# Supporting Information

## Markov et al. 10.1073/pnas.0812138106

**SI Text**

**Searching Strategy.** Orthologs from the proteins of interest (i.e., all mentioned vertebrate, arthropod, and nematode steroidogenic enzymes) were searched by blasting again the mentioned databases. The phylogenetic position of the organisms that were screened in this study is indicated in Fig. S2, with information on the data quality. All sequence hits were retrieved, and the dataset was cleaned using the following criteria: sequences lacking a conserved family motif (e.g., SDR cofactor binding site TG***G*G) were discarded, and truncated sequences were also discarded when there were found to be members of paralog groups. The majority of sequences from *Schmidtea mediterranea*, *Trichinella spiralis*, and *Aplysia californica* were eliminated during this step, because ab initio predictions were shown to be less accurate than the EST-based predictions that were available for other species. Sequences were checked by eye in SEAVIEW (2) to eliminate too divergent positions and unaligned regions before phylogenetic reconstruction.

**Protein Sequences.** Additionally, as *SI Text* we provide expanded trees corresponding to Figs. 2 and 3. For each tree, all sequence accession numbers are grouped in Dataset S1, (1 sheet per tree). The given accession numbers are GenBank IDs, jgi IDs (only numbers, mainly for *Capitella*, *Nematostella*, *Helobdella*, *Lottia*, *Trichoplax, Daphnia*, and some sequences from *Branchiostoma floridae*), or Ensembl IDs (sequences beginning with *ENS0000*). Two additional sequences are not in those databases:

- Pr5 from *Aplysia californica* is a homemade GENSCAN (1) prediction from the GenBank contig AASC01065054.1.
- 002 from *Amphimedon queenslandica*, which is a manually annotated prediction based upon traces, that is available on the online website from David Nelson:

http://drnelson.utmem.edu/biblioC.html.

Proteins that are marked with * are those for which corrected intron-exon boundaries were manually performed by D. Nelson and that are available on the indicated Web page.

**Detection of Annotation Errors.** Illustrating the difficulties in assessing orthology when partial data sets are used, the recently cloned *Branchiostoma belcheri* protein BAF61103.1, originally described as CYP11 (3) is in fact a member of CYP374, a distant paralog group of deuterostome CYPs, which was lost in vertebrates (Fig. S4). This shows that experimental data concerning the enzymatic activities of the CYPs can be biased by a wrong identification linked to a partial phylogenetic analysis (4). Similarly, the *Branchiostoma belcheri* BAF61104.1, that is described in the same paper as a CYP17 is clearly not an ortholog of the vertebrate and *Branchiostoma floridae* CYP17s, but a paralog from a subfamily where the gene may have been lost in vertebrates too.

**A Nomenclature Note About Fig. 4.** We propose the name "ecdysosteroid" to name steroids from ecdysozoans because the classically used name "ecdysteroids" is used to describe steroids from arthropods and steroids from plants that have the same structure.

1. Burge C, Karlin S (1997) Prediction of complete gene structures in human genomic DNA. *J Mol Bio* 268:78–94.
2. Galtier N, Gouy M, Gautier C (1996) SEAVIEW and PHYLO_WIN: Two graphic tools for sequence alignment and molecular phylogeny. *Comput Appl Biosci* 12:543–548.
3. Mizuta T, Kubokawa K (2007) Presence of sex steroids and cytochrome P450 (CYP) genes in amphioxus. *Endocrinology* 148:3554–3565.
4. Markov G, Lecointre G, Demeneix B, Laudet V (2008) The ''street light syndrome'', or how protein taxonomy can bias experimental manipulations. *Bioessays* 30:349–357.

Numbering of the carbons in steroid skeleton

Cholesterol

Aldosterone

Cortisol

Estradiol

Dehydrotestosterone
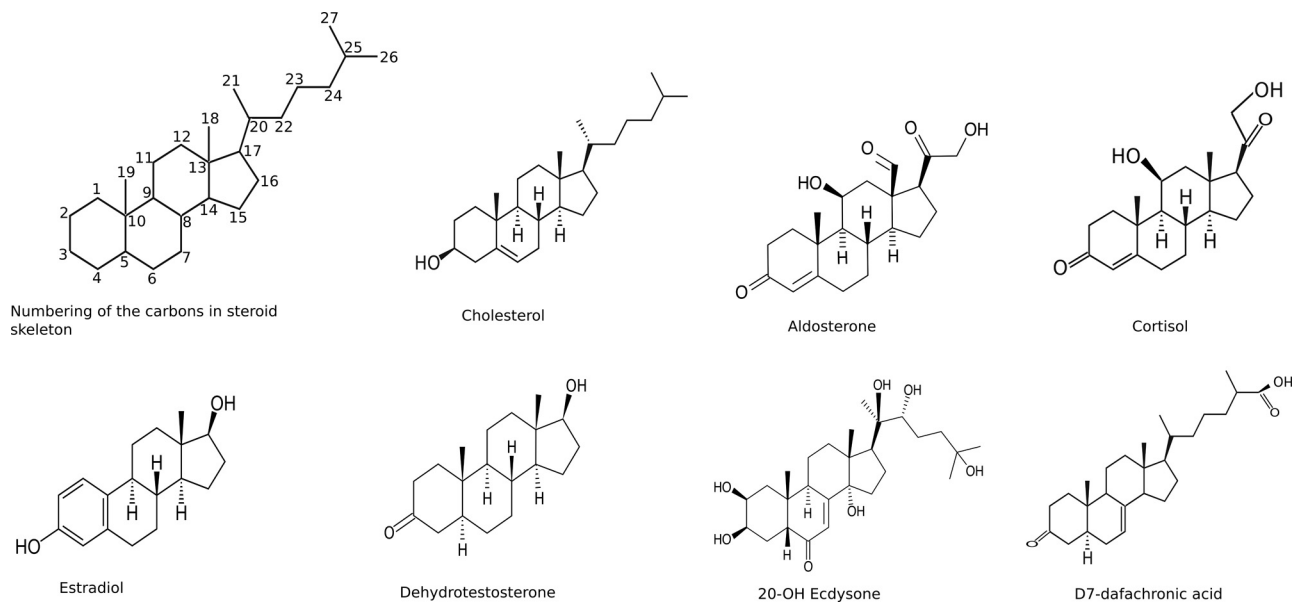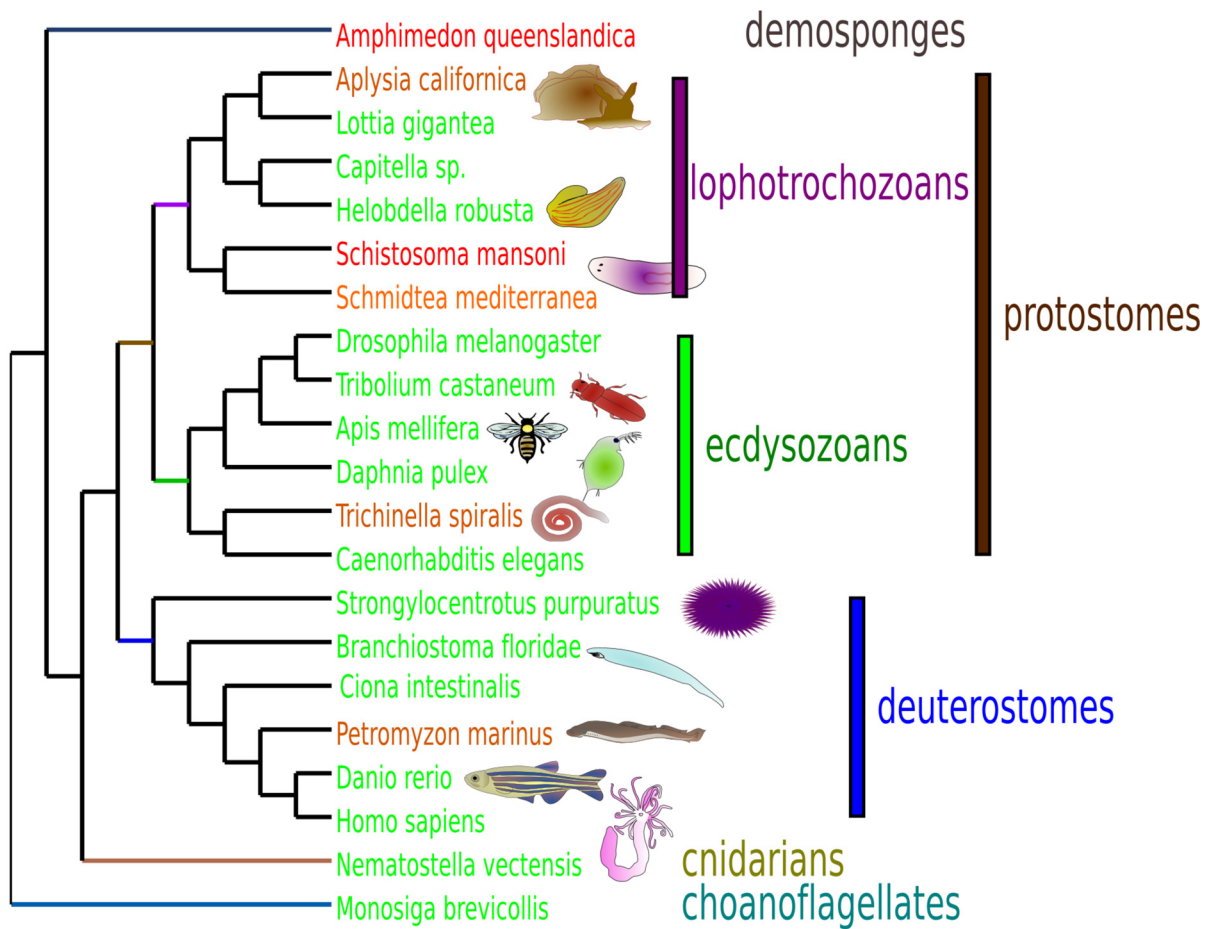
20-OH Ecdysone

D7-dafachronic acid

**Fig. S1.** Structure of the sterol ring, of cholesterol, and of the steroid hormones of human, *Drosophila melanogaster* and *Caenorhabditis elegans*. The sterol ring numbering is indicated. Aldosterone, cortisone, estradiol and dihydrotestosterone are human steroid hormones. 20-OH ecdysone is the main steroid hormone in *Drosophila melanogaster*, whereas delta-7-dafachronic acid is 1 of the 2 steroids in *Caenorhabditis elegans*.

**Fig. S2.** Genomic data used in this study. The genomic model species that were screened in this study are indicated, with complementary information about their phylogenetic relationships and about the quality of their genome data. Species in green are those for which EST-based gene predictions are available. The genome of species in orange is provided as contigs, that were used for ab initio predictions. The genome of species in red is available only as traces.
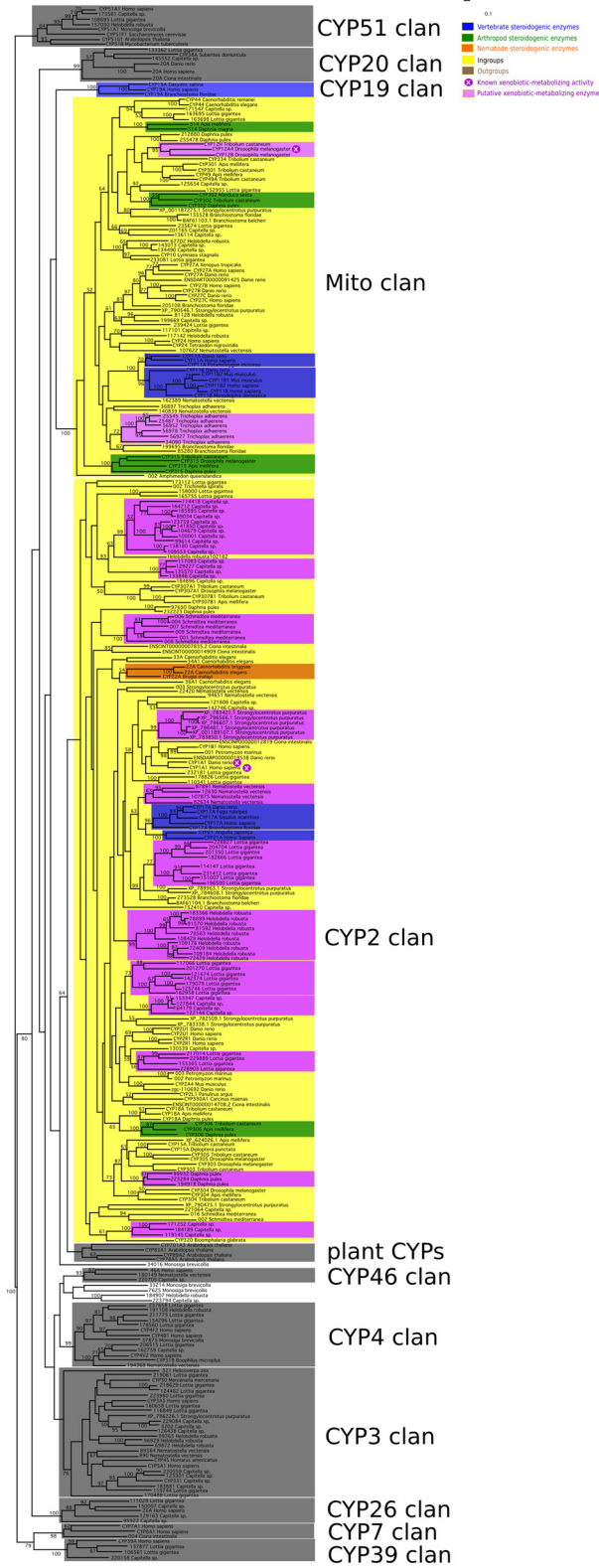
**Fig. S3.** Phylogeny of the CYP family. A maximum-likelihood analysis of the CYP family. Vertebrate steroidogenic proteins are highlighted in blue, arthropod steroidogenic proteins are in green and nematode steroidogenic proteins are in orange. Enzymes with known xenobiotic-metabolizing activity are indicated by circled ''X'', and proteins resulting from abundant lineage-specific duplication, that are thus candidate xenobiotic-metabolising enzymes, are highlighted by purple boxes. For details about the mito clan, see also Fig. S4.
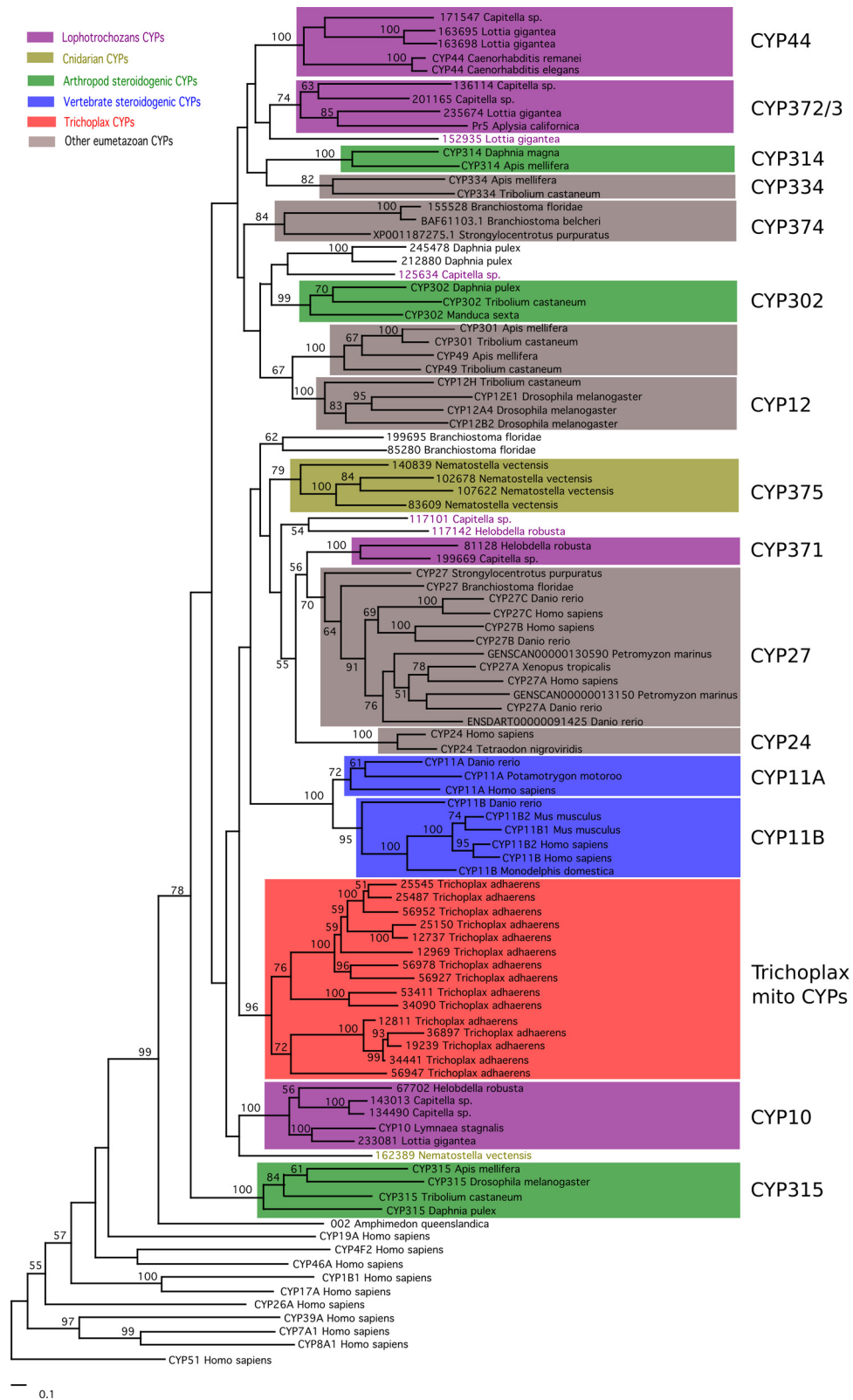
**Fig. S4.** Phylogeny of the mitochondrial CYP clan. A maximum-likelihood analysis of the mitochondrial CYP clan. Proteins are named according to classical CYP nomenclature when an official name exists. At least 3 groups of paralogs with unknown activity were found in lophotrochozoans (CYP10, CYP372/CYP373, and CYP371). Vertebrate steroidogenic CYPs are highlighted in blue, arthropod steroidogenic CYPs are in green, lophotrochozoan mito CYPs are in purple, cnidarian mito CYPs are in light brown, Trichoplax mito CYPs are in red, and other supported mito CYP clades are in gray.
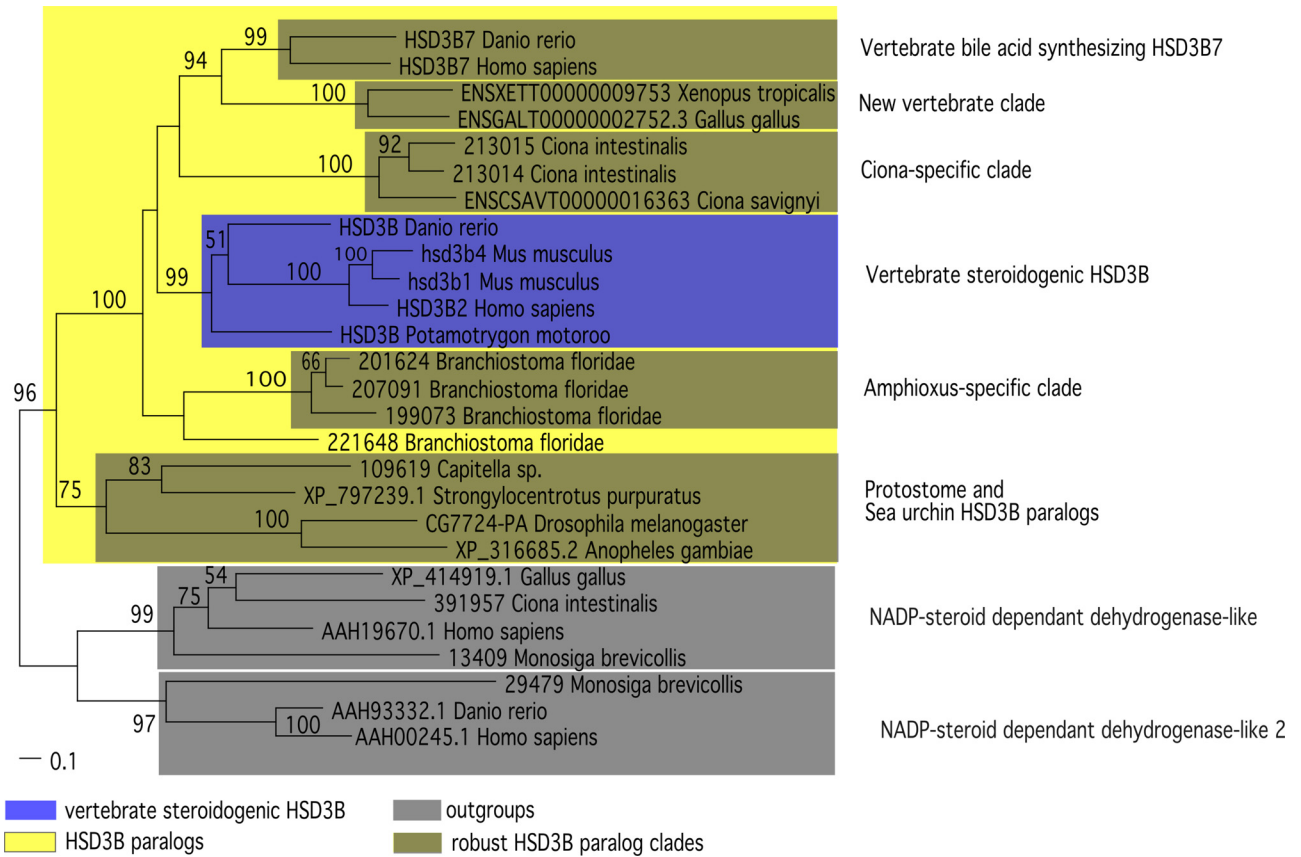
**Fig. S5.** Phylogeny of the SDR family. Our analysis reveals the existence of 23 strongly supported families that were arbitrarily named from 1 to 23, and are highlighted in gray. Groups of vertebrate steroidogenic enzymes are in blue, and other member of the same subfamily are in yellow. Enzymes with known retinoid-metabolizing activities are indicated by circled ''R.''

**Fig. S6.** Phylogeny of the HSD3B family. A maximum-likelihood analysis of the HSD3B family. Groups of vertebrate steroidogenic enzymes are in blue, other members of the same subfamily are in yellow. Outgroups are in gray. Robust HSD3B paralogs clades (those who are indicated by red dots in Fig. 2) are highlighted in yellow+gray.
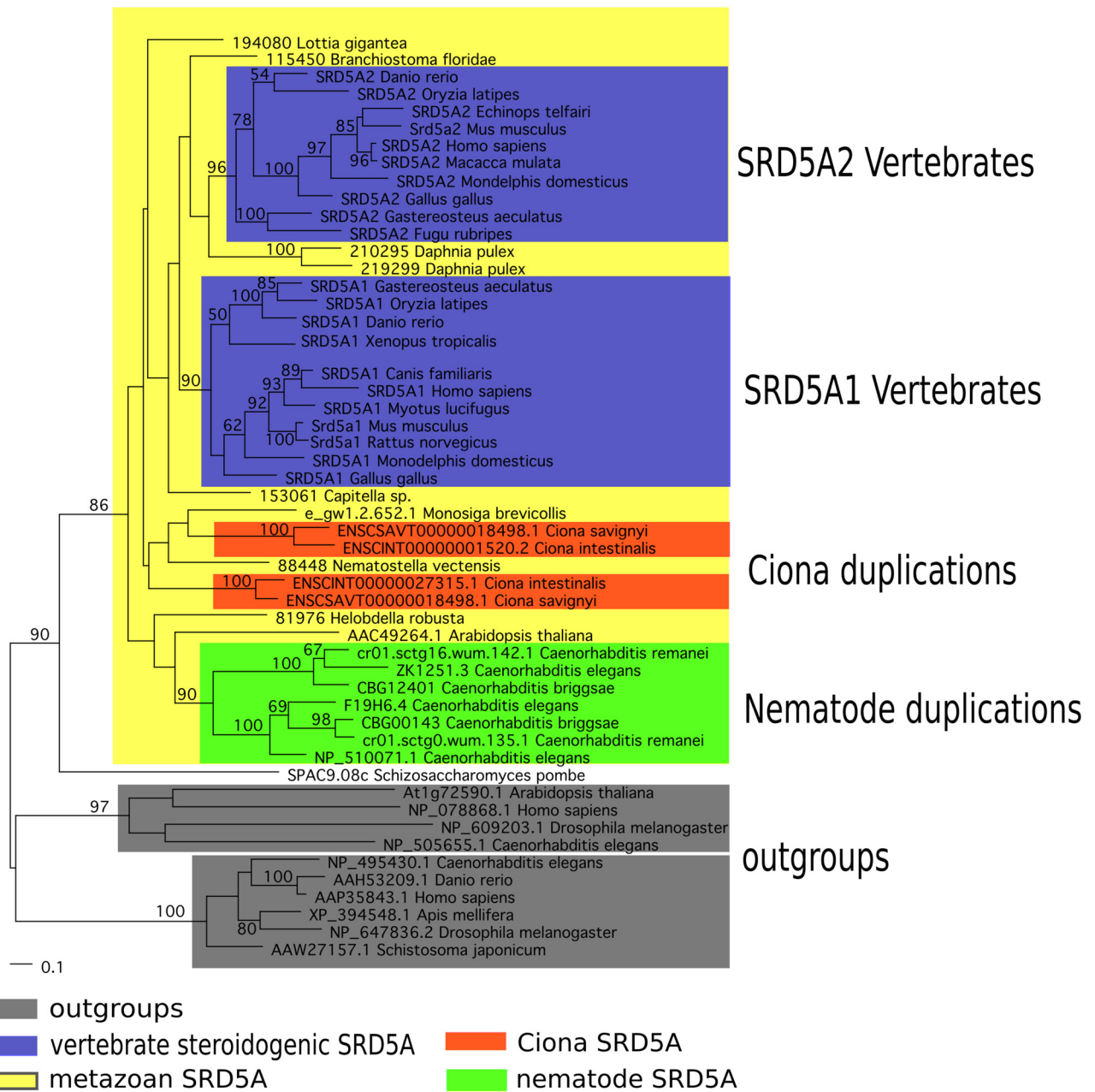
**Fig. S7.** Phylogeny of the SRD5A family. A maximum-likelihood analysis of the SRD5A family. Vertebrate SRD5A are highlighted in blue, nematode duplicated SRD5 are highlighted in green and ciona duplications are highlighted in orange. The metazoan SDR5A family is in yellow and the outgroups in gray.

## Other Supporting Information Files

Dataset S1 (XLS)