# Геном эукариот

## Н.Н. Колесников

Institute of Cytology and Genetics SB RAS, Novosibirsk, Russia

X       Y       4

21.9    19.9      40.9      3.1 1.2

2L       2R       3L       3R

22.2    6.2   12.1    20.3    23.4    9.2   8.3    27.9

▓ Heterochromatin
☐ Euchromatin
○ Centromere
▲ Physical map gap

# TABLE 2   Annotations in Release 3

| Description[a] | Euchromatin 116.8 Mb | 12 Mb of heterochromatin |
|---|---|---|
| Protein-coding genes | 13,379 | 297 |
| tRNA genes | 290 | 0 |
| microRNA genes | 23 | ND |
| snRNA genes | 34 | ND |
| snoRNA genes | 28 | ND |
| Pseudogenes | 17 | ND |
| Misc. noncoding RNA | 28 | 2 |
| rRNA genes | — | 6 |
| Transposons | 1,572 | ND |
| Total protein-coding genes | 13,379 | 297 |
| Total length of euchromatin/heterochromatin | 116.8 Mb | 12 Mb |
| Exons | 60,897 | 1109 |
| Protein-coding exons[b] | 54,934 | 999 |
| Length of genome in exons | 27.8 Mb (24%) | 382 kb (3%) |
| Introns | 48,257 | 803 |
| Genes with 5′ UTR | 10,224 (76%) | 152 (51%) |
| Transcripts with 5′ UTR | 14,706 (81%) | 215 (57%) |
| Average 5′ UTR length | 265 nucleotides | 217 nt |
| Genes with 3′ UTR | 9,646 (72%) | 119 (40%) |
| Transcripts with 3′ UTR | 14,012 (77%) | 172 (46%) |
| Average 3′ UTR length | 442 nucleotides | 311 nt |
| Average ratio of length of CDS/transcript[c] | 0.75 | 0.79 |
| Total protein-coding transcripts | 18,106 | 379 |
| Genes with alternative transcripts | 2,729 (20%) | 49 (16%) |
| Average number of transcripts per alternatively spliced gene | 2.75 | 2.8 |
| Total number alternative transcripts | 4,743 | 88 |
| Unique peptides | 15,848 | 351 |
| Gene-prediction data only | 815 | 49 |
| BLASTX/TBLASTX homologies | 11,936 | 167 |
| ESTs and DGC cDNA sequencing reads | 10,498 | 134 |
| GenBank accessions | 5,104 | 22 |
| ARGS (RefSeq) | 795 | 0 |
| Error report submissions | 825 | 7 |
| Full insert DGC cDNAs | 9,297 | 58 |

[a]Abbreviations: UTR, untranslated region; CDS, (protein)-coding sequence; R2, Release 2; R3, Release 3; ND, not determined. All statistics are for protein-coding genes only. The numbers reflect the FlyBase annotation database of
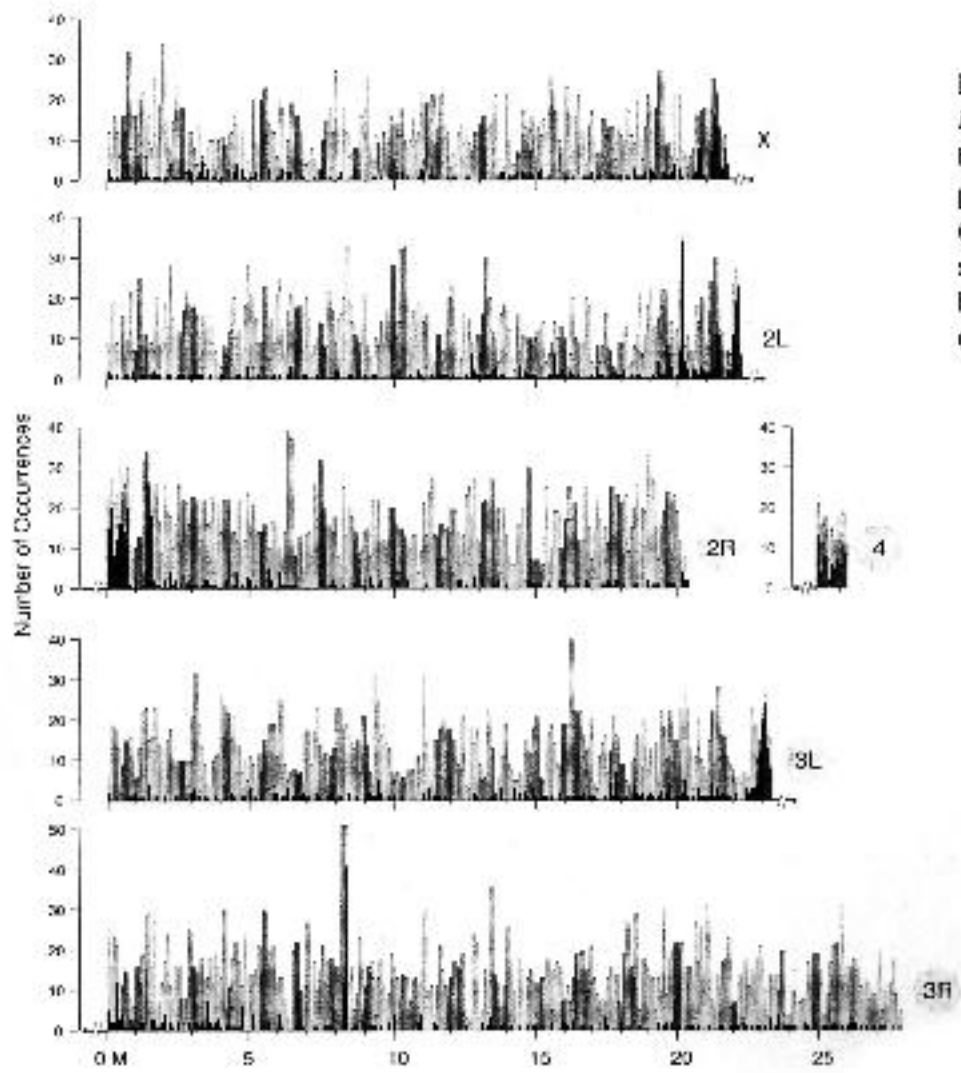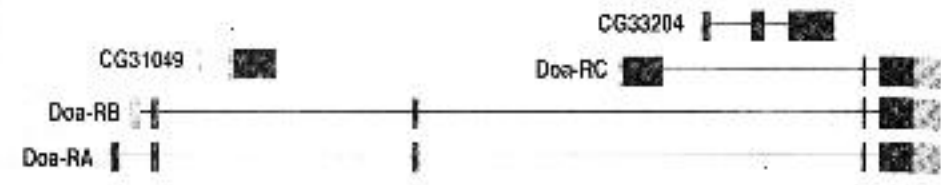
**Figure 2** Distribution of protein-coding genes and transposable elements in the *Drosophila melanogaster* euchromatic genome. Each chromosome arm is represented by a black horizontal line with a circle indicating its centromere. The number of transposable elements (*black*) and protein-coding genes (*gray*) is shown for 100-kilobase (kb) windows along each chromosome arm. Although there is local variation, gene density is rather uniform along the chromosome arms, whereas transposable elements are highly enriched at the centromere-proximal regions of each arm. A scale in megabases (Mbs) is shown at the bottom of the figure.

Spn1:CG9456

CG9455

B. Nested genes

CG33204

CG31049

Doa-RC

Doa-RB

Doa-RA

C. Interleaved genes

CG5500

ro:CG6348

D. Dicistronic gene

ORF1

ORF2

CG31188

E. Alternatively spliced genes

Vanaso:CG32315

alpha-Spec:CG1977

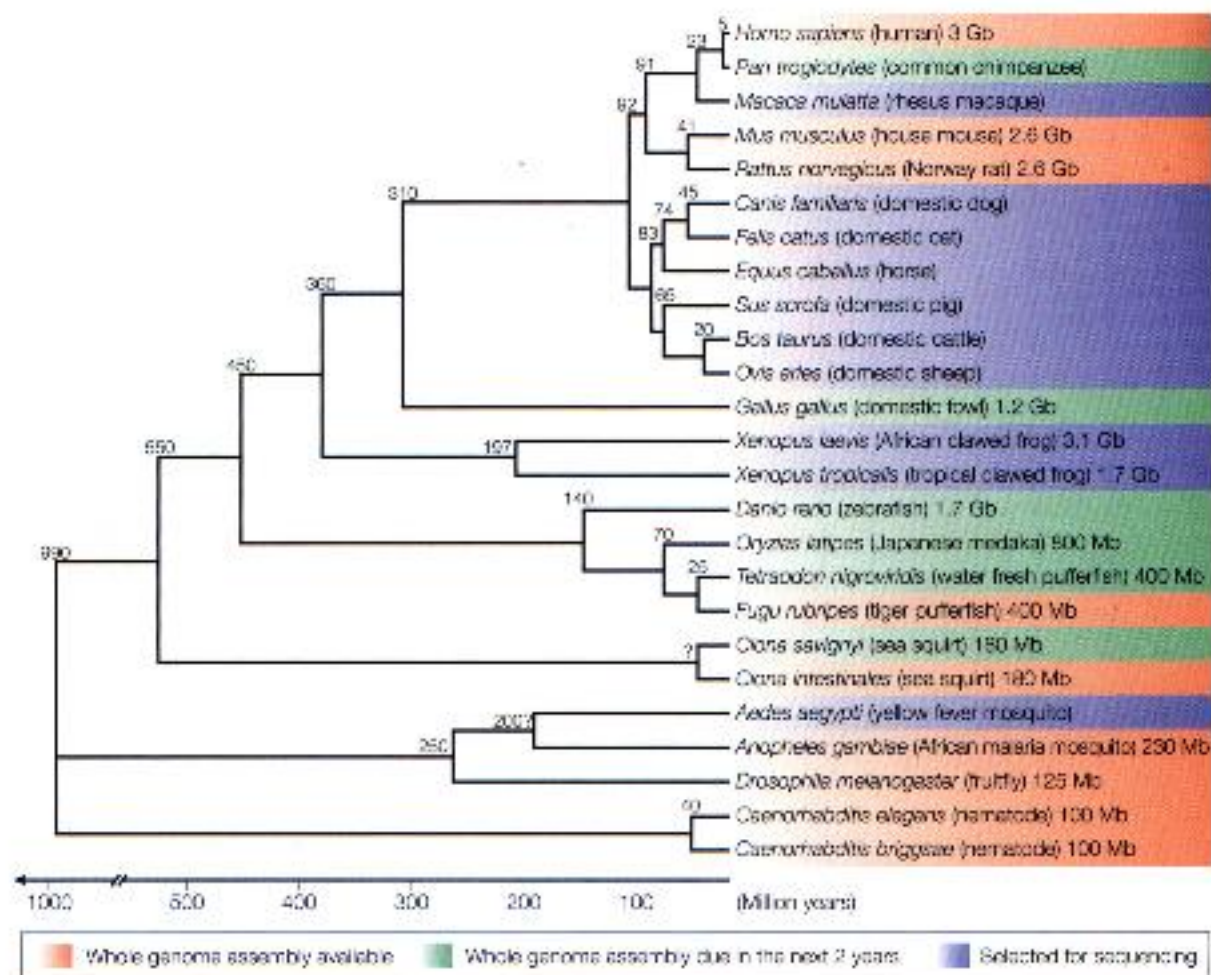F. Trans-spliced gene

mod(mdg4):CG7836

cis-splice

trans-splice

Figure 1 | **Evolutionary relationship between metazoans that are sequenced or due for sequencing.** The simplified phylogenetic relationships between the metazoans for which the complete, or nearly complete, genome sequences are available or will be available soon. Evolutionary distances (in million years) and genome sizes are based on REFS. 10,13,58–100.
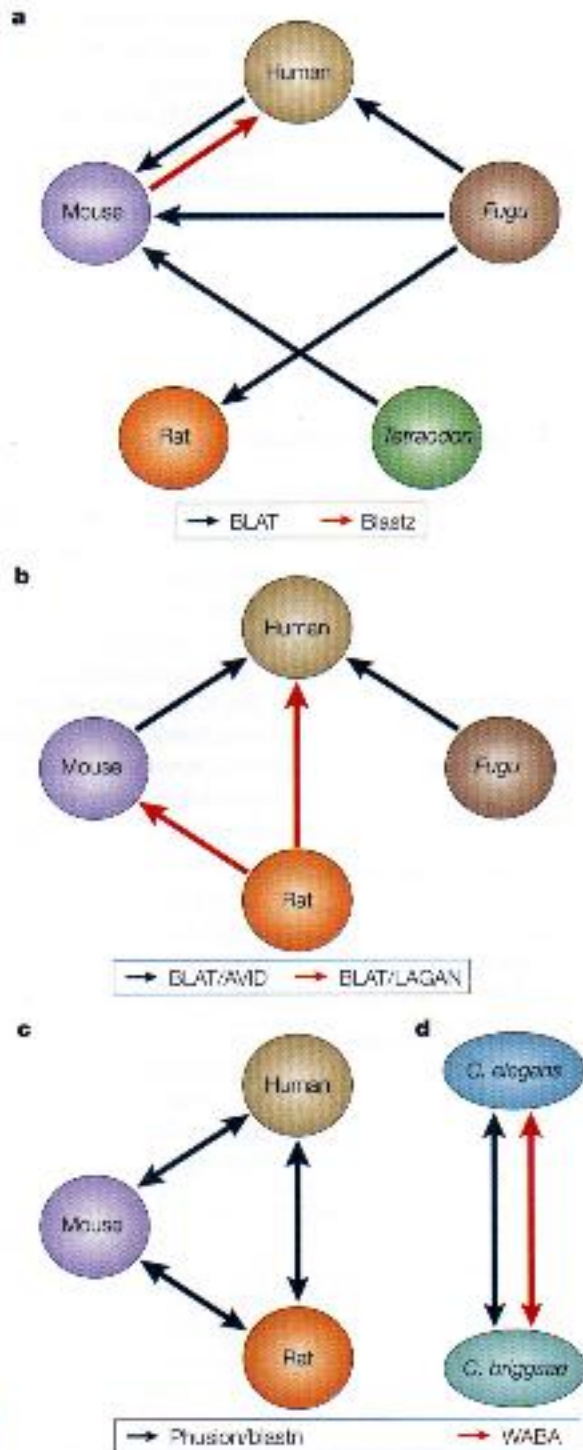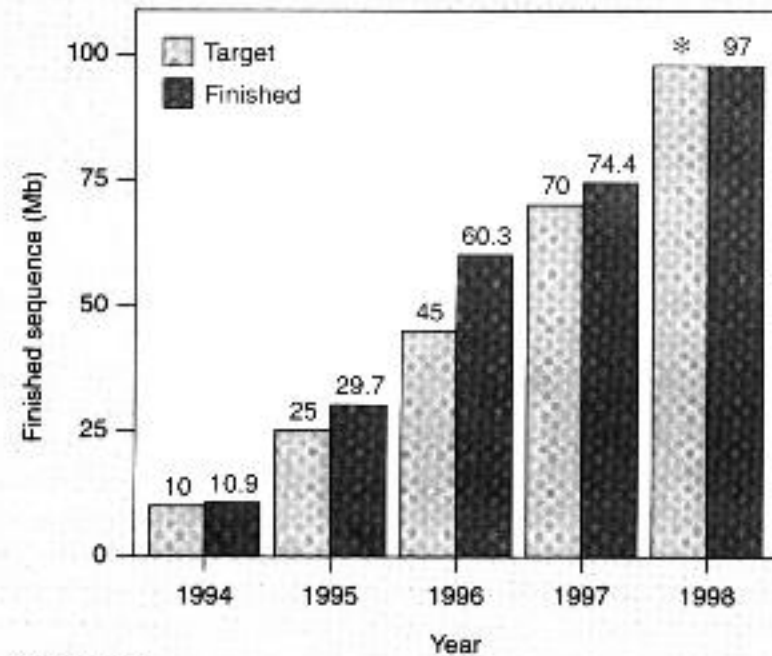
Figure 2 | **Whole-genome alignments available online.**
Whole-genome precomputed alignments are available from the web sites of several different groups. Each group uses a different alignment method. The figure gives an overview of what is provided by each group and the methods used. The direction of the arrows indicate the species that are used as the reference for the display. For example, in the rat–mouse pair in panel **b**, the arrow pointing to mouse indicates that rat alignments can be seen on the mouse genome (but not the opposite), and in the human–mouse pair in panel **c**, the arrows pointing to both species indicate that human alignments can been seen on the mouse genome, and vice versa. Alignments can be found at the University of California, Santa Cruz, UCSC Genome Browser and the Penn State Bioinformatics Group (**a**), the Berkeley Genome Pipeline and the ECR Browser (**b**), the Ensembl Genome Browser (**c** and **d**) and WormBase (**d**). The figure represents a snapshot of the alignments that are available at the time of writing.
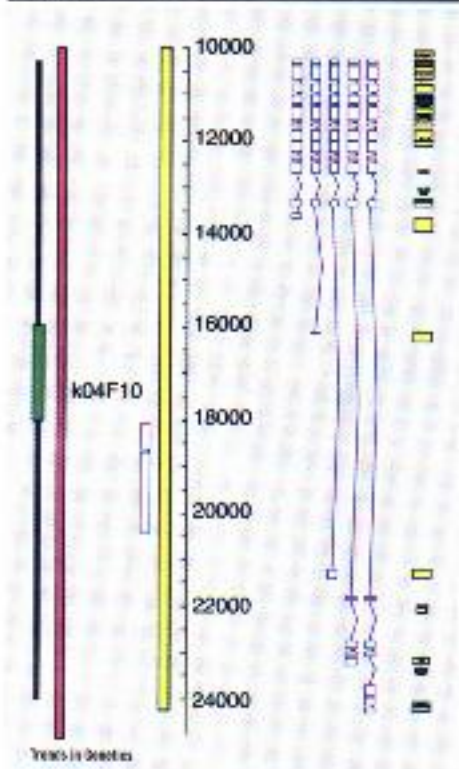
**FIGURE 1. Sequencing project**

Sequencing targets and progress in the *C. elegans* genome sequencing project. Asterisk indicates target completion date.
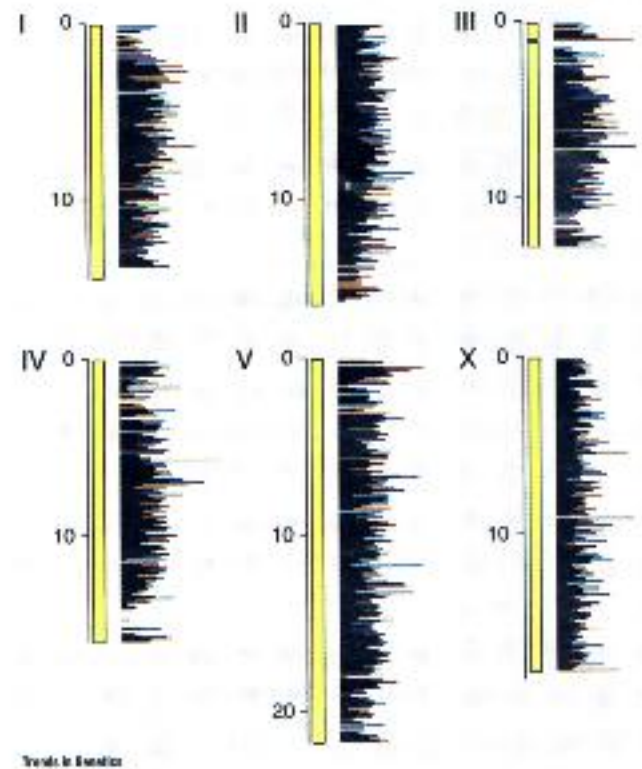
**FIGURE 4. Alternatively spliced genes**

An ACeDB sequence display of a portion of the C. elegans cosmid K04F10 shows the multiple, alternatively spliced forms of the gene bli-4. From left are bli-4E, bli-4A, bli-4B, bli-4C and bli-4D. All the alternatively spliced forms are confirmed either by EST matches or experimentally[a]. The displayed sequence features are described in the legend for Fig. L.
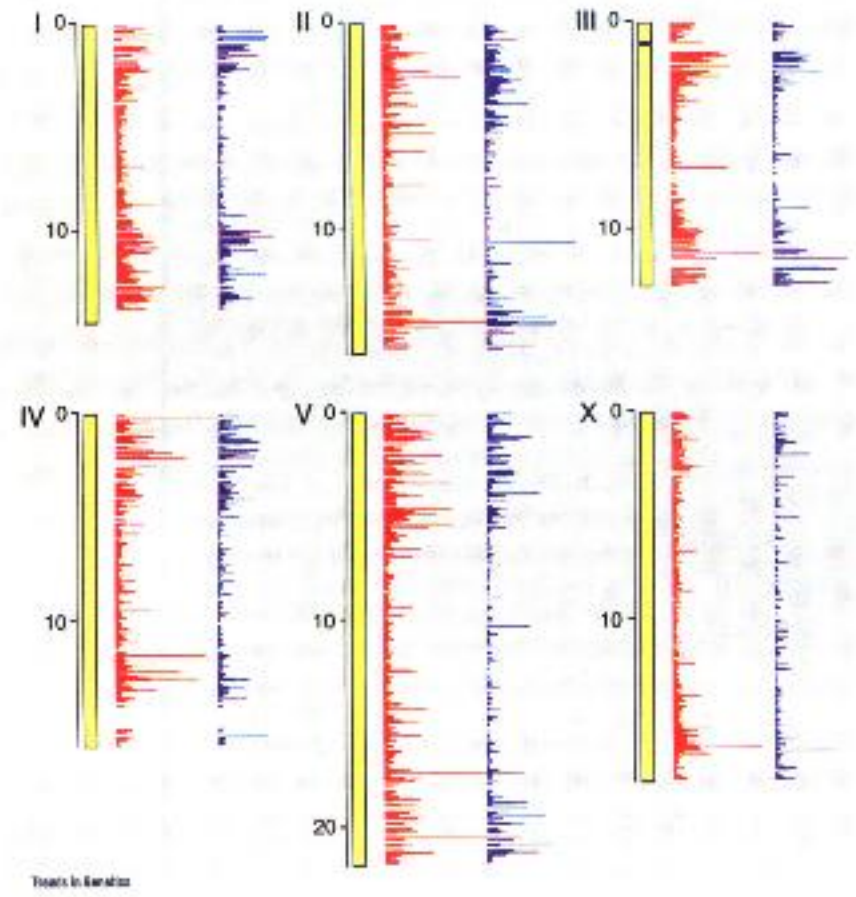
## FIGURE 5. Gene distribution

The distribution of predicted genes is plotted along each chromsome. The vertical yellow bars represent the clonal physical map of the genome (in Mbl.
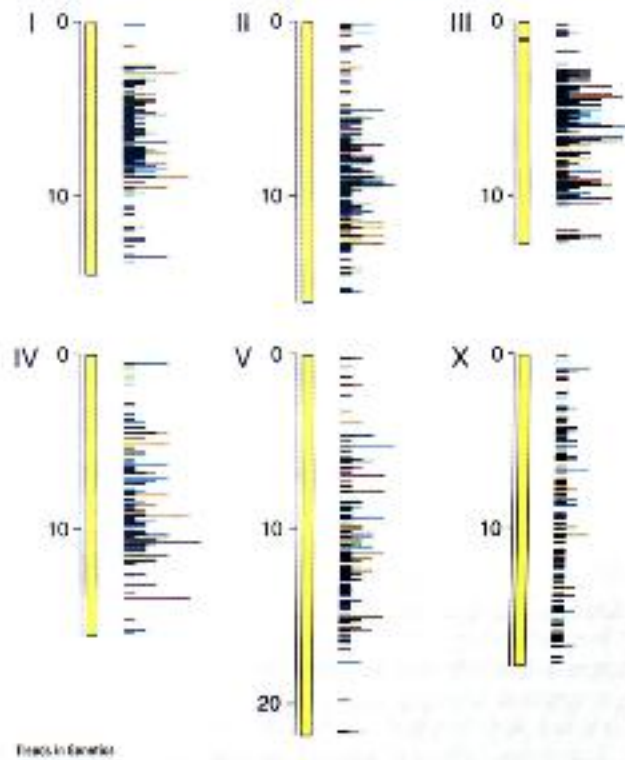
## FIGURE 6. Tandem and inverted repeats

Distribution of local tandem and inverted repeats along each of the chromosomes. Inverted repeats are shown in red while tandem repeats are blue. Both kinds of repeats are more frequent on the arms of the autosomes than in the central gene-rich regions, while they appear more uniformly distributed on the X chromosome.
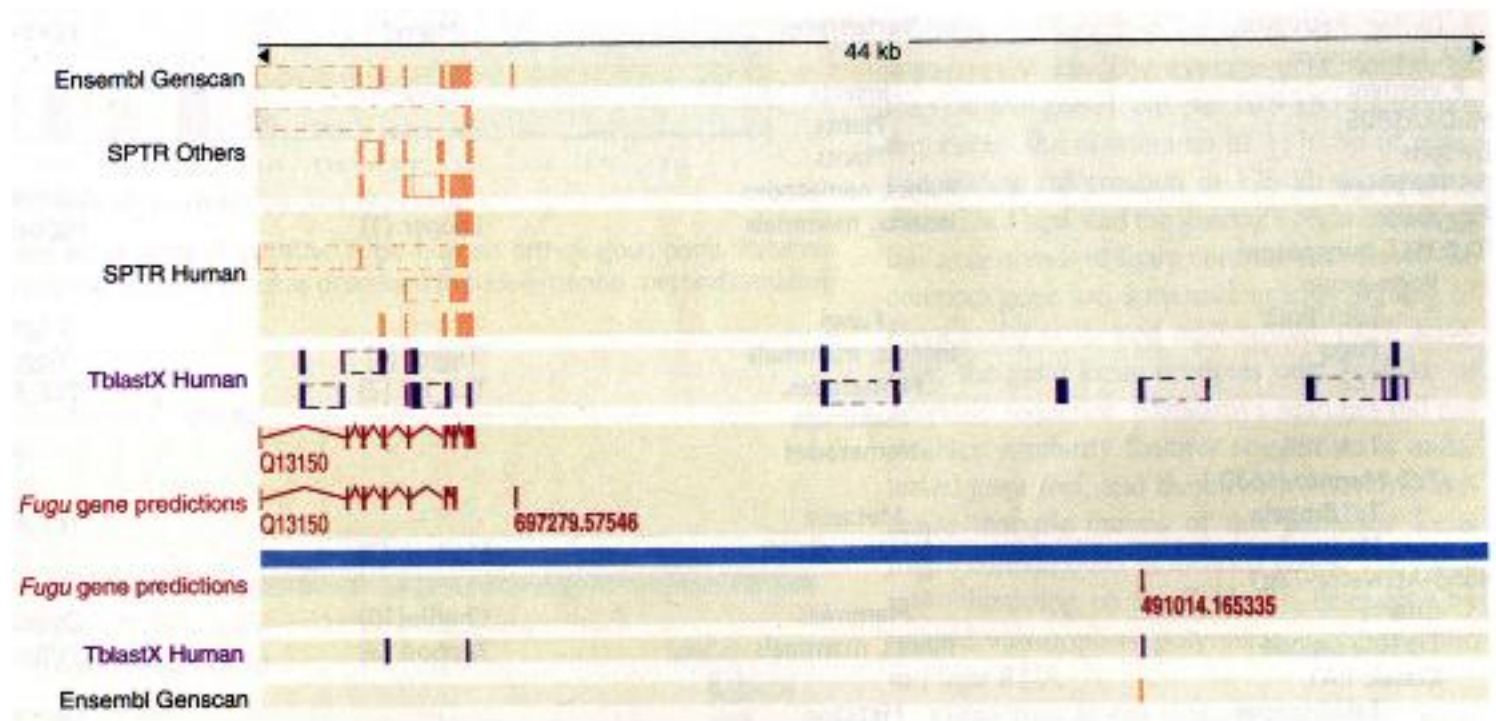
**FIGURE 8. Comparing *C. elegans* with *S. cerevisiae***

Distribution along each chromosome of the genes that are conserved between *S. cerevisiae* and *C. elegans*. These genes are clustered and coincide with the locations of genes with EST matches.
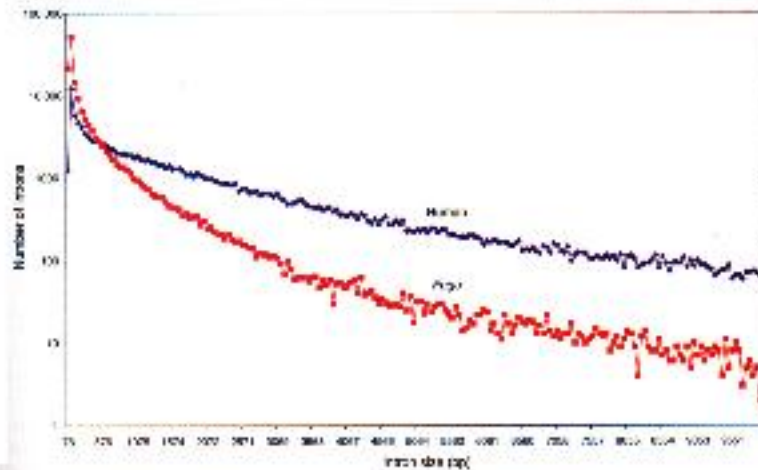
Fig. 2. Comparative frequency distribution of intron sizes in Fugu and human.
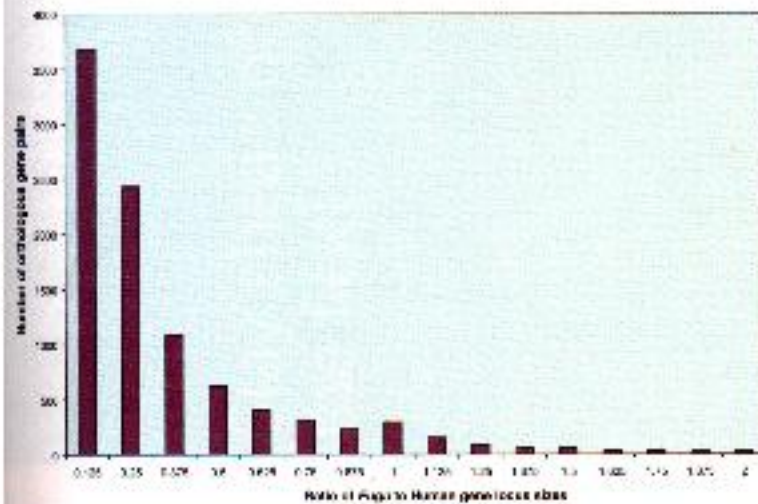


Fig. 3. Distribution of ratios for gene locus sizes of putative Fugu-human orthologous pairs. Putative Fugu-human orthologous gene pairings were determined as described in supplemental methods relating to conservation of synteny.
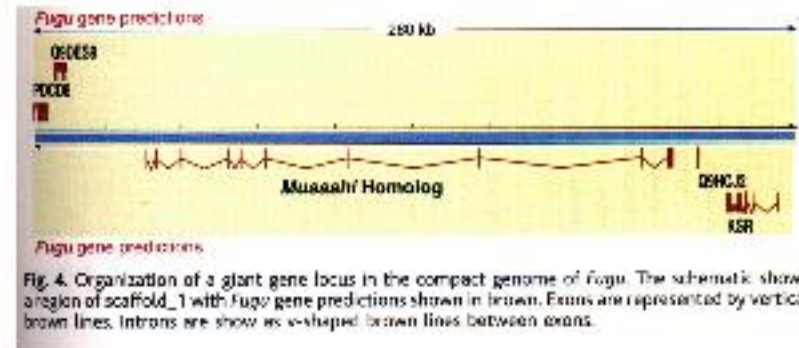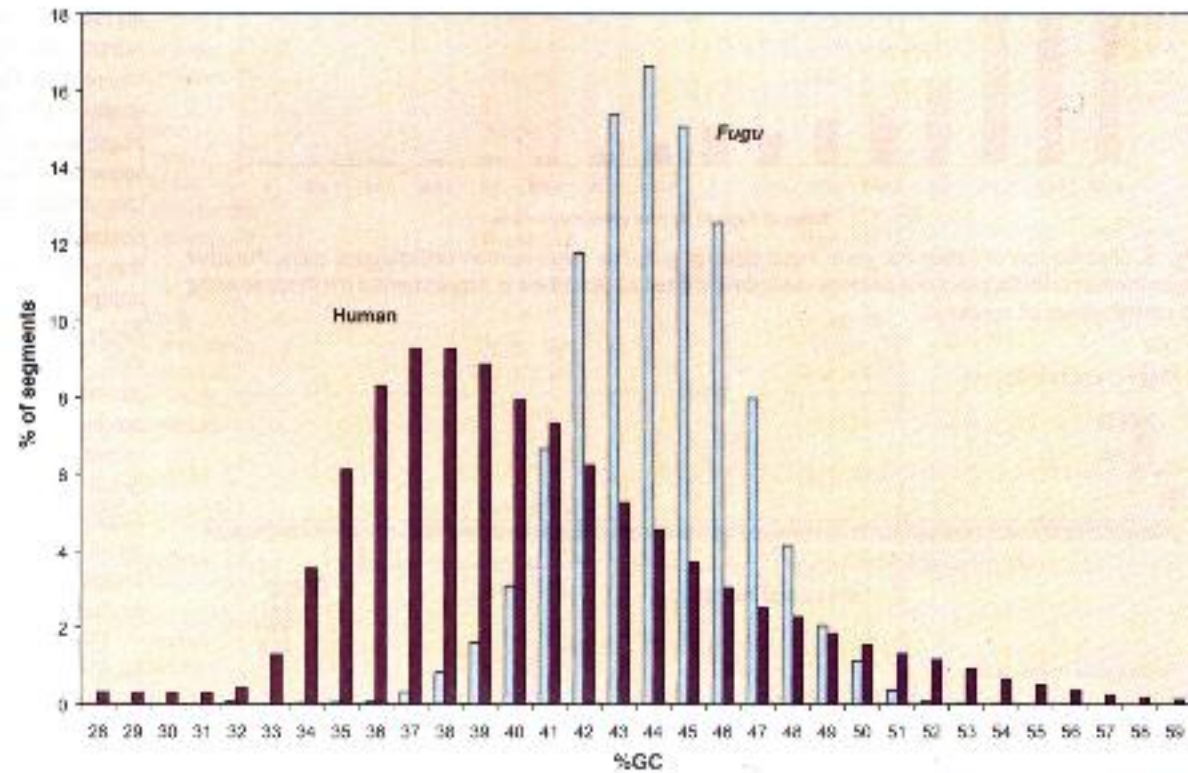


Fig. 4. Organization of a giant gene locus in the compact genome of Fugu. The schematic shows a region of scaffold_1 with Fugu gene predictions shown in brown. Exons are represented by vertical brown lines. Introns are show as v-shaped brown lines between exons.

**Fig. 5.** Distribution of GC content in the *Fugu* and human genomes. Sliding windows of 50 kb were used; similar conclusions were derived with windows of 25 and 100 kb (not shown).
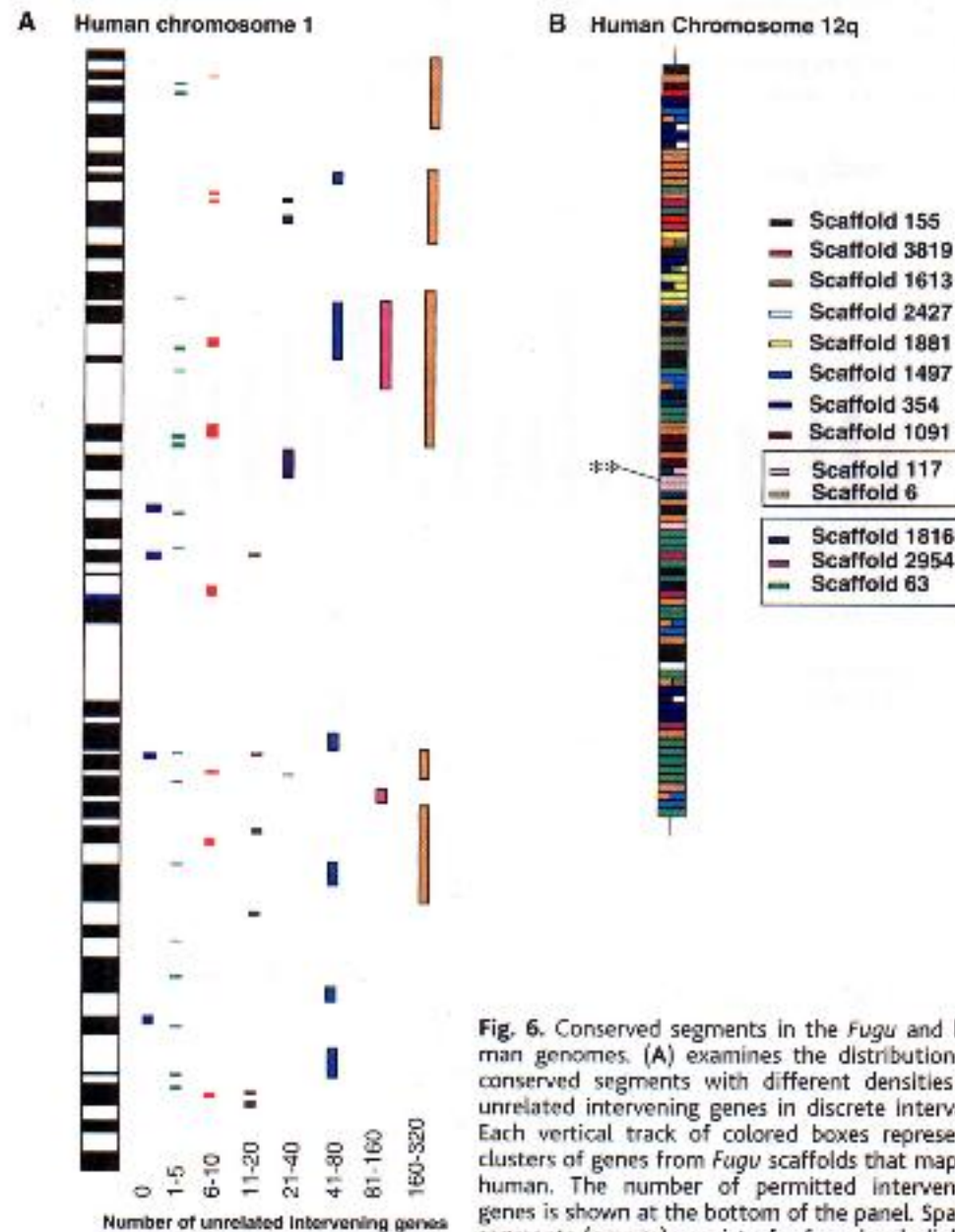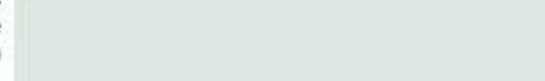
**A** Human chromosome 1

**B** Human Chromosome 12q

- ━ Scaffold 155
- ━ Scaffold 3819
- ━ Scaffold 1613
- ▭ Scaffold 2427
- ▭ Scaffold 1881
- ━ Scaffold 1497
- ━ Scaffold 354
- ━ Scaffold 1091

- ▭ Scaffold 117
- ▭ Scaffold 6

- ▭ Scaffold 1816
- ▭ Scaffold 2954
- ▭ Scaffold 63

Number of unrelated intervening genes

0  1-5  6-10  11-20  21-40  41-80  81-160  160-320

**Fig. 6.** Conserved segments in the *Fugu* and human genomes. (**A**) examines the distribution of conserved segments with different densities of unrelated intervening genes in discrete intervals. Each vertical track of colored boxes represents clusters of genes from *Fugu* scaffolds that map to human. The number of permitted intervening genes is shown at the bottom of the panel. Sparse segments (orange) consist of a few closely linked *Fugu* genes whose orthologs are spread over large chromosomal distances in human. Very sparse segments (>320 intervening genes) are not shown on this panel. (**B**) A detailed view of human chromosome 12q to illustrate shuffling gene order. Colored boxes represent individual genes from *Fugu* scaffolds whose orthologs on human chromosome 12q were determined through alignment by hand. The order of the orthologs along the human chromsome is shown, with the corresponding *Fugu* scaffold of origin in the key on the right. The scaffolds shown grouped together in boxes in the key are known to be linked in *Fugu*. ** indicates the position of the *Hox-c* complex on this chromosome, represented by scaffolds 117, 1327, and 1458 (the latter two are not shown in the key). Where a human gene has equally matching (co-orthologous) *Fugu* genes, this is shown as a double- or triple-colored box.
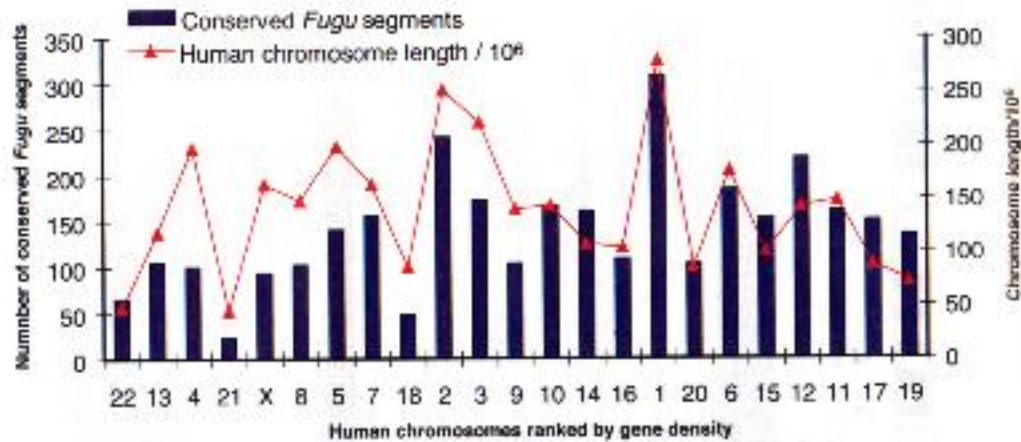
**Fig. 7.** Distribution of conserved segments of *Fugu* on human chromosomes ranked by gene density. The figure shows the relation between the number of conserved segments of *Fugu* on human chromosomes, the length of human chromosomes, and their gene density. Chromosome 22 is the most gene poor, chromosome 19 the most gene dense. There is no apparent relation between human chromosomal gene density and the number of segments. The distribution of conserved segments varies with human chromosomal length.
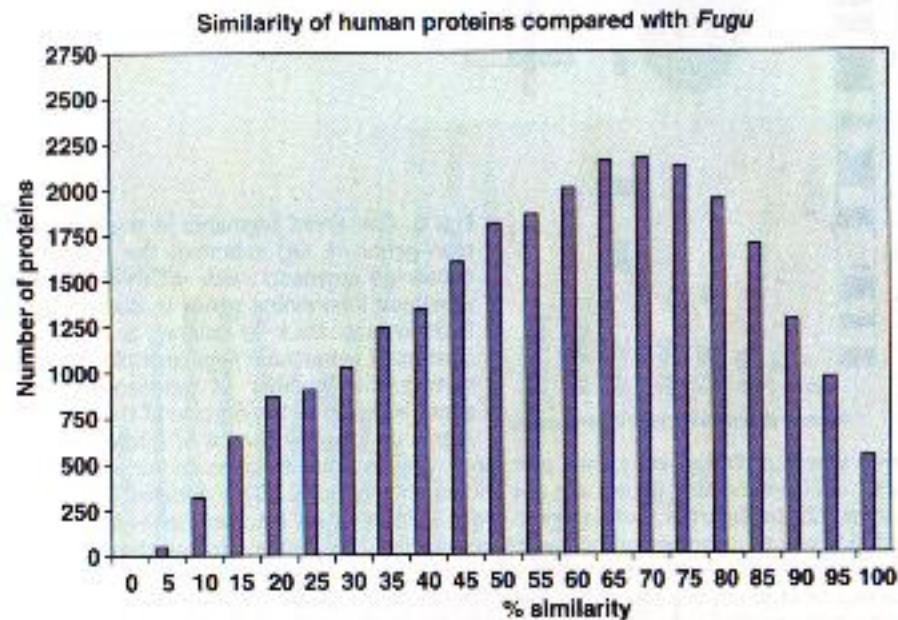


**Fig. 8.** Distribution of protein similarities between *Fugu* and human proteomes. Global similarities were calculated as the sum of similarities in all nonoverlapping HSPs using a BLOSUM62 matrix over the query (human) sequence length.
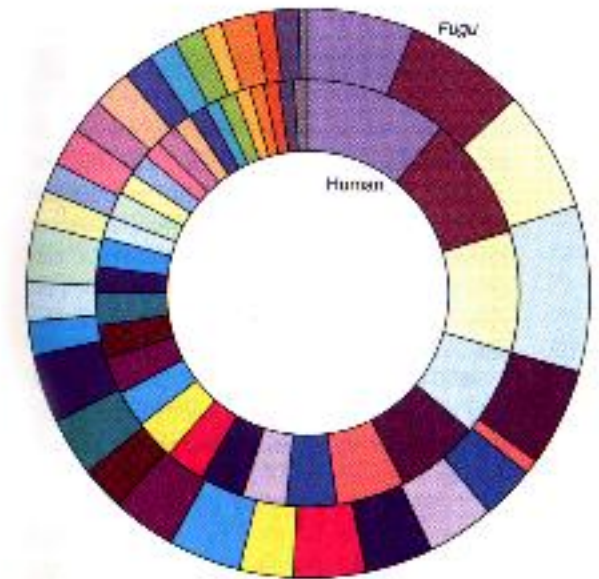


**Fig. 9.** Protein domains of *Fugu* and human. The schematic shows the number of gene loci with corresponding Interpro domains in *Fugu* and human for the 32 most populous families.
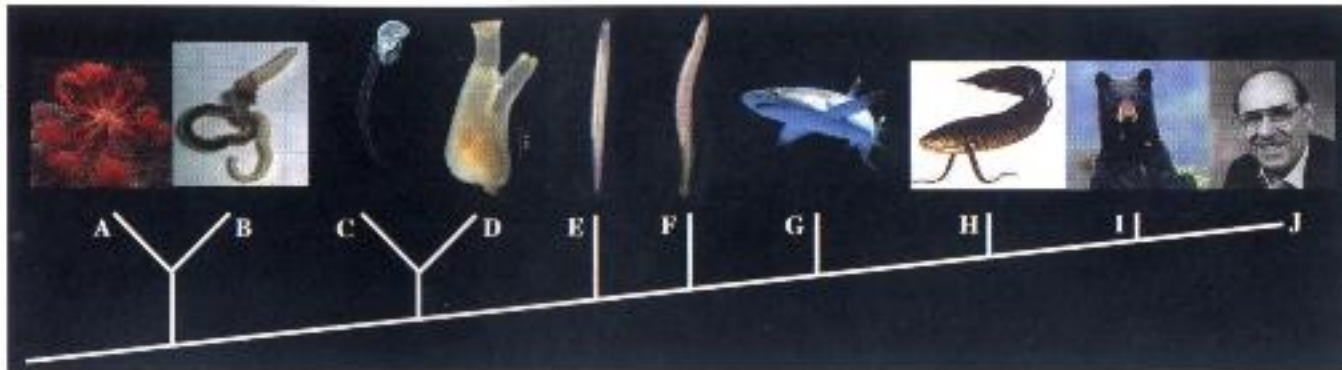
**Figure 1.** A deuterostome phylogeny. A, an echinoderm, a starfish. B, a hemichordate. C, a larvacean, *Oikopleura*. D, an ascidian, *Ciona*. E, a cephalochordate, amphioxus. F, a jawless fish, a lamprey. G, a gnathostome, a shark. H, a sarcopterygian fish, a lunglish. I, a mammal, a bear. J, a human, Yogi Berra.
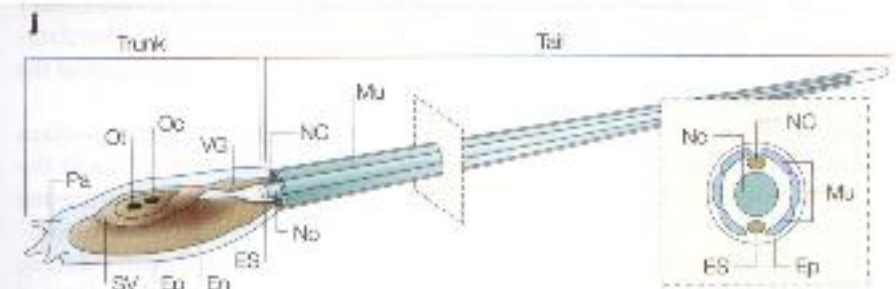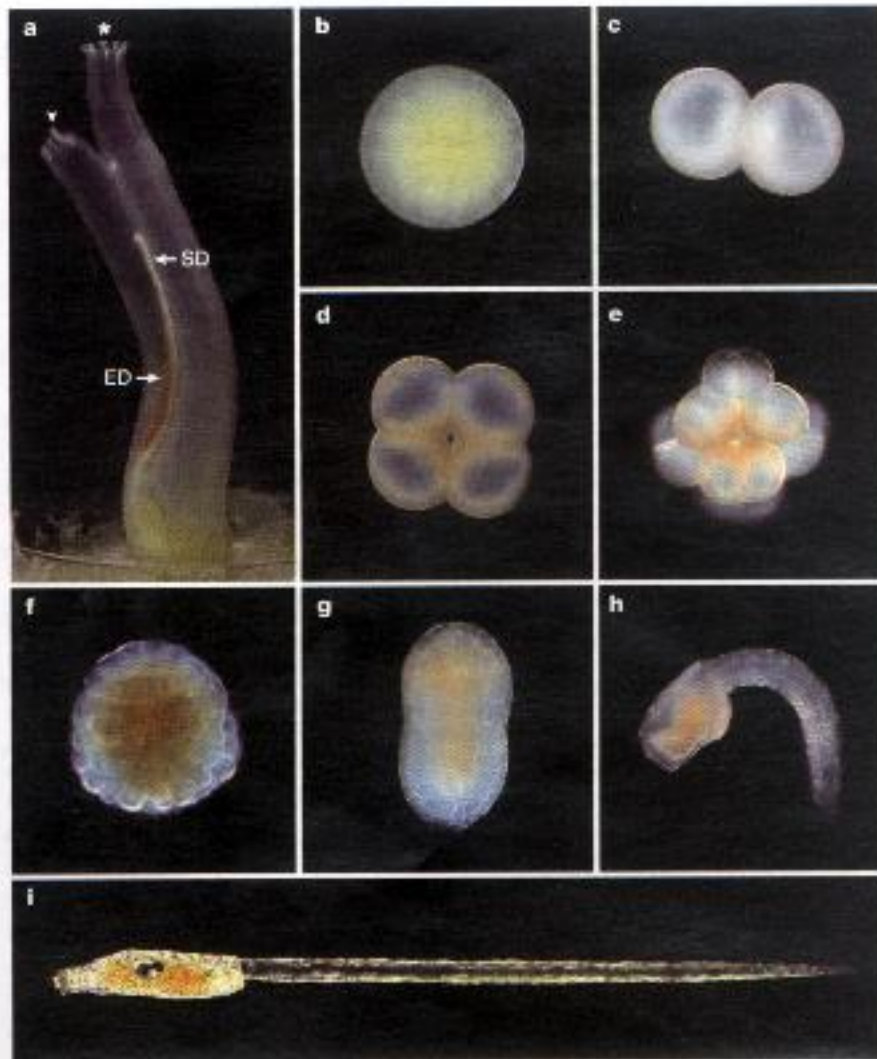
Figure 1 | Tadpole-type larva of the ascidian *Ciona intestinalis*. Panel **a** shows adults with incurrent siphons (arrowhead) and outcurrent siphons (asterisk) for circulating seawater from which the ascidians filter plankton and nutrient materials. The white duct is the sperm duct (SD), and the orange duct is the egg duct (ED). Panels **b–i** show the stages of embryogenesis: fertilized egg (**b**), 2-cell embryo (**c**), 4-cell embryo (**d**), 16-cell embryo (**e**), gastrula (**f**), early-tailbud embryo (**g**), mid-tailbud embryo (**h**) and tadpole larva (**i** and **j**). Embryos were dechorionated to show their outer morphology. Mesenchyme in the posterior region of the trunk is not shown. En, endoderm; Ep, epidermis; ES, endodermal strand; Mu, muscle; NC, nerve cord; No, notochord; Oc, ocellus; Ot, otolith; Pa, palps; SV, sensory vesicle; VG, visceral ganglion.
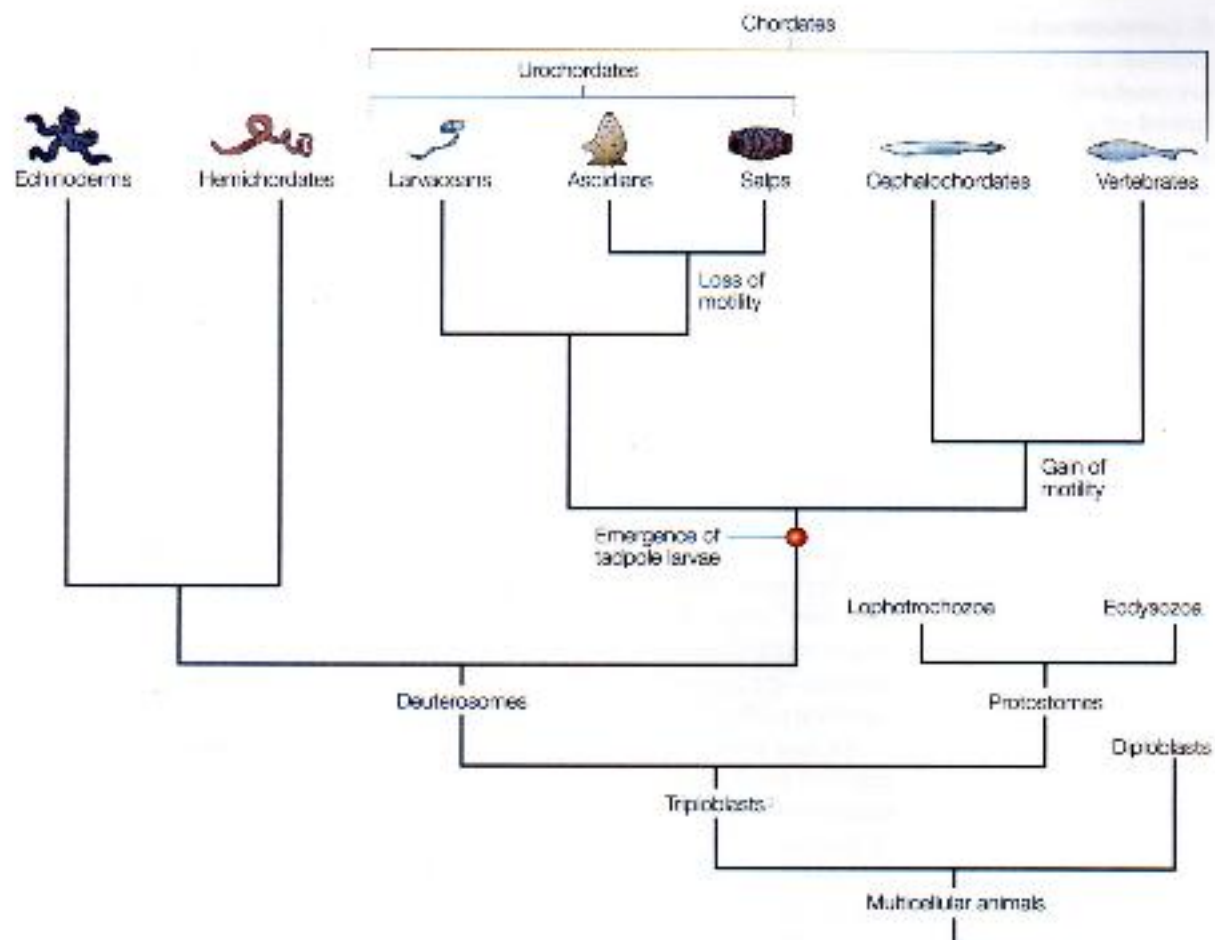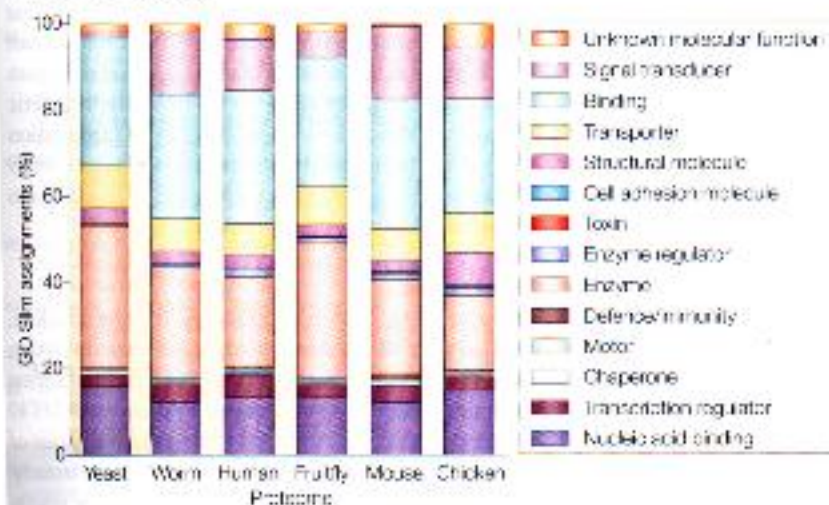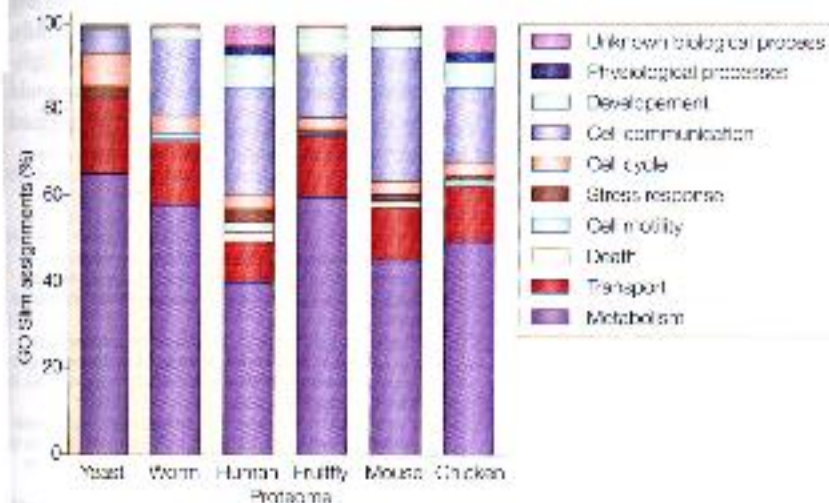
Figure 2 | **The evolution of the chordates.** The chordates comprise the urochordates, the cephalochordates and the vertebrates. They are thought to have evolved from a common ancestor shared with the non-chordate deuterostomes (the echinoderms and the hemichordates). The emergence of tadpole-type larvae was a key event in the evolution of the chordates.

**a** Molecular function

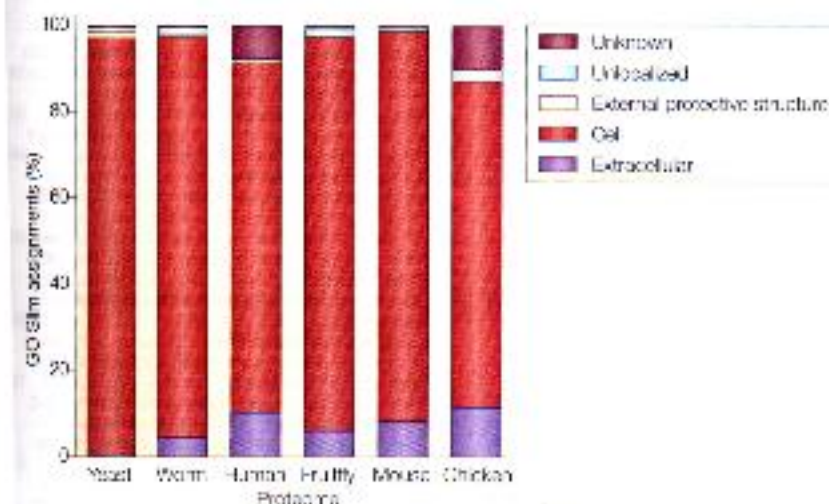**b** Biological processes

**c** Cellular component

Figure 2 | **Gene Ontology functional assignments to eukaryotic proteomes that have been completely sequenced.** The relative fractions of assigned protein functions using the Gene Ontology (GO) Slim classification for *Saccharomyces cerevisiae* (yeast), *Caenorhabditis elegans* (worm), *Homo sapiens* (human), *Drosophila melanogaster* (fruitfly) and *Mus musculus* (mouse) taken from the European Bioinformatics Institute (EBI) proteome web site (see online links) compared with a GO assignment on a known assembly of 350,000 chicken ESTs (see online link to BBSRC Chicken EST Project) are shown.
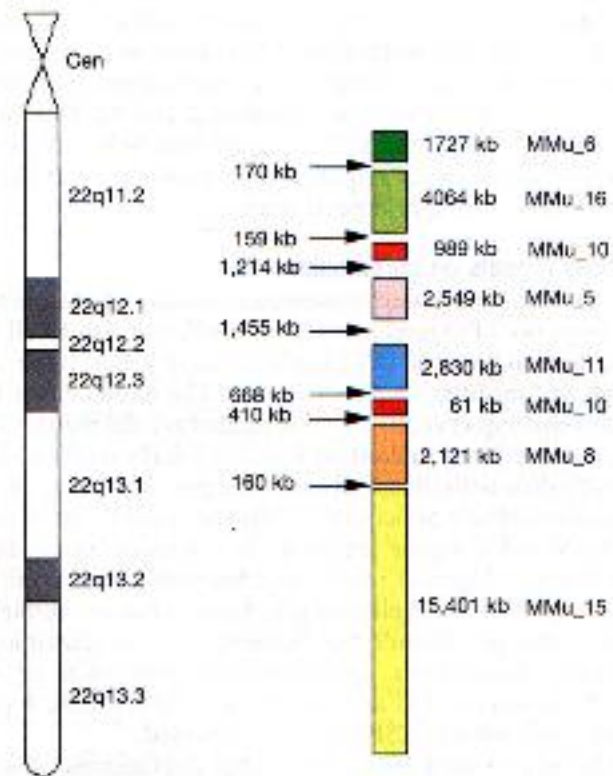
**Figure 5** Regions of conserved synteny between human chromosome 22 and the mouse genome. Regions of mouse chromosomes with conserved synteny to human chromosome 22 are shown as adjacent coloured blocks, determined by the mouse map position of mouse orthologues to human chromosome 22 genes. The size of human chromosome 22 corresponding to each mouse chromosomal region is indicated in kb, as well as the size of the gap between the last orthologue in each conserved block. These data are available at http://www.sanger.ac.uk/Chr22/Mouse.
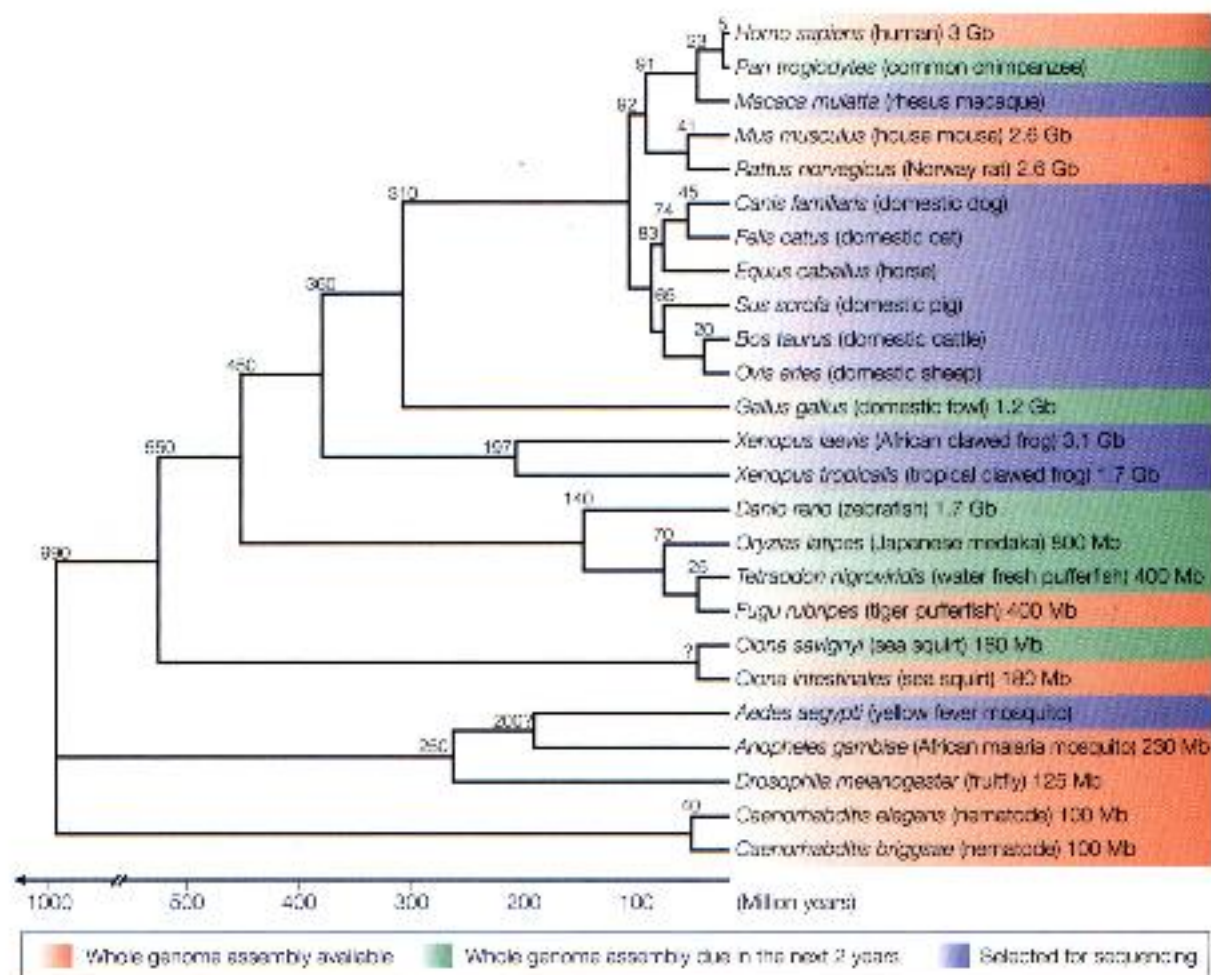
Figure 1 | **Evolutionary relationship between metazoans that are sequenced or due for sequencing.** The simplified phylogenetic relationships between the metazoans for which the complete, or nearly complete, genome sequences are available or will be available soon. Evolutionary distances (in million years) and genome sizes are based on REFS 10, 13, 58–100.