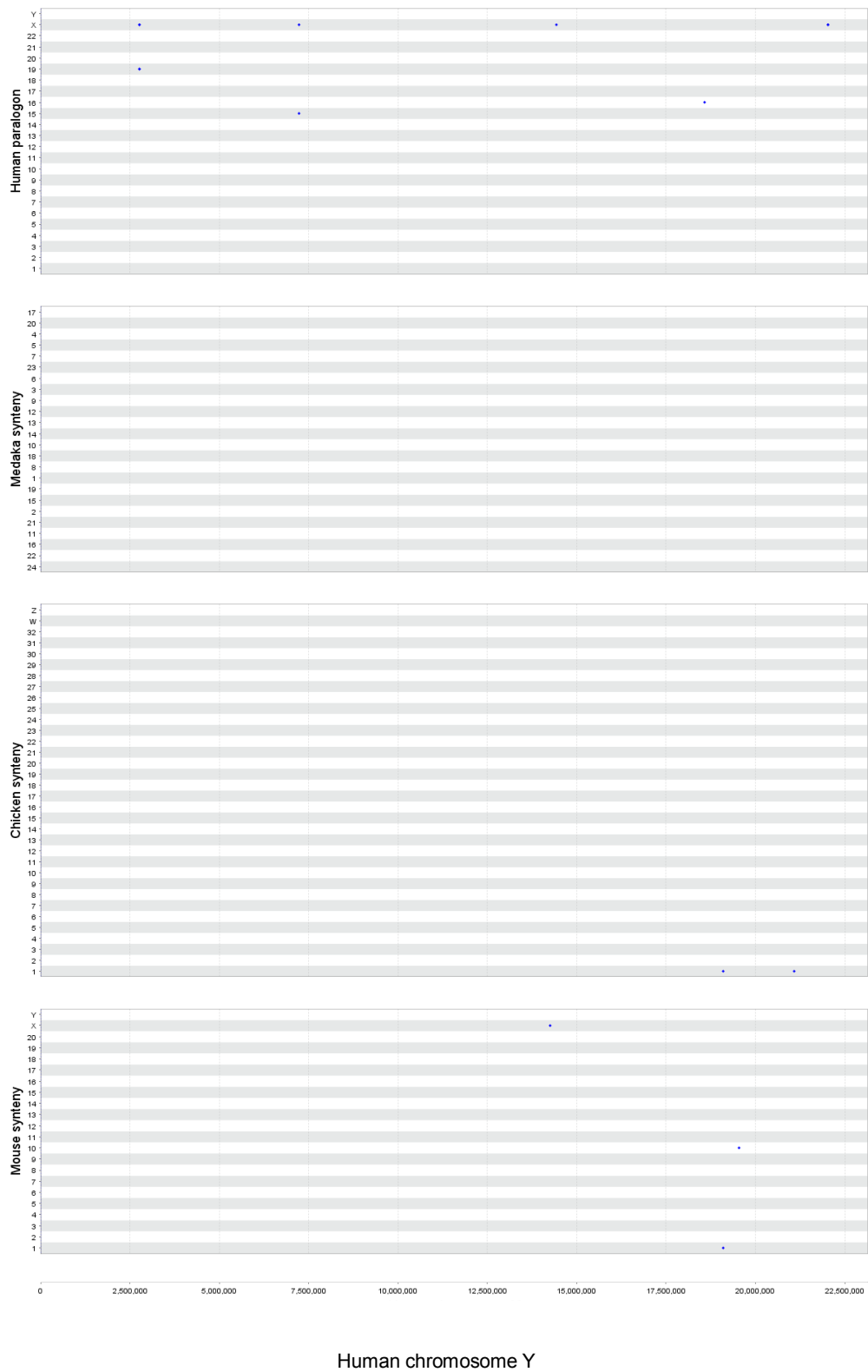


Human chromosome 21

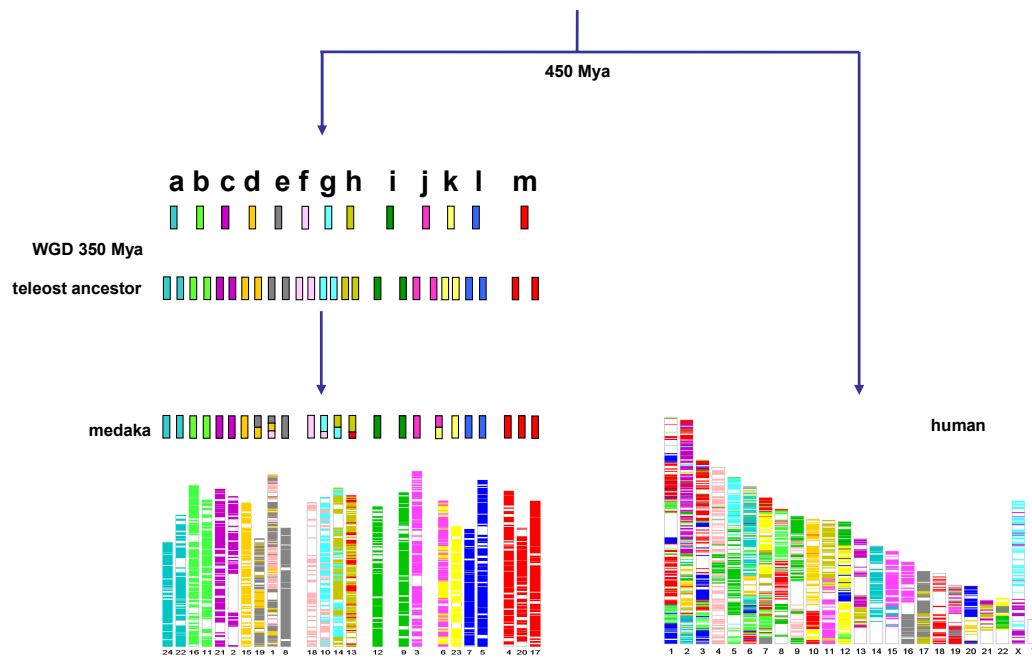


Human chromosome 22



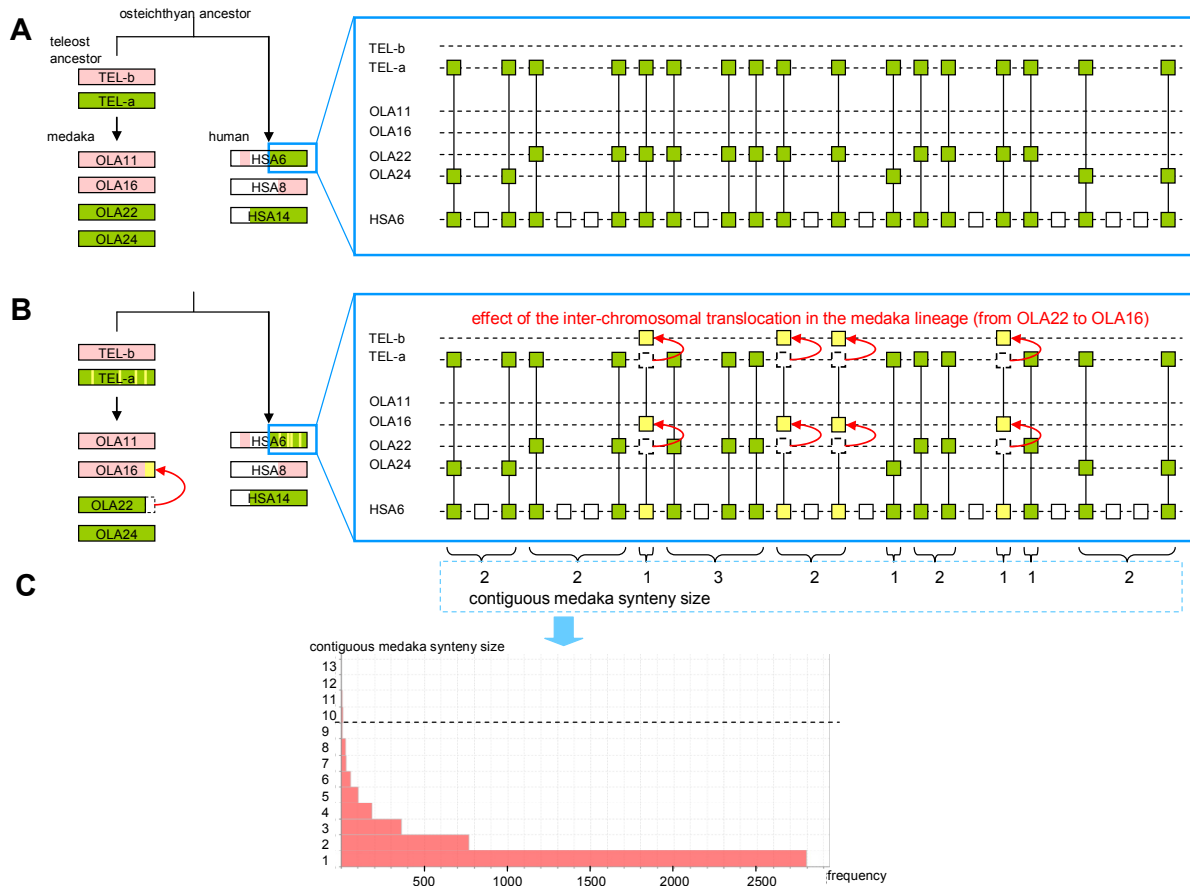


Supplementary Figure S3



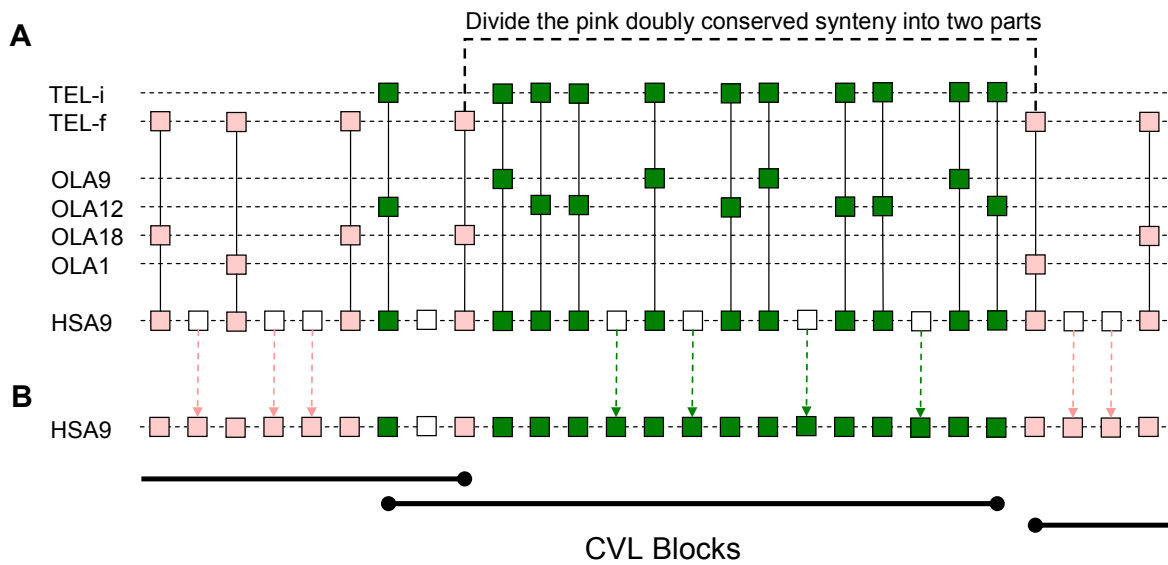
Doubly conserved synteny (DCS) regions between the human and medaka genomes. Thirteen proto-chromosomes before the teleost whole-genome duplication that yielded the ancestral teleost karyotype are colored differently for visualizing DCS regions between the medaka and human genomes.

Supplementary Figure S4



Effects of rearrangements in doubly conserved synteny (DCS) regions. A. In an ideal case of CVL block construction, a translocation occurred in the human lineage. The green chromosomal region is distributed in HSA6 and HSA14, and human genes in the green region of HSA6 have medaka orthologs in OLA22 or OLA24. These genes are properly assigned to reconstructed teleost ancestor chromosome a. **B.** Effect of translocation in the medaka lineage. The four yellow genes were originally located in the teleost ancestor chromosome TEL-a, but after whole-genome duplication (WGD) were translocated from OLA22 to OLA16 in the medaka lineage. In this case, the yellow human genes have medaka orthologs in OLA16 and are erroneously assigned to TEL-b. **C.** Distribution of contiguous medaka synteny sizes. The vertical axis shows the size of contiguous human genes that have orthologs in the same medaka chromosome and are not interrupted by orthologs to different medaka chromosomes; the horizontal axis shows the frequency of the synteny size in the human genome. Contiguous medaka synteny size is very small because of the massive gene loss after WGD that occurred randomly between duplicated chromosomes.

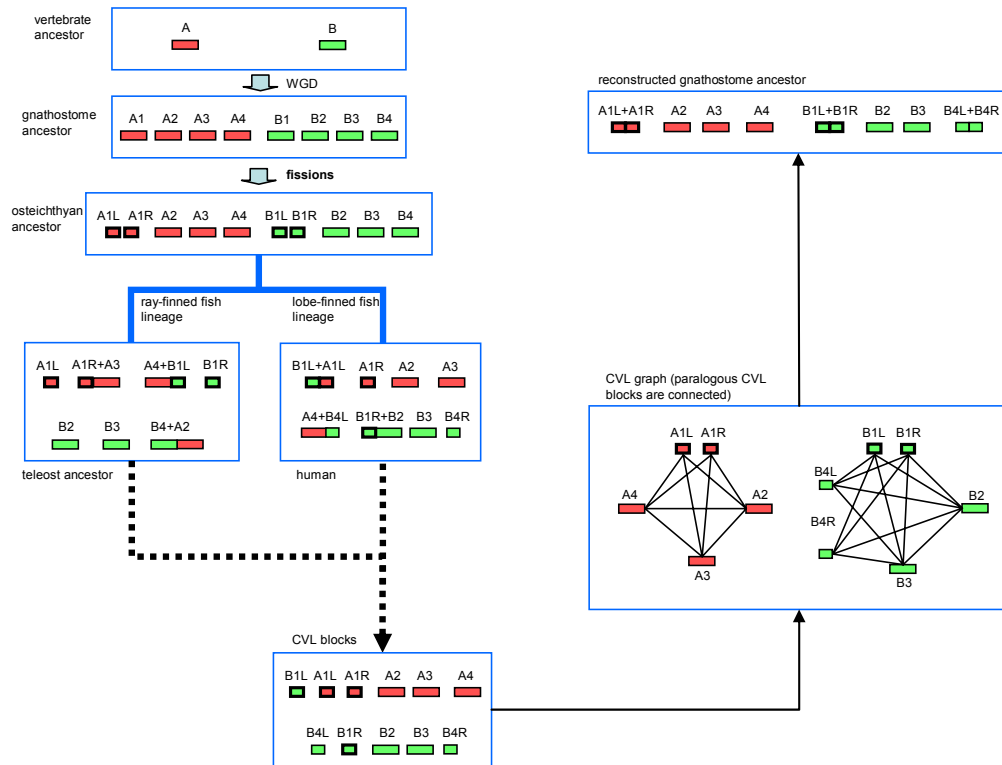
Supplementary Figure S5



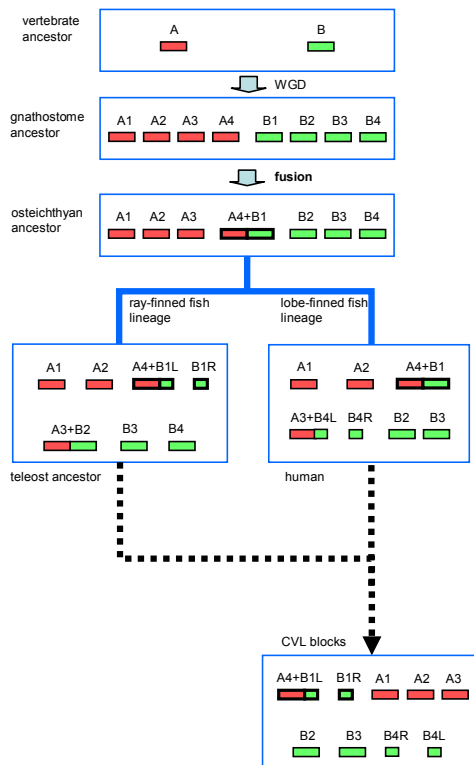
Breaking a doubly conserved synteny (DCS) region into two blocks. The Figure illustrates a subregion of human genes on human chromosome 9. **A.** Small boxes indicate genes; pink genes have orthologs in medaka chromosome OLA18 or OLA1 that originated from the proto-chromosome “f” in the teleost ancestor. Similarly, green genes are associated with the proto-chromosome “i” via OLA9 and OLA12. White genes have no orthologs and are not assigned to any medaka chromosome. Because 10 green genes divide the series of pink genes, they are treated as three distinct CVL blocks. **B.** After the generation of the three CVL blocks, some white genes were incorporated into the surrounding CVL blocks.

Supplementary Figure S6

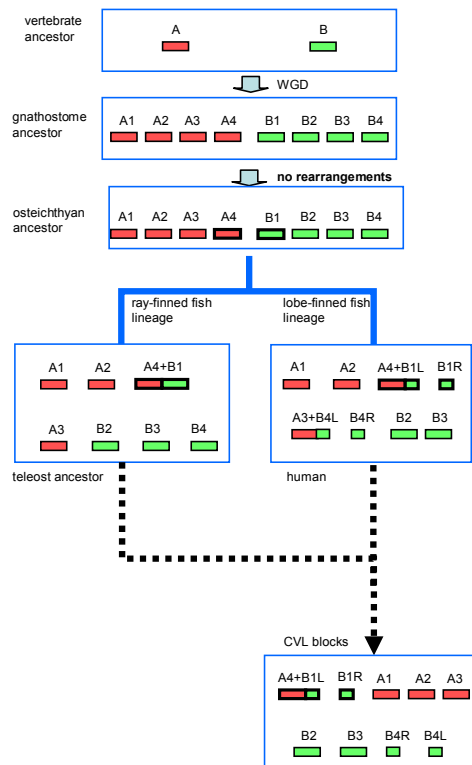
A: effect of fissions in early vertebrate evolution

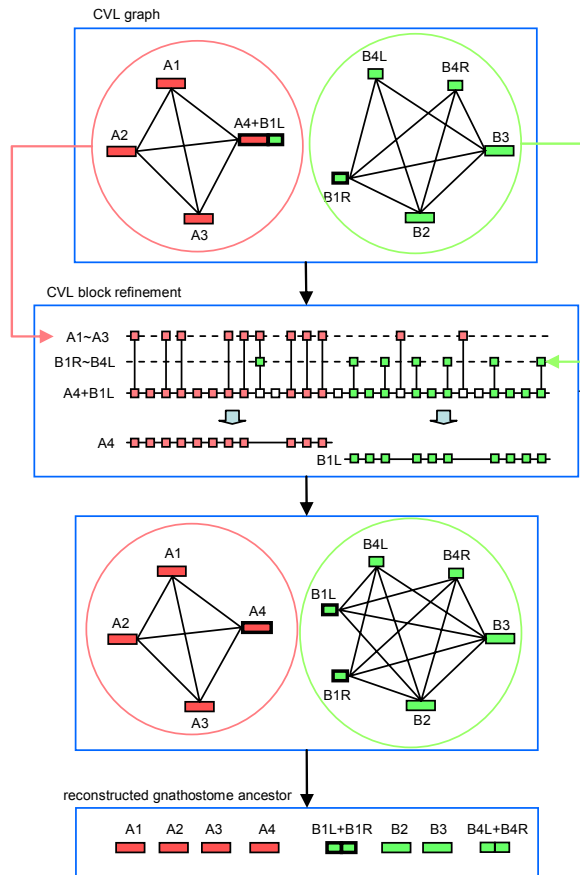
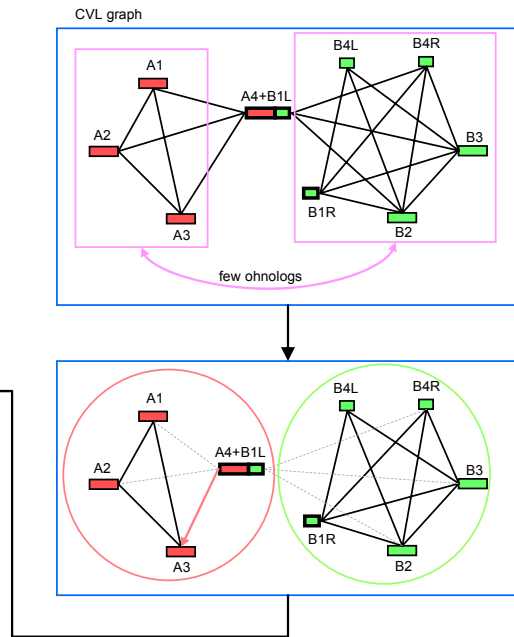


B: effect of fusions in early vertebrate evolution



C: effect of independent fusions in the teleost ancestor and human lineages



D: refinement of CVL blocks**E: refinement of connected components**

Effect of ancestral rearrangements. **A.** Chromosome fission events took place before the osteichthyan ancestor. We can reconstruct the CVL blocks by segmenting the human genome according to teleost ancestor synteny. **B.** Fusion events before the osteichthyan ancestor are likely to produce improper CVL blocks originating from multiple proto-chromosomes in the vertebrate ancestor. **C.** Similarly, erroneous CVL blocks may result from extensive rearrangements in the teleost and human lineages. **D.** Fused CVL blocks can be identified by mapping red and green ohnologs to each CVL block. **E.** Fused CVL blocks may incorrectly assign CVL blocks derived from two vertebrate proto-chromosomes to one connected component. This error can be fixed by checking the number of ohnologs between two subcomponents.

Supplementary Figure S7

Reconstruction of proto-chromosomes using CVL blocks. Ten reconstructed proto-chromosomes, A–J, are displayed. **A.** The most significant five candidates of the duplicated sister chromosomes in each proto-chromosome. The first column shows the significance of each candidate, and the remaining columns present groups of CVL blocks that constitute sister chromosomes. **B.** The table shows how genes in the vertebrate proto-chromosome are distributed in the human, teleost ancestor, and chicken genomes. The first two columns present the identifiers of CVL blocks, and the number of genes in the CVL block. The next three columns display the syntenic chromosomes in the human, teleost ancestor, and chicken genomes. The remaining columns present the numbers of orthologs shared between the CVL blocks and individual chicken chromosomes. **C.** For the most significant reconstruction of sister chromosomes, dots are plotted for orthologs in the matrix of CVL blocks grouped by sister chromosomes. **D.** The figure presents another representation of orthologs in which CVL blocks are placed on a circle, and curved lines indicate orthologs between CVL blocks. **E.** The strength of the correlation between CVL blocks in terms of orthologs is visualized using a graph in which each node represents a CVL block labeled with its identifier, its human chromosome number, its teleost proto-chromosome number, and its size. The connected lines between CVL block nodes represent the significance of the number of shared orthologs. Bold lines indicate significant relations with the probability <0.0001 , while dotted lines indicate the probability <0.001 .

A

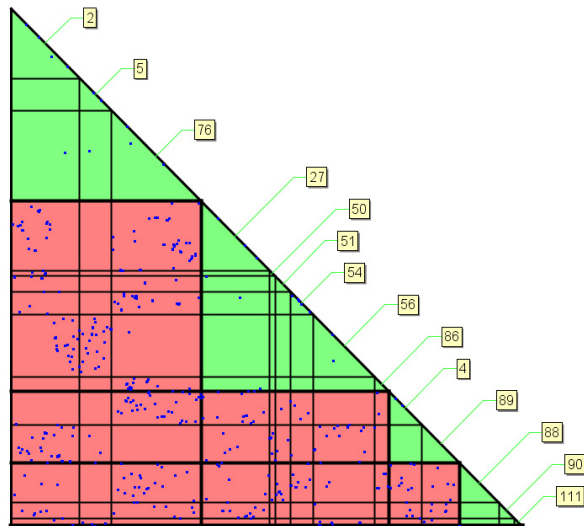
a

significance	gnathostome subgroup 1	gnathostome subgroup 2	gnathostome subgroup 3	gnathostome subgroup 4	gnathostome subgroup 5
2.72E-37	2,5,76	27,50,51,54,56,86	4,89	88,90,111	
1.42E-36	2,5,76	27,50,51,54,56,86	4,89	88,90	111
4.32E-36	2,5,76	27,50,51,54,56,86	4,89	88,111	90
8.88E-36	2,5,76	27,50,51,54,56,86	4,89	88	90,111
1.17E-35	2,5,76	27,51,54,56,86	4,89	88,90,111	50

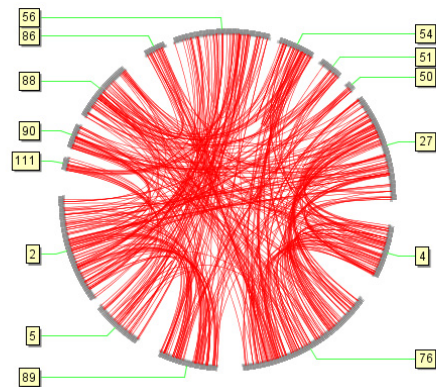
b

ID	gene	human	teleost	chicken	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	W	Z	Un
2	305	1	m	8	0	0	0	0	0	0	0	0	224	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
5	142	1	m	8,Un	0	0	0	0	0	0	0	0	92	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	4	
76	399	15	j	10	0	0	0	0	0	0	0	0	0	203	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
27	306	5	i	10,13,WZ,Un	0	0	0	0	0	0	0	0	0	4	0	0	8	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	12	146	5
50	24	9	i	Z	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	12	0	
51	67	9	f	Z	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	21	0	
54	103	9	i	Z,Un	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	53	5
56	275	9	i	17,Z	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	152	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	4	0
86	60	18	i	14,Un	0	0	0	0	0	0	0	0	0	0	0	0	0	8	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	12	0
4	148	1	b	Un	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	17	0
89	169	19	e		0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
88	175	19	m	28	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
90	68	19	m	28	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
111	35	6	b		0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	

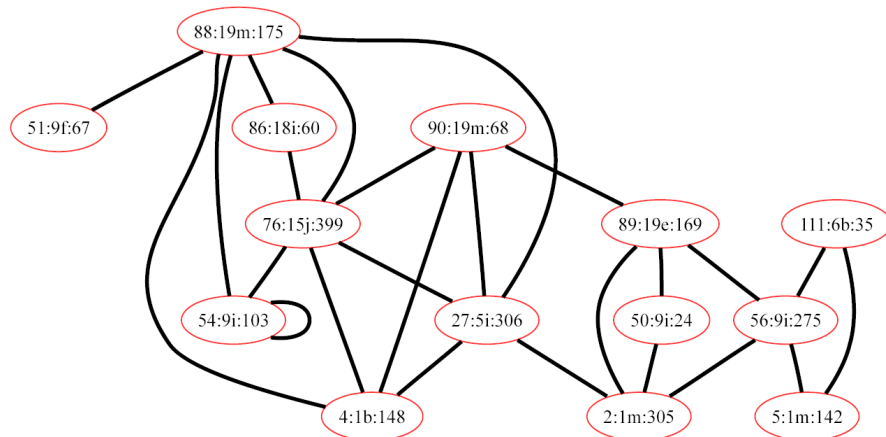
c



d



e



B

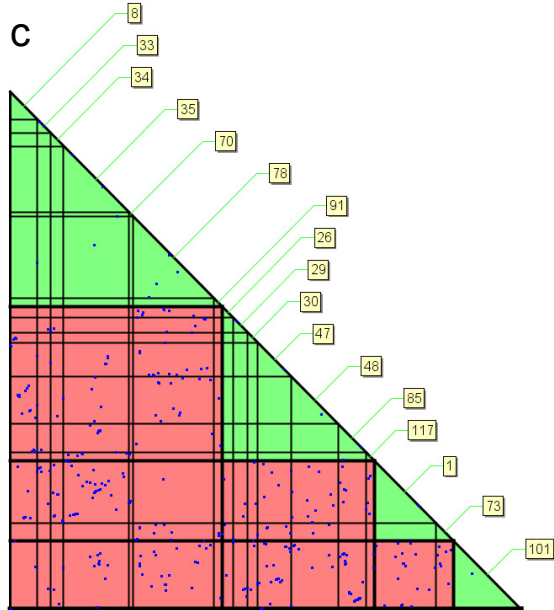
a

significance	gnathostome subgroup 1	gnathostome subgroup 2	gnathostome subgroup 3	gnathostome subgroup 4	gnathostome subgroup 5
4.62E-30	8,33,34,35,70,78,91	26,29,30,47,48,85,117	1,73	101	
4.92E-30	8,33,34,35,78,91	26,29,30,47,48,85,117	1,73	70,101	
9.02E-30	8,33,34,35,78,91	26,29,30,47,48,70,85	1,73	101	117
1.31E-29	8,33,34,35,78,91	26,29,30,47,48,85,117	1,73	101	70
2.60E-29	8,33,34,35,70,78,91	26,29,30,47,48,85	1,73	101	117

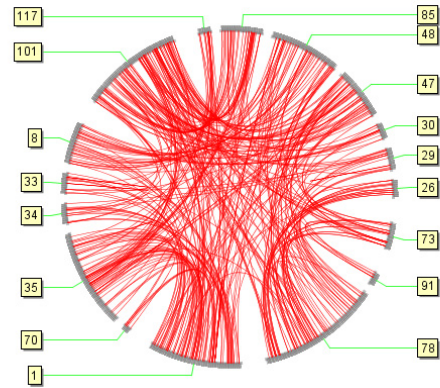
b

ID	gene	human	teleost	chicken	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	W	Z	Un
8	105	2	a	3	0	0	64	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
33	51	6	a	3	0	0	28	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
34	50	6	d	3	0	0	35	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
35	251	6	a	3,Un	0	0	176	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	6	
70	18	13	c	1	8	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
78	312	16	j	11	0	0	0	0	0	0	0	0	0	0	152	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
91	31	19	j	11	0	0	0	0	0	0	0	0	0	0	25	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
26	43	5	b	2	0	15	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
29	55	6	m	2	0	24	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
30	39	6	b	2	0	32	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
47	130	8	m	2	0	74	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
48	180	8	b	2	0	116	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
85	107	18	m	2,Un	0	54	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	4	
117	32	X	g	4	0	0	8	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
1	239	1	b	23	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	125	0	0	0	0	0	0	0	0	0	0	0	0
73	67	14	m		0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
101	261	20	l	20	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	143	0	0	0	0	0	0	0	0	0	0	0	0	0	0

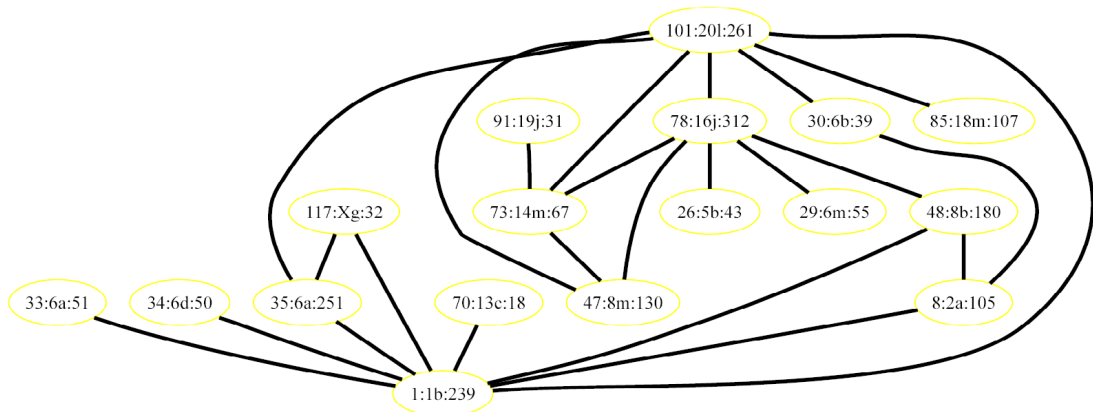
c



d



e



C

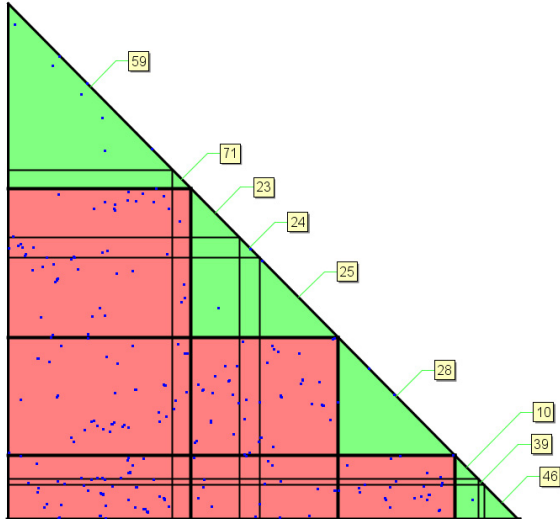
a

significance	gnathostome subgroup 1	gnathostome subgroup 2	gnathostome subgroup 3	gnathostome subgroup 4	gnathostome subgroup 5
9.35E-20	59,71	23,24,25	28	10,39,46	
1.14E-19	59,71	23,24,25	28	39,46	10
1.32E-19	59,71	23,24,25	28	46	10,39
1.97E-19	59,71	23,24,25	28	10,46	39
5.30E-19	59,71	10,23,24,25	28	39,46	

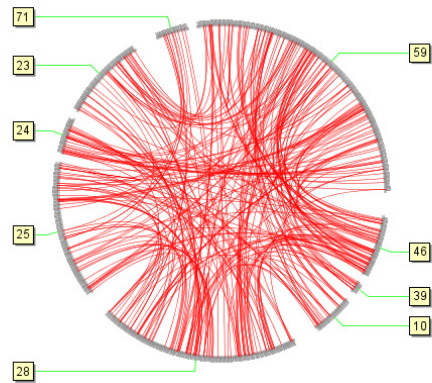
b

ID	gene	human	teleost	chicken	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	W	Z	Un
59	500	10	d	6	0	0	0	0	0	244	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
71	56	13	h	1	37	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
23	148	4	f	4	0	0	0	73	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
24	61	4	i	4	0	0	0	29	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
25	239	4	f	4	0	0	0	153	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
28	355	5	g	13	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
10	72	2	i	22,Un	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	4	0	0	0	0	0	0	0	0	0	0	0	6	
39	18	7	i		0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
46	106	8	i	22,Un	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	15	0	0	0	0	0	0	0	0	0	0	0	30	

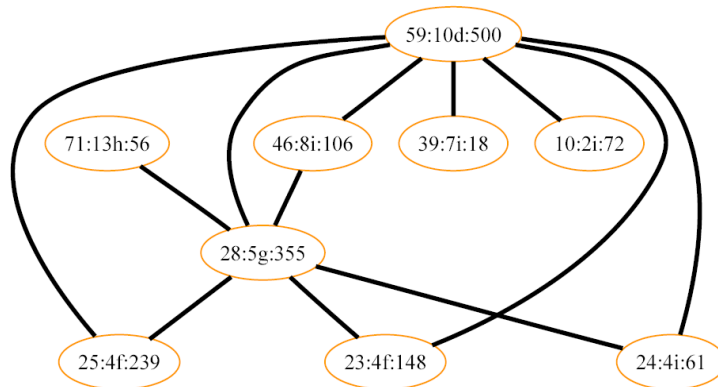
c



d



e



D

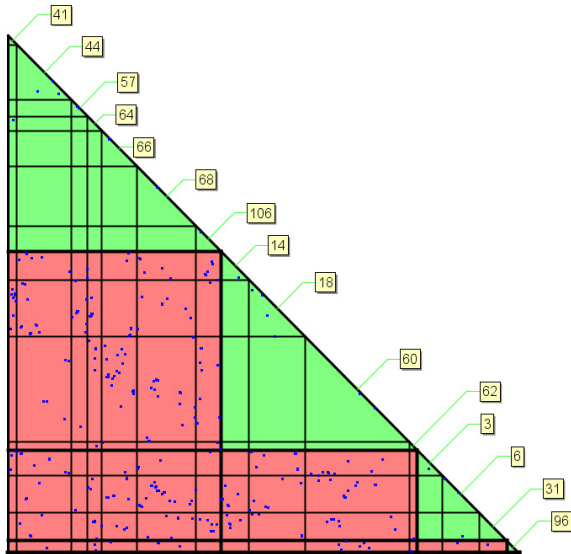
a

significance	gnathostome subgroup 1	gnathostome subgroup 2	gnathostome subgroup 3	gnathostome subgroup 4	gnathostome subgroup 5
2.36E-30	41,44,57,64,66,68,106	14,18,60,62	3,6,31	96	
4.52E-30	41,44,57,64,66,68,106	14,18,60,62	3,6,31,96		
3.61E-29	41,44,57,64,66,68,106	14,18,60,62	3,6,96	31	
4.47E-29	44,57,64,66,68,106	14,18,60,62	3,6,31	41,96	
4.53E-29	41,44,57,64,66,68,106	14,18,60,62	3,6	31	96

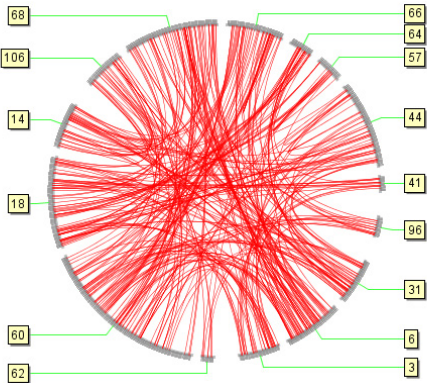
b

ID	gene	human	teleost	chicken	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	W	Z	Un
41	27	7	k	1	22	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
44	163	7	k	1,14	90	0	0	0	0	0	0	0	0	0	0	0	0	11	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
57	49	10	k	1,14	13	0	0	0	0	0	0	0	0	0	0	0	0	6	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
64	43	12	k	1	24	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
66	106	12	k	1	77	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
68	177	12	k	1	136	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
106	76	22	k	1	37	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
14	84	3	l	12	0	0	0	0	0	0	0	0	0	0	0	0	55	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
18	169	3	l	12,Un	0	0	0	0	0	0	0	0	0	0	0	0	83	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	5
60	313	11	j	3,5	0	0	20	0	155	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
62	23	11	j	5	0	0	0	0	10	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
3	77	1	l	26	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		
6	109	1	l	26,32	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	17	0	0	
31	84	6	l	26	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	30	0	0	0	0	0	0		
96	37	19	e		0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	

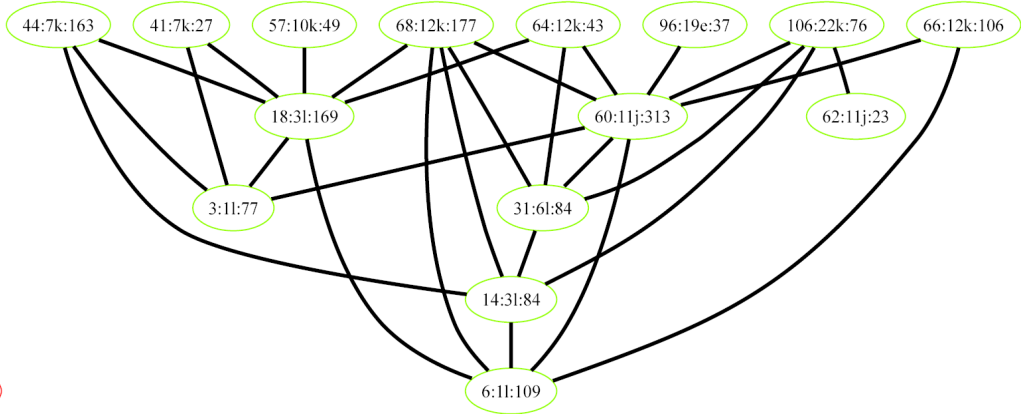
c



d



e



E

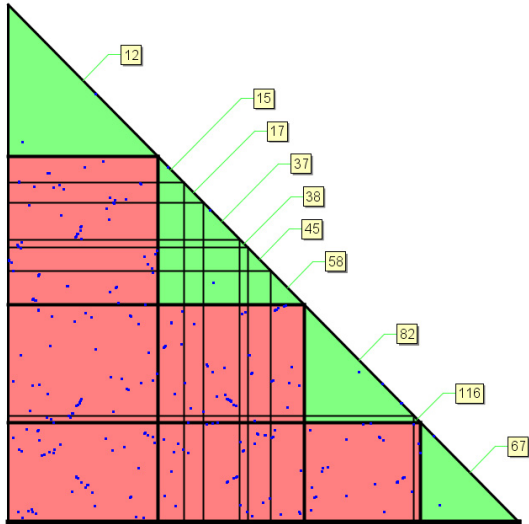
a

significance	gnathostome subgroup 1	gnathostome subgroup 2	gnathostome subgroup 3	gnathostome subgroup 4	gnathostome subgroup 5
4.00E-20	12	15,17,37,38,45,58	82,116	67	
2.21E-19	12	15,17,37,38,45,58	82	67	116
4.99E-19	12	15,17,37,45,58	82,116	67	38
9.48E-19	12	15,37,38,45,58	17,67	82,116	
9.94E-19	15,17,37,38,45,58,116	12	82	67	

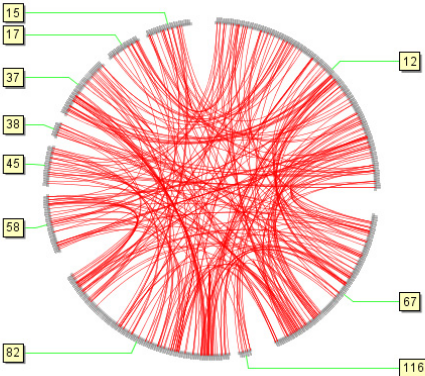
b

ID	gene	human	teleost	chicken	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	W	Z	Un
12	368	2	c	7	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
15	65	3	b	2	0	48	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
17	48	3	m	2	0	22	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
37	90	7	b	2	0	63	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
38	20	7	m	2	0	12	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
45	56	7	m	2	0	29	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
58	82	10	m	2	0	56	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
82	271	17	e	27,Un	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	10
116	17	17	e	27	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
67	242	12	i	1,Un	6	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	6

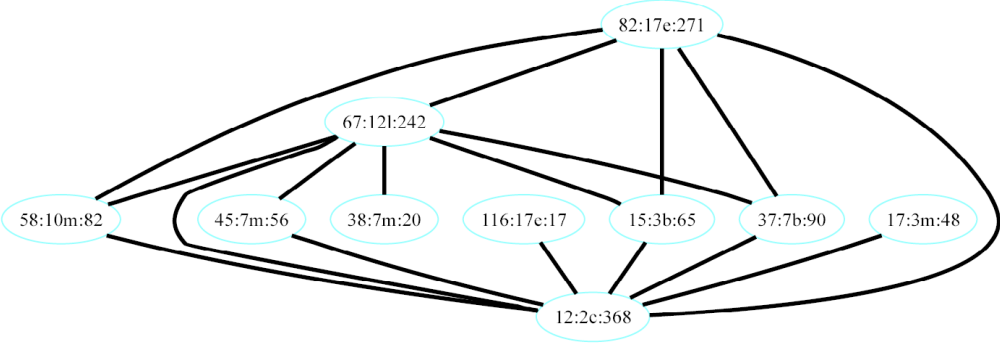
c



d



e



F

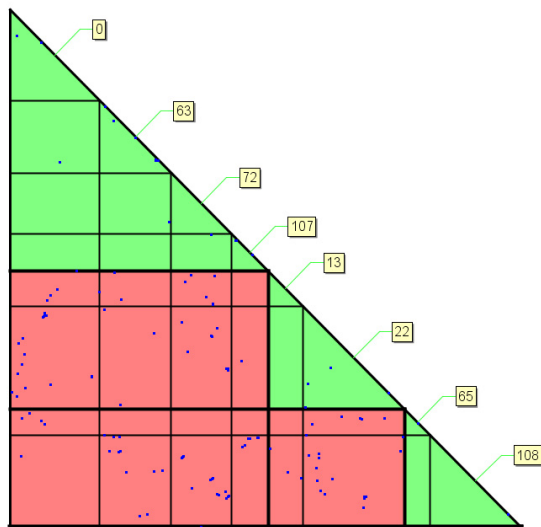
a

significance	gnathostome subgroup 1	gnathostome subgroup 2	gnathostome subgroup 3	gnathostome subgroup 4	gnathostome subgroup 5
1.02E-06	0.63,72,107	13,22	65,108		
1.83E-06	0.63,65,72,107	13,22	108		
6.43E-06	0.63,72,107	13,22	108	65	
1.19E-05	0.63,72	13,22,107	65,108		
1.30E-05	0.63,72,107	13,65,108	22		

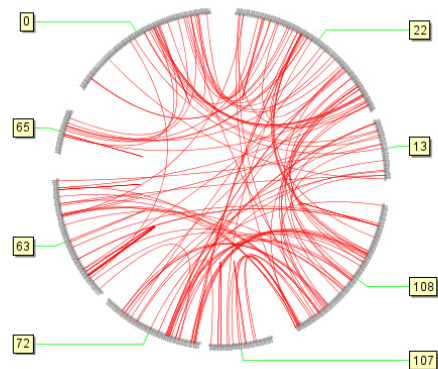
b

ID	gene	human	teleost	chicken	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	W	Z	Un
0	215	1	l	21_Un	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	100	0	0	0	0	0	0	0	0	0	0	0	0	0	4
63	173	11	h	1	92	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
72	144	13	c	1	86	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
107	89	X	c	1	56	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
13	83	2	m	9	0	0	0	0	0	0	0	0	37	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
22	246	3	m	2,9	0	6	0	0	0	0	0	0	144	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
65	61	12	b	1	26	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
108	219	X	g	4	0	0	0	64	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	

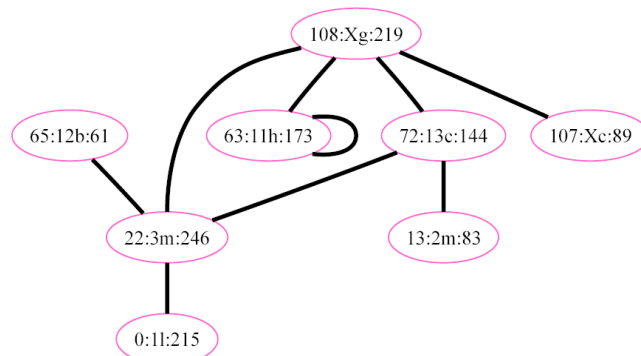
c



d



e



G

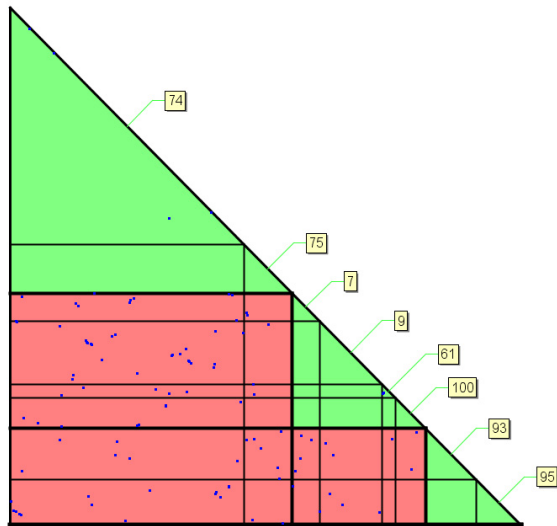
a

significance	gnathostome subgroup 1	gnathostome subgroup 2	gnathostome subgroup 3	gnathostome subgroup 4	gnathostome subgroup 5
1.12E-15	74,75	7,9,61,100	93,95		
7.01E-15	74,75	7,9,100	93,95	61	
1.40E-14	74,75	7,9,61,100	93	95	
1.54E-14	74,75	7,9,100	61,93,95		
2.79E-14	74,75	9,61,100	93,95	7	

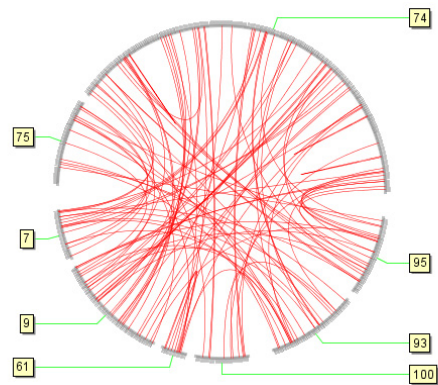
b

ID	gene	human	teleost	chicken	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	W	Z	Un
74	408	14	a	5	0	0	0	0	275	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
75	84	15	a	5	0	0	0	0	57	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
7	48	1	d	3	0	0	16	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
9	109	2	d	3	0	0	53	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
61	24	11	h		0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
100	52	20	d	3,14	0	0	8	0	0	0	0	0	0	0	0	0	0	13	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
93	89	19	h		0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
95	78	19	h		0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

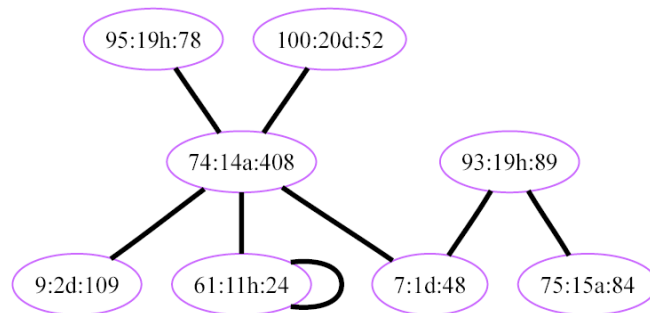
c



d



e



H

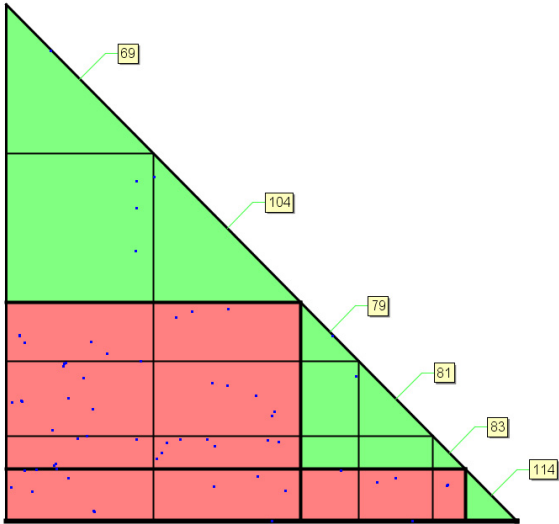
a

significance	gnathostome subgroup 1	gnathostome subgroup 2	gnathostome subgroup 3	gnathostome subgroup 4	gnathostome subgroup 5
1.59E-08	69,104	79,81,83	114		
1.19E-07	69,104	79,114	81,83		
1.26E-07	69,104	81,83	79	114	
1.74E-07	69,104	79,83	81	114	
2.29E-07	69,104	79,81	114	83	

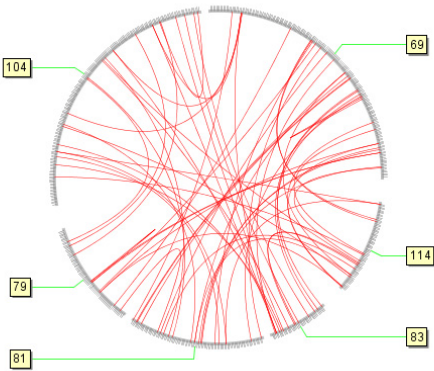
b

ID	gene	human	teleost	chicken	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	W	Z	Un	
69	186	12	i	15	0	0	0	0	0	0	0	0	0	0	0	0	0	0	110	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
104	186	22	i	15	0	0	0	0	0	0	0	0	0	0	0	0	0	0	68	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
79	73	17	h	19	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	41	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
81	94	17	h	19	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	45	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
83	41	17	h	19	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	27	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
114	65	16	e		0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	

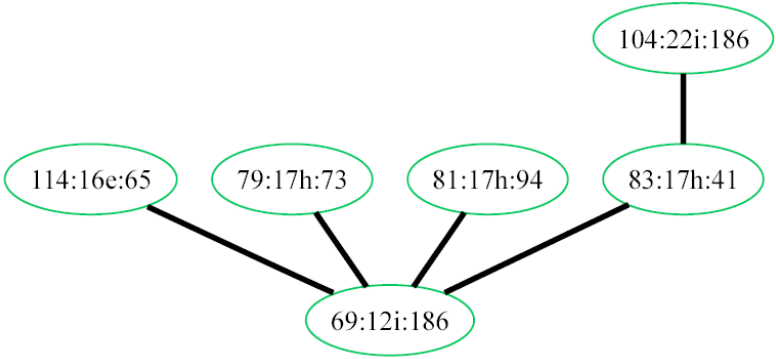
c



d



e



I

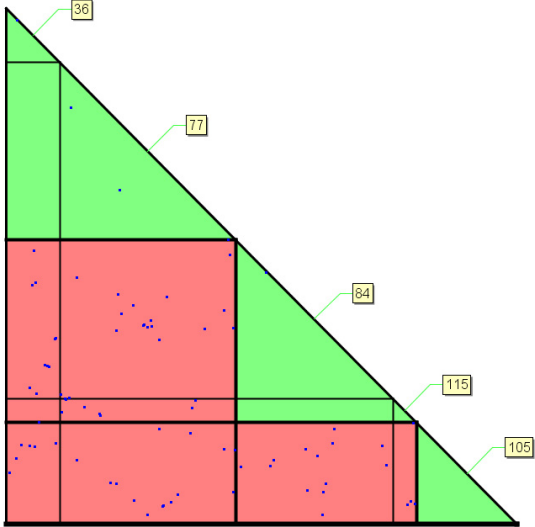
a

significance	gnathostome subgroup 1	gnathostome subgroup 2	gnathostome subgroup 3	gnathostome subgroup 4	gnathostome subgroup 5
4.77E-11	36,77	84,115	105		
9.50E-10	36,77	84	105	115	
6.81E-08	84,115	77	105	36	
9.01E-07	77	84	105	36	115
2.23E-06	77	84	105	36,115	

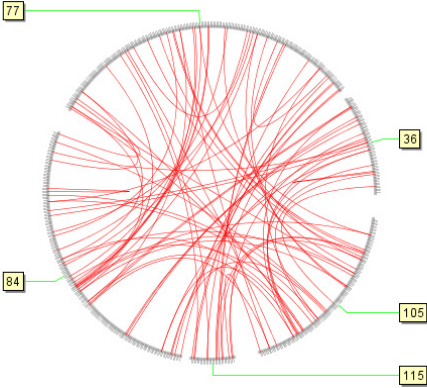
b

ID	gene	human	teleost	chicken	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	W	Z	Un
36	66	7	e	14	0	0	0	0	0	0	0	0	0	0	0	0	0	38	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
77	215	16	e	14,Un	0	0	0	0	0	0	0	0	0	0	0	0	0	86	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	4
84	193	17	e	3,18	0	0	39	0	0	0	0	0	0	0	0	0	0	0	0	0	68	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
115	29	17	e	3,18	0	0	13	0	0	0	0	0	0	0	0	0	0	0	0	8	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
105	123	22	e	1	80	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

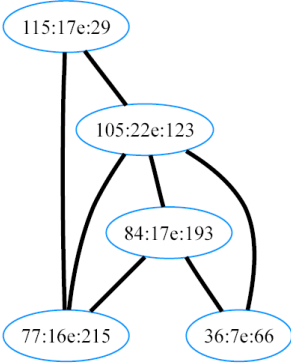
c



d



e



J

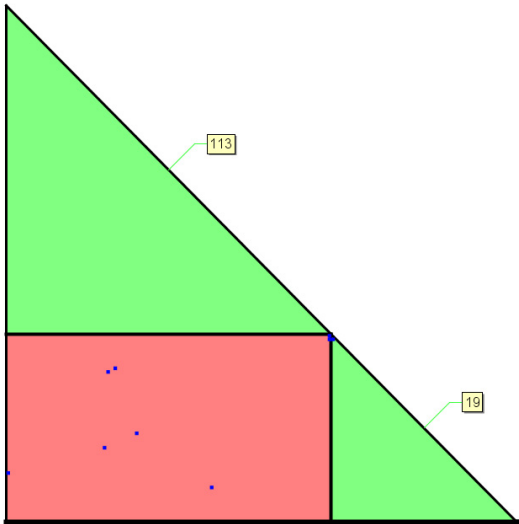
a

significance	gnathostome subgroup 1	gnathostome subgroup 2	gnathostome subgroup 3	gnathostome subgroup 4	gnathostome subgroup 5
0.011706765	113	19			

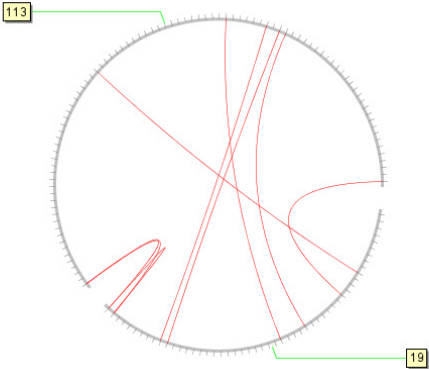
b

ID	gene	human	teleost	chicken	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	W	Z	Un
113	91	11	h	24	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	32	0	0	0	0	0	0	0	0	0	0	0	0
19	52	3	m	1	26	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

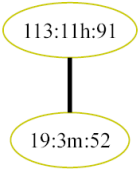
c



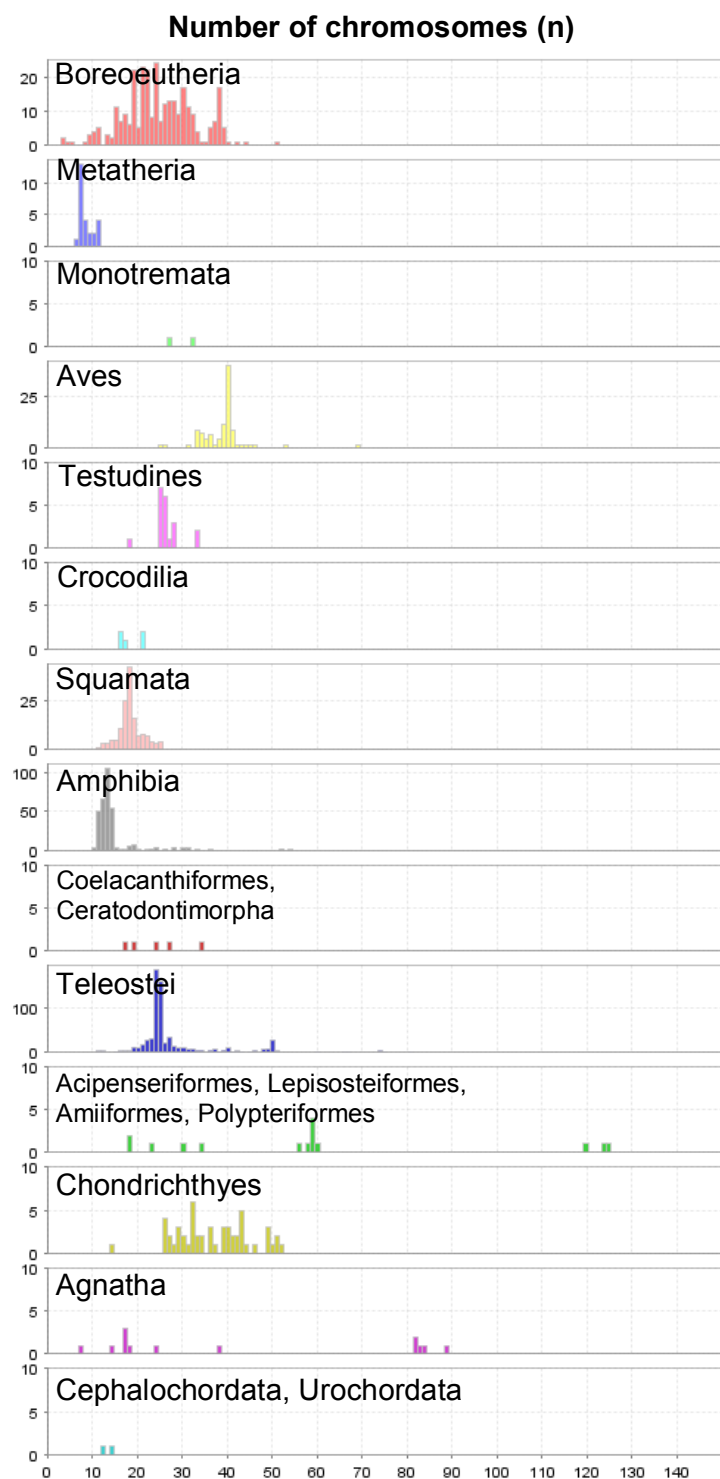
d



e

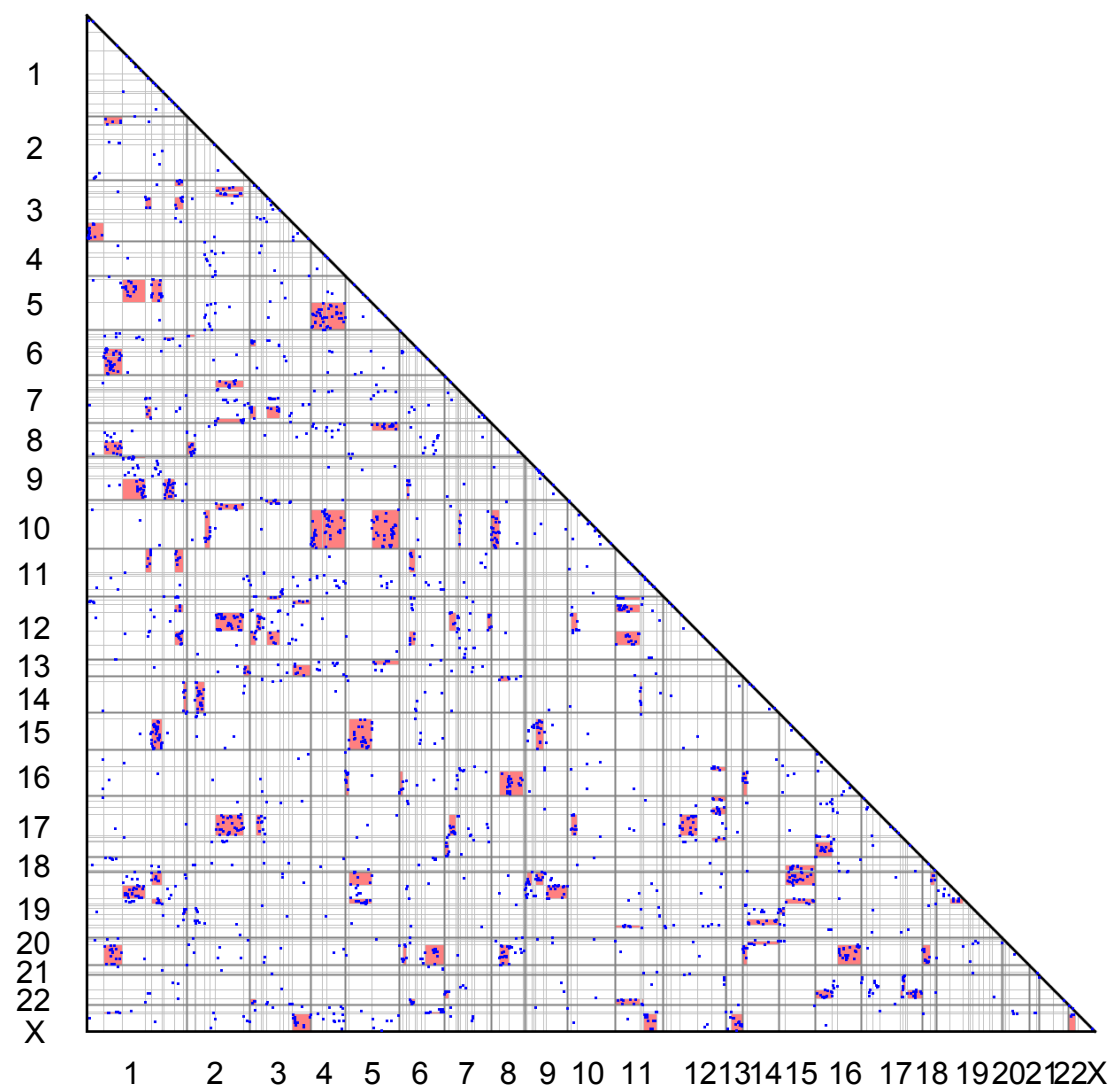


Supplementary Figure S8



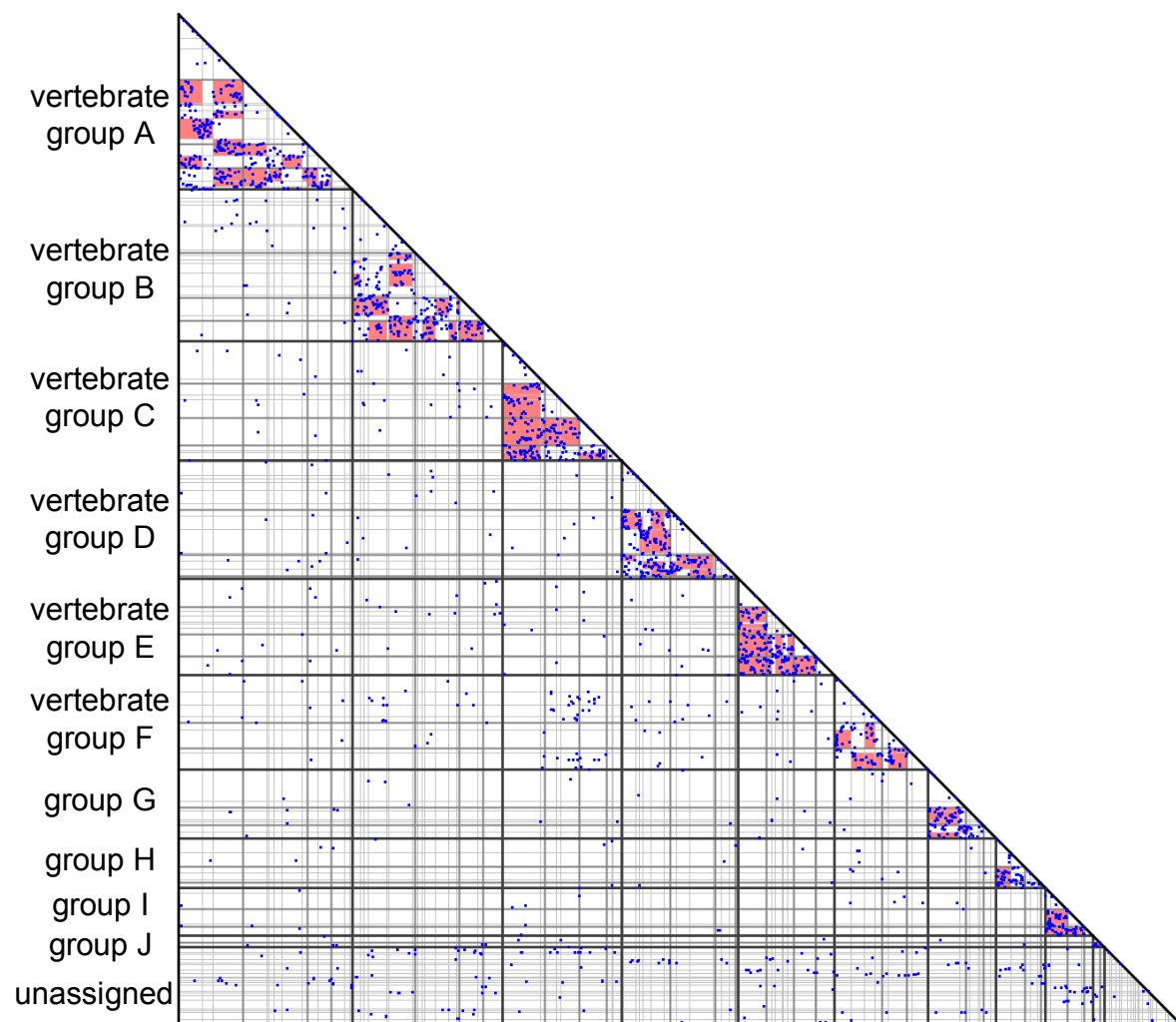
Distribution of the number of chromosomes in individual lineages. Some species possess additional lineage-specific whole-genome duplications.

Supplementary Figure S9



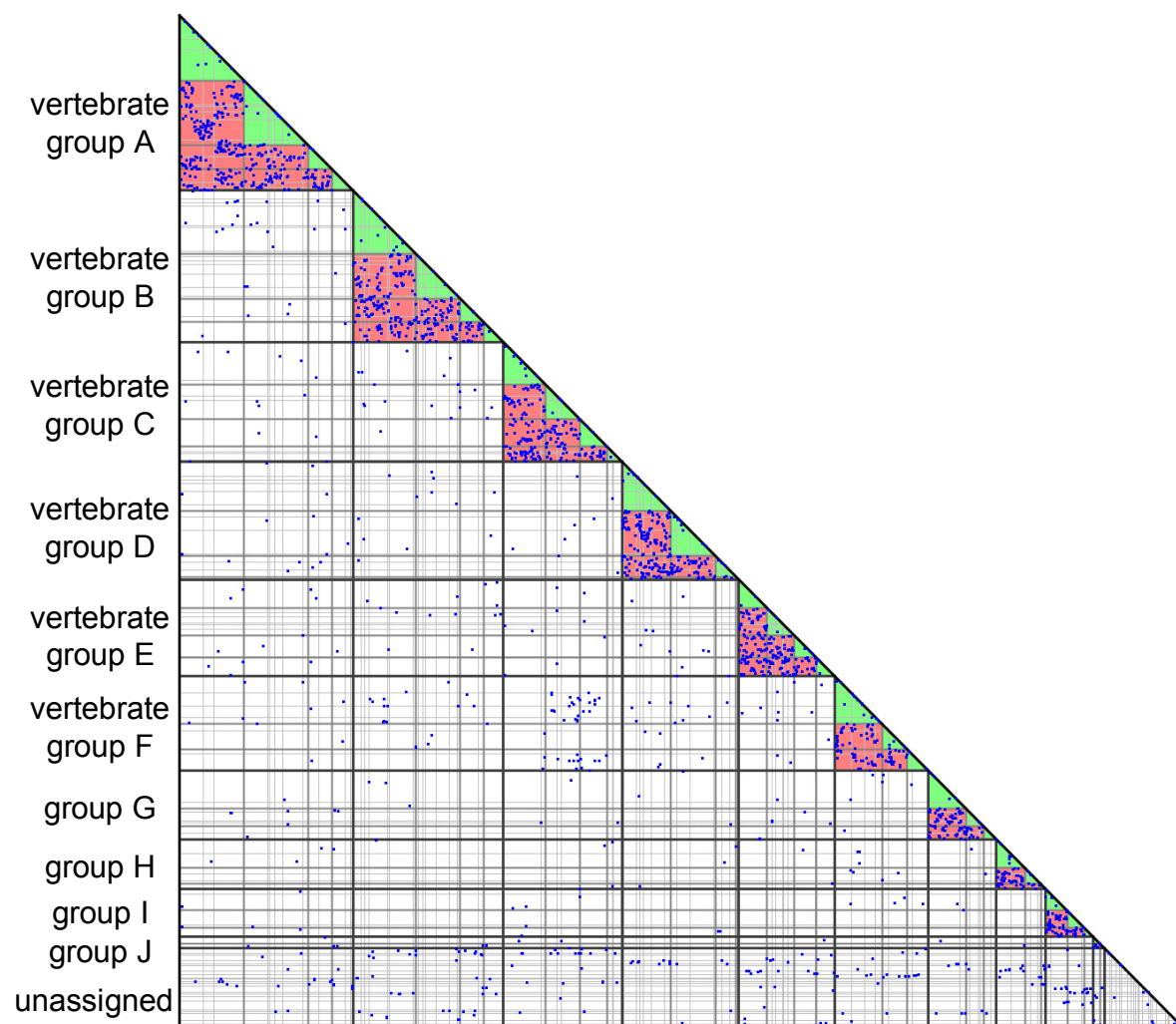
Enlarged copy of Fig. 2E.

Supplementary Figure S10



Enlarged copy of Fig. 2F.

Supplementary Figure S11



Enlarged copy of Fig. 2G.

10 Supplementary tables

Supplementary Table S1

CVL blocks in proto-chromosomes of the vertebrate ancestor. Each CVL block in the proto-chromosome is associated with its identifier, the number of human genes, the number of human genes, the number of ohnologs, human chromosome, teleost ancestor chromosome, gnathostome ancestor chromosome, syntenic chicken chromosomes, and the total syntenic block size for each chicken chromosome.

vertebrate ancestor	gnathostome ancestor	block ID	gene	ohnolog	human	teleost ancestor	chicken	Gga1	Gga2	Gga3	Gga4	Gga5	Gga6	Gga7	Gga8	Gga9	Gga10	Gga11	Gga12	Gga13	Gga14	Gga15	Gga16	Gga17	Gga18	Gga19	Gga20	Gga21	Gga22	Gga23	Gga24	Gga25	Gga26	Gga27	Gga28	Gga29	Gga30	Gga31	Gga32	Gga33	Gga34	Gga35	Gga36	Gga37	Gga38	Gga39	Gga40	Gga41	Gga42	Gga43	Gga44	Gga45	Gga46	Gga47	Gga48	Gga49	Gga50	Gga51	Gga52	Gga53	Gga54	Gga55	Gga56	Gga57	Gga58	Gga59	Gga60	Gga61	Gga62	Gga63	Gga64	Gga65	Gga66	Gga67	Gga68	Gga69	Gga70	Gga71	Gga72	Gga73	Gga74	Gga75	Gga76	Gga77	Gga78	Gga79	Gga80	Gga81	Gga82	Gga83	Gga84	Gga85	Gga86	Gga87	Gga88	Gga89	Gga90	Gga91	Gga92	Gga93	Gga94	Gga95	Gga96	Gga97	Gga98	Gga99	Gga100	Gga101	Gga102	Gga103	Gga104	Gga105	Gga106	Gga107	Gga108	Gga109	Gga110	Gga111	Gga112	Gga113	Gga114	Gga115	Gga116	Gga117	Gga118	Gga119	Gga120	Gga121	Gga122	Gga123	Gga124	Gga125	Gga126	Gga127	Gga128	Gga129	Gga130	Gga131	Gga132	Gga133	Gga134	Gga135	Gga136	Gga137	Gga138	Gga139	Gga140	Gga141	Gga142	Gga143	Gga144	Gga145	Gga146	Gga147	Gga148	Gga149	Gga150	Gga151	Gga152	Gga153	Gga154	Gga155	Gga156	Gga157	Gga158	Gga159	Gga160	Gga161	Gga162	Gga163	Gga164	Gga165	Gga166	Gga167	Gga168	Gga169	Gga170	Gga171	Gga172	Gga173	Gga174	Gga175	Gga176	Gga177	Gga178	Gga179	Gga180	Gga181	Gga182	Gga183	Gga184	Gga185	Gga186	Gga187	Gga188	Gga189	Gga190	Gga191	Gga192	Gga193	Gga194	Gga195	Gga196	Gga197	Gga198	Gga199	Gga200	Gga201	Gga202	Gga203	Gga204	Gga205	Gga206	Gga207	Gga208	Gga209	Gga210	Gga211	Gga212	Gga213	Gga214	Gga215	Gga216	Gga217	Gga218	Gga219	Gga220	Gga221	Gga222	Gga223	Gga224	Gga225	Gga226	Gga227	Gga228	Gga229	Gga230	Gga231	Gga232	Gga233	Gga234	Gga235	Gga236	Gga237	Gga238	Gga239	Gga240	Gga241	Gga242	Gga243	Gga244	Gga245	Gga246	Gga247	Gga248	Gga249	Gga250	Gga251	Gga252	Gga253	Gga254	Gga255	Gga256	Gga257	Gga258	Gga259	Gga260	Gga261	Gga262	Gga263	Gga264	Gga265	Gga266	Gga267	Gga268	Gga269	Gga270	Gga271	Gga272	Gga273	Gga274	Gga275	Gga276	Gga277	Gga278	Gga279	Gga280	Gga281	Gga282	Gga283	Gga284	Gga285	Gga286	Gga287	Gga288	Gga289	Gga290	Gga291	Gga292	Gga293	Gga294	Gga295	Gga296	Gga297	Gga298	Gga299	Gga300	Gga301	Gga302	Gga303	Gga304	Gga305	Gga306	Gga307	Gga308	Gga309	Gga310	Gga311	Gga312	Gga313	Gga314	Gga315	Gga316	Gga317	Gga318	Gga319	Gga320	Gga321	Gga322	Gga323	Gga324	Gga325	Gga326	Gga327	Gga328	Gga329	Gga330	Gga331	Gga332	Gga333	Gga334	Gga335	Gga336	Gga337	Gga338	Gga339	Gga340	Gga341	Gga342	Gga343	Gga344	Gga345	Gga346	Gga347	Gga348	Gga349	Gga350	Gga351	Gga352	Gga353	Gga354	Gga355	Gga356	Gga357	Gga358	Gga359	Gga360	Gga361	Gga362	Gga363	Gga364	Gga365	Gga366	Gga367	Gga368	Gga369	Gga370	Gga371	Gga372	Gga373	Gga374	Gga375	Gga376	Gga377	Gga378	Gga379	Gga380	Gga381	Gga382	Gga383	Gga384	Gga385	Gga386	Gga387	Gga388	Gga389	Gga390	Gga391	Gga392	Gga393	Gga394	Gga395	Gga396	Gga397	Gga398	Gga399	Gga400	Gga401	Gga402	Gga403	Gga404	Gga405	Gga406	Gga407	Gga408	Gga409	Gga410	Gga411	Gga412	Gga413	Gga414	Gga415	Gga416	Gga417	Gga418	Gga419	Gga420	Gga421	Gga422	Gga423	Gga424	Gga425	Gga426	Gga427	Gga428	Gga429	Gga430	Gga431	Gga432	Gga433	Gga434	Gga435	Gga436	Gga437	Gga438	Gga439	Gga440	Gga441	Gga442	Gga443	Gga444	Gga445	Gga446	Gga447	Gga448	Gga449	Gga450	Gga451	Gga452	Gga453	Gga454	Gga455	Gga456	Gga457	Gga458	Gga459	Gga460	Gga461	Gga462	Gga463	Gga464	Gga465	Gga466	Gga467	Gga468	Gga469	Gga470	Gga471	Gga472	Gga473	Gga474	Gga475	Gga476	Gga477	Gga478	Gga479	Gga480	Gga481	Gga482	Gga483	Gga484	Gga485	Gga486	Gga487	Gga488	Gga489	Gga490	Gga491	Gga492	Gga493	Gga494	Gga495	Gga496	Gga497	Gga498	Gga499	Gga500	Gga501	Gga502	Gga503	Gga504	Gga505	Gga506	Gga507	Gga508	Gga509	Gga510	Gga511	Gga512	Gga513	Gga514	Gga515	Gga516	Gga517	Gga518	Gga519	Gga520	Gga521	Gga522	Gga523	Gga524	Gga525	Gga526	Gga527	Gga528	Gga529	Gga530	Gga531	Gga532	Gga533	Gga534	Gga535	Gga536	Gga537	Gga538	Gga539	Gga540	Gga541	Gga542	Gga543	Gga544	Gga545	Gga546	Gga547	Gga548	Gga549	Gga550	Gga551	Gga552	Gga553	Gga554	Gga555	Gga556	Gga557	Gga558	Gga559	Gga560	Gga561	Gga562	Gga563	Gga564	Gga565	Gga566	Gga567	Gga568	Gga569	Gga570	Gga571	Gga572	Gga573	Gga574	Gga575	Gga576	Gga577	Gga578	Gga579	Gga580	Gga581	Gga582	Gga583	Gga584	Gga585	Gga586	Gga587	Gga588	Gga589	Gga590	Gga591	Gga592	Gga593	Gga594	Gga595	Gga596	Gga597	Gga598	Gga599	Gga600	Gga601	Gga602	Gga603	Gga604	Gga605	Gga606	Gga607	Gga608	Gga609	Gga610	Gga611	Gga612	Gga613	Gga614	Gga615	Gga616	Gga617	Gga618	Gga619	Gga620	Gga621	Gga622	Gga623	Gga624	Gga625	Gga626	Gga627	Gga628	Gga629	Gga630	Gga631	Gga632	Gga633	Gga634	Gga635	Gga636	Gga637	Gga638	Gga639	Gga640	Gga641	Gga642	Gga643	Gga644	Gga645	Gga646	Gga647	Gga648	Gga649	Gga650	Gga651	Gga652	Gga653	Gga654	Gga655	Gga656	Gga657	Gga658	Gga659	Gga660	Gga661	Gga662	Gga663	Gga664	Gga665	Gga666	Gga667	Gga668	Gga669	Gga670	Gga671	Gga672	Gga673	Gga674	Gga675	Gga676	Gga677	Gga678	Gga679	Gga680	Gga681	Gga682	Gga683	Gga684	Gga685	Gga686	Gga687	Gga688	Gga689	Gga690	Gga691	Gga692	Gga693	Gga694	Gga695	Gga696	Gga697	Gga698	Gga699	Gga700	Gga701	Gga702	Gga703	Gga704	Gga705	Gga706	Gga707	Gga708	Gga709	Gga710	Gga711	Gga712	Gga713	Gga714	Gga715	Gga716	Gga717	Gga718	Gga719	Gga720	Gga721	Gga722	Gga723	Gga724	Gga725	Gga726	Gga727	Gga728	Gga729	Gga730	Gga731	Gga732	Gga733	Gga734	Gga735	Gga736	Gga737	Gga738	Gga739	Gga740	Gga741	Gga742	Gga743	Gga744	Gga745	Gga746	Gga747	Gga748	Gga749	Gga750	Gga751	Gga752	Gga753	Gga754	Gga755	Gga756	Gga757	Gga758	Gga759	Gga760	Gga761	Gga762	Gga763	Gga764	Gga765	Gga766	Gga767	Gga768	Gga769	Gga770	Gga771	Gga772	Gga773	Gga774	Gga775	Gga776	Gga777	Gga778	Gga779	Gga780	Gga781	Gga782	Gga783	Gga784	Gga785	Gga786	Gga787	Gga788	Gga789	Gga790	Gga791	Gga792	Gga793	Gga794	Gga795	Gga796	Gga797	Gga798	Gga799	Gga800	Gga801	Gga802	Gga803	Gga804	Gga805	Gga806	Gga807	Gga808	Gga809	Gga810	Gga811	Gga812	Gga813	Gga814	Gga815	Gga816	Gga817	Gga818	Gga819	Gga820	Gga821	Gga822	Gga823	Gga824	Gga825	Gga826	Gga827	Gga828	Gga829	Gga830	Gga831	Gga832	Gga833	Gga834	Gga835	Gga836	Gga837	Gga838	Gga839	Gga840	Gga841	Gga842	Gga843	Gga844	Gga845	Gga846	Gga847	Gga848	Gga849	Gga850	Gga851	Gga852	Gga853	Gga854	Gga855	Gga856	Gga857	Gga858	Gga859	Gga860	Gga861	Gga862	Gga863	Gga864	Gga865	Gga866	Gga867	Gga868	Gga869	Gga870	Gga871	Gga872	Gga873	Gga874	Gga875	Gga876	Gga877	Gga878	Gga879	Gga880	Gga881	Gga882	Gga883	Gga884	Gga885	Gga886	Gga887	Gga888	Gga889	Gga890	Gga891	Gga892	Gga893	Gga894	Gga895	Gga896	Gga897	Gga898	Gga899	Gga900	Gga901	Gga902	Gga903	Gga904	Gga905	Gga906	Gga907	Gga908	Gga909	Gga910	Gga911	Gga912	Gga913	Gga914	Gga915	Gga916	Gga917	Gga918	Gga919	Gga920	Gga921	Gga922	Gga923	Gga924	Gga925	Gga926	Gga927	Gga928	Gga929	Gga930	Gga931	Gga932	Gga933	Gga934	Gga935	Gga936	Gga937	Gga938	Gga939	Gga940	Gga941	Gga942	Gga943	Gga944	Gga945	Gga946	Gga947	Gga948	Gga949	Gga950	Gga951	Gga952	Gga953	Gga954	Gga955	Gga956	Gga957	Gga958	Gga959	Gga960	Gga961	Gga962	Gga963	Gga964	Gga965	Gga966	Gga967	Gga968	Gga969	Gga970	Gga971	Gga972	Gga973	Gga974	Gga975	Gga976	Gga977	Gga978	Gga979	Gga980	Gga981	Gga982	Gga983	Gga984	Gga985	Gga986	Gga987	Gga988	Gga989	Gga990	Gga991	Gga992	Gga993	Gga994	Gga995	Gga996	Gga997	Gga998	Gga999	Gga1000	Gga1001	Gga1002	Gga1003	Gga1004	Gga1005	Gga1006	Gga1007	Gga1008	Gga1009	Gga1010	Gga1011	Gga1012	Gga1013	Gga1014	Gga1015	Gga1016	Gga1017	Gga1018	Gga1019	Gga1020	Gga1021	Gga1022	Gga1023	Gga1024	Gga1025	Gga1026	Gga1027	Gga1028	Gga1029	Gga1030	Gga1031	Gga1032	Gga1033	Gga1034	Gga1035	Gga1036	Gga1037	Gga1038	Gga1039	Gga1040	Gga1041	Gga1042	Gga1043	Gga1044	Gga1045	Gga1046	Gga1047	Gga1048	Gga1049	Gga1050	Gga1051	Gga1052	Gga1053	Gga1054	Gga1055	Gga1056	Gga1057	Gga1058	Gga1059	Gga1060	Gga1061	Gga1062	Gga1063	Gga1064	Gga1065	Gga1066	Gga1067	Gga1068	Gga1069	Gga1070	Gga1071	Gga1072	Gga1073	Gga1074	Gga1075	Gga1076	Gga1077	Gga1078	Gga1079	Gga1080	Gga1081	Gga1082	Gga1083	Gga1084	Gga1085	Gga1086	Gga1087	Gga1088	Gga1089	Gga1090	Gga1091	Gga1092	Gga1093	Gga1094	Gga1095	Gga1096	Gga1097	Gga1098	Gga1099	Gga1100	Gga1101	Gga1102	Gga1103	Gga1104	Gga1105	Gga1106	Gga1107	Gga1108	Gga1109	Gga1110	Gga1111	Gga1112	Gga1113	Gga1114	Gga1115	Gga1116	Gga1117	Gga1118	Gga1119	Gga1120	Gga1121	Gga1122	Gga1123	Gga1124	Gga1125	Gga1126	Gga1127	Gga1128	Gga1129	Gga1130	Gga1131	Gga1132	Gga1133	Gga1134	Gga1135	Gga1136	Gga1137	Gga1138	Gga1139	Gga1140	Gga1141	Gga1142	Gga1143	Gga1144	Gga1145	Gga1146	Gga1147	Gga1148	Gga1149	Gga1150	Gga1151	Gga1152	Gga1153	Gga1154	Gga1155	Gga1156	Gga1157	Gga1158	Gga1159	Gga1160	Gga1161	Gga1162	Gga1163	Gga1164	Gga1165	Gga1166	Gga1167	Gga1168	Gga1169	Gga1170	Gga1171	Gga1172	Gga1173	Gga1174	Gga1175	Gga1176	Gga1177	Gga1178	Gga1179	Gga1180	Gga1181	Gga1182	Gga1183	Gga1184	Gga1185	Gga1186	Gga1187	Gga1188	Gga1189	Gga1190	Gga1191	Gga1192	Gga1193	Gga1194	Gga1195	Gga1196	Gga1197	Gga1198	Gga1199	Gga1200	Gga1201	Gga1202	Gga1203	Gga1204	Gga1205	Gga1206	Gga1207	Gga1208	Gga1209	Gga1210	Gga1211	Gga1212	Gga1213	Gga1214	Gga1215	Gga1216	Gga1217	Gga1218	Gga1219	Gga1220	Gga1221	Gga1222	Gga1223	Gga1224	Gga1225	Gga1226	Gga1227	Gga1228	Gga1229	Gga1230	Gga1231	Gga1232	Gga1233	Gga1234	Gga1235	Gga1236	Gga1237	Gga1238	Gga1239	Gga1240	Gga1241	Gga1242	Gga1243	Gga1244	Gga1245	Gga1246	Gga1247	Gga1248	Gga1249	Gga1250	Gga1251	Gga1252	Gga1253	Gga1254	Gga1255	Gga1256	Gga1257	Gga1258	Gga1259	Gga1260	Gga1261	Gga1262	Gga1263	Gga1264	Gga1265	Gga1266	Gga1267	Gga1268	Gga1269	Gga1270	Gga1271	Gga1272	Gga1273	Gga1274	Gga1275	Gga1276	Gga1277	Gga1278	Gga1279	Gga1280	Gga1281	Gga1282	Gga1283	Gga1284	Gga1285	Gga1286	Gga1287	Gga1288	Gga1289	Gga1290	Gga1291	Gga1292	Gga1293	Gga1294	Gga1295	Gga1296	Gga1297	Gga1298	Gga1299	Gga1300	Gga1301	Gga1302	Gga1303	Gga1304	Gga1305	Gga1306	Gga1307	Gga1308	Gga1309	Gga1310	Gga1311	Gga1312	Gga1313	Gga1314	Gga1315	Gga1316	Gga1317	Gga1318	Gga1319	Gga1320	Gga132
---------------------	----------------------	----------	------	---------	-------	------------------	---------	------	------	------	------	------	------	------	------	------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	--------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	--------

CVL blocks in proto-chromosomes of the osteichthyan ancestor.

60

CVL blocks in proto-chromosomes of the amniote ancestor.

61

Supplementary Table S4

Effect of changing key parameter values on reconstruction of the vertebrate ancestral genome

Case	1st parameter	2nd parameter	3rd parameter	Number of CVL blocks	Number of vertebrate groups	Number of quadruplicated vertebrate groups	Number of gnathostome subgroups	Number of osteichthyan proto-chromosomes	Number of amniote proto-chromosomes
1	10	1.E-04	1.E-02	118	10	5	34	31	26
2	7	1.E-04	1.E-02	139	9	6	32	28	24
3	13	1.E-04	1.E-02	105	10	4	35	28	26
4	10	5.E-04	1.E-02	118	9	6	32	31	26
5	10	5.E-05	1.E-02	118	10	5	34	30	26
6	10	1.E-04	5.E-02	121	9	5	31	28	25
7	10	1.E-04	1.E-03	114	10	6	35	30	25

1st parameter: the minimum threshold on the number of genes in a CVL block

2nd parameter: the maximum threshold on the significance (probability) that two CVL blocks are paralogous

3rd parameter: the maximum threshold for the Mann-Whitney U-test to decide whether a CVL block is divided

11 References

- Benton, M.J. and Donoghue, P.C.J. 2007. Paleontological evidence to date the tree of Life. *Mol. Biol. Evol.* **24**: 26-53.
- Blair, J.E. and Hedges, S.B. 2005. Molecular phylogeny and divergence times of deuterostome animals. *Mol. Biol. Evol.* **22**: 2275-2284.
- Blomme, T., Vandepoele, K., De Bodt, S., Simillion, C., Maere, S., and Van de Peer, Y. 2006. The gain and loss of genes during 600 million years of vertebrate evolution. *Genome Biol.* **7**: R43.
- Bourque, G., Zdobnov, E.M., Bork, P., Pevzner, P.A., and Tesler, G. 2005. Comparative architectures of mammalian and chicken genomes reveal highly variable rates of genomic rearrangements across different lineages. *Genome Res.* **15**: 98-110.
- Crow, K.D., Stadler, P.F., Lynch, V.J., Amemiya, C., and Wagner, G.P. 2006. The "fish-specific" Hox cluster duplication is coincident with the origin of teleosts. *Mol. Biol. Evol.* **23**: 121-136.
- Dehal, P. and Boore, J.L. 2005. Two rounds of whole genome duplication in the ancestral vertebrate. *PLoS Biol.* **3**: e314.
- Dehal, P.S. and Boore, J.L. 2006. A phylogenomic gene cluster resource: the Phylogenetically Inferred Groups (PhIGs) database. *BMC Bioinformatics* **7**: 201.
- Hedges, S.B. and Poling, L.L. 1999. A molecular phylogeny of reptiles. *Science* **283**: 998-1001.
- Hoegg, S., Brinkmann, H., Taylor, J.S., and Meyer, A. 2004. Phylogenetic timing of the fish-specific genome duplication correlates with the diversification of teleost fish. *J. Mol. Evol.* **59**: 190-203.
- Inoue, J.G., Miya, M., Venkatesh, B., and Nishida, M. 2005. The mitochondrial genome of Indonesian coelacanth *Latimeria menadoensis* (Sarcopterygii: Coelacanthiformes) and divergence time estimation between the two coelacanths. *Gene* **349**: 227-235.
- Jaillon, O., Aury, J.M., Brunet, F., Petit, J.L., Stange-Thomann, N., Mauceli, E., Bouneau, L., Fischer, C., Ozouf-Costaz, C., Bernot, A. et al. 2004. Genome duplication in the teleost fish *Tetraodon nigroviridis* reveals the early vertebrate proto-karyotype. *Nature* **431**: 946-957.
- Kellis, M., Birren, B.W., and Lander, E.S. 2004. Proof and evolutionary analysis of ancient genome duplication in the yeast *Saccharomyces cerevisiae*. *Nature* **428**: 617-624.
- Kohn, M., Hogel, J., Vogel, W., Minich, P., Kehrer-Sawatzki, H., Graves, J.A., and Hameister, H. 2006. Reconstruction of a 450-My-old ancestral vertebrate protokaryotype. *Trends Genet.* **22**: 203-210.
- Meyer, A. and Zardoya, R. 2003. Recent advances in the (molecular) phylogeny of vertebrates. *Annu. Rev. Ecol. Evol. S.* **34**: 311-338.

- Naruse, K., Tanaka, M., Mita, K., Shima, A., Postlethwait, J., and Mitani, H. 2004. A medaka gene map: the trace of ancestral vertebrate proto-chromosomes revealed by comparative gene mapping. *Genome Res.* **14**: 820-828.
- Postlethwait, J.H., Woods, I.G., Ngo-Hazelett, P., Yan, Y.L., Kelly, P.D., Chu, F., Huang, H., Hill-Force, A., and Talbot, W.S. 2000. Zebrafish comparative genomics and the origins of vertebrate chromosomes. *Genome Res.* **10**: 1890-1902.
- Springer, M.S., Murphy, W.J., Eizirik, E., and O'Brien, S.J. 2003. Placental mammal diversification and the Cretaceous-Tertiary boundary. *Proc. Natl Acad. Sci. USA* **100**: 1056-1061.
- Stadler, P.F., Fried, C., Prohaska, S.J., Bailey, W.J., Misof, B.Y., Ruddle, F.H., and Wagner, G.P. 2004. Evidence for independent Hox gene duplications in the hagfish lineage: a PCR-based gene inventory of *Eptatretus stoutii*. *Mol. Phylogenet. Evol.* **32**: 686-694.
- Woodburne, M.O., Rich, T.H., and Springer, M.S. 2003. The evolution of tribospheny and the antiquity of mammalian clades. *Mol. Phylogenet. Evol.* **28**: 360-385.
- Woods, I.G., Wilson, C., Friedlander, B., Chang, P., Reyes, D.K., Nix, R., Kelly, P.D., Chu, F., Postlethwait, J.H., and Talbot, W.S. 2005. The zebrafish gene map defines ancestral vertebrate chromosomes. *Genome Res.* **15**: 1307-1314.
- Yamanoue, Y., Miya, M., Inoue, J.G., Matsuura, K., and Nishida, M. 2006. The mitochondrial genome of spotted green pufferfish *Tetraodon nigroviridis* (Teleostei: Tetraodontiformes) and divergence time estimation among model organisms in fishes. *Genes Genet. Syst.* **81**: 29-39.