

SUPPLEMENTAL MATERIAL

**Increased taxon sampling reveals thousands of hidden orthologs in
flatworms**

José M. Martín-Durán, Joseph F. Ryan, Bruno C. Vellutini, Kevin Pang, Andreas Hejnol

TABLE OF CONTENTS

- Supplemental Methods
- Supplemental Figures S1–S9
- Supplemental Tables 1–10
- Supplemental References
- Supplemental File 1 (scripts and Leapfrog code)
- Supplemental File 2 (HumRef2015 dataset and protein alignments)

Supplemental Methods

Macrostomum lignano transcriptome

Adult and juveniles of *M. lignano* were kept under laboratory conditions as described elsewhere (Rieger et al. 1988). Animals starved for four days were homogenized and used as source material to isolate total RNA with the TRI Reagent (Life Technologies) following the manufacturer's recommendations. A total of 1 µg was used for Illumina paired-end library preparation at the GeneCore facilities (EMBL) and sequencing in a HiSeq 2000 platform. Paired-end reads were cleaned for adaptors using Trimmomatic v.0.35 (Bolger et al. 2014) and resulting reads were assembled *de novo* with Trinity v.r20140717 using default settings (Grabherr et al. 2011).

Data set preparation

We downloaded the Human RefSeq FASTA file from the NCBI FTP site last updated on March 25, 2015

(ftp://ftp.ncbi.nlm.nih.gov/refseq/H_sapiens/H_sapiens/protein/protein.fa.gz). We also downloaded the gene2accession data file from NCBI, which was last updated on July 3, 2015 (<ftp://ftp.ncbi.nlm.nih.gov/gene/DATA/gene2accession.gz>). We then used the reduce_refseq script (Supplemental File 1) to generate a non-redundant Human RefSeq FASTA file with the following command: (reduce_refseq --fasta=protein.fa.gz --gene2accession=gene2accession.gz > HumRef2015.fa). This script prints only the first isoform for each Gene ID in the RefSeq FASTA file. The resulting file (available in Supplemental File 2) will be hereafter referred to as HumRef2015. Additionally, we assembled *de novo* a transcriptome for *M. lignano* (Supplemental Methods), and downloaded the 28 RNA-seq *de novo* assemblies from (Laumer et al. 2015) and six additional *S. mediterranea* datasets from PlanMine v1.0 (Brandl et al. 2016) on May 29,

2015. On July 14, 2015 we downloaded *Schistosoma mansoni*, *Hymenolepis microstoma*, and *Girardia tigrina* gene models from the Sanger FTP sites (*S. mansoni*: <ftp://ftp.sanger.ac.uk/pub/project/pathogens/Schistosoma/mansoni/>; *H. microstoma*: <ftp://ftp.sanger.ac.uk/pub/project/pathogens/Hymenolepis/microstoma/>; *G. tigrina*: <https://drive.google.com/open?id=0B-oXGfLmcaXQUUhVS29IbC1CeHM>). Further details on datasets are available in Supplemental Table 1. All flatworm transcriptomes were constructed from adult tissue, except for *S. ellipticus*, *P. vittatus*, *G. applanata* and *Prorhynchus* sp. I, which included embryonic stages.

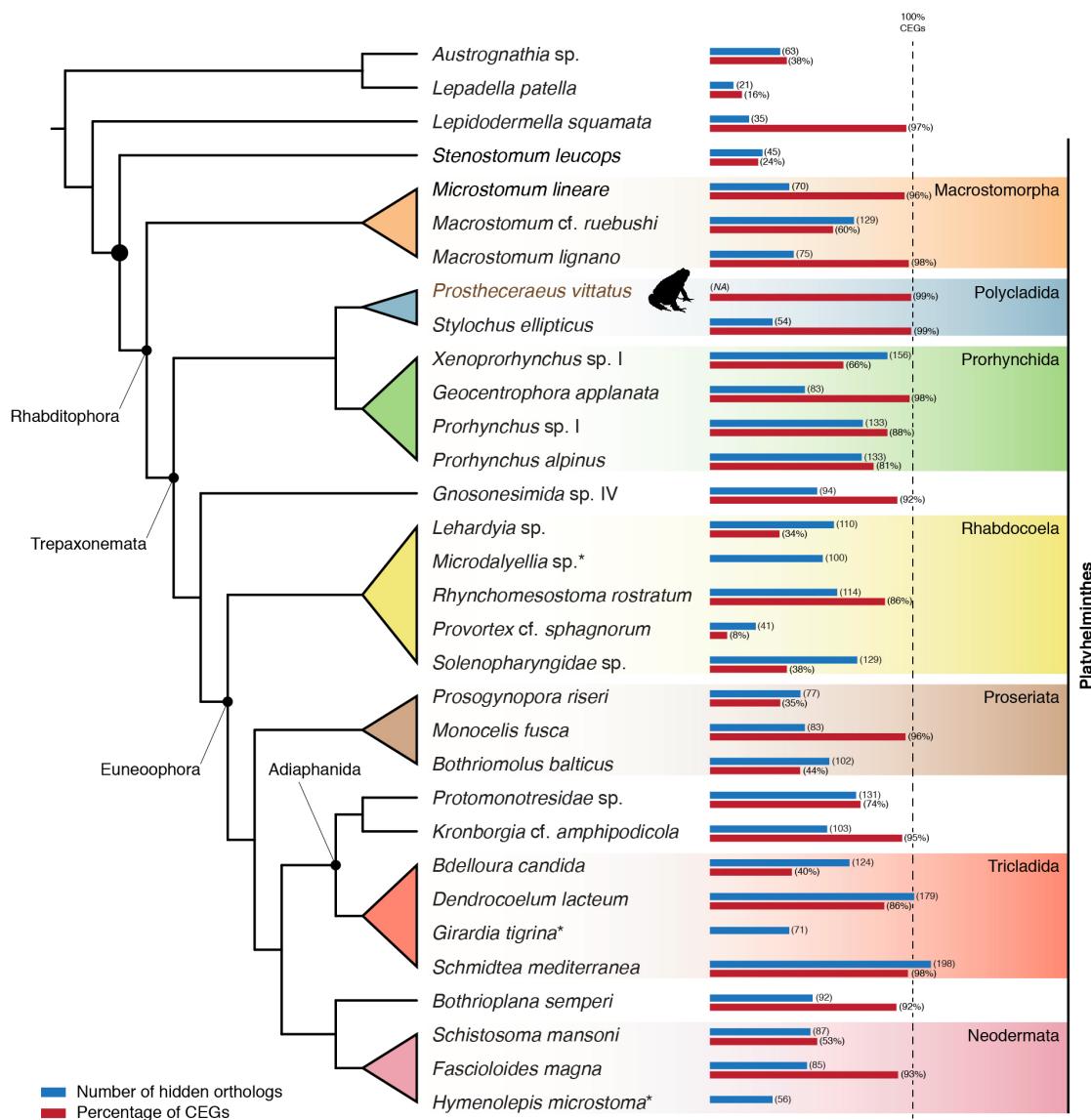
CEGMA analysis, transcriptome quality assessment, and statistics

Transcriptome completeness was evaluated with CEGMA (Parra et al. 2007; Parra et al. 2009). We could not run the CEGMA pipeline in the transcriptomes of *G. tigrina*, *Microdalyellia* sp. and *H. microstoma* due to an untraceable error. We calculated the contig metrics for each transcriptome assembly with TransRate (Smith-Unna et al. 2015). Principal component analysis was performed in R (R Core Team 2015) and plotted using the ggplot2 package (Wickham 2009). For this analysis, branch lengths were the root-to-tip distances inferred from a maximum likelihood phylogeny previously reported (Laumer et al. 2015).

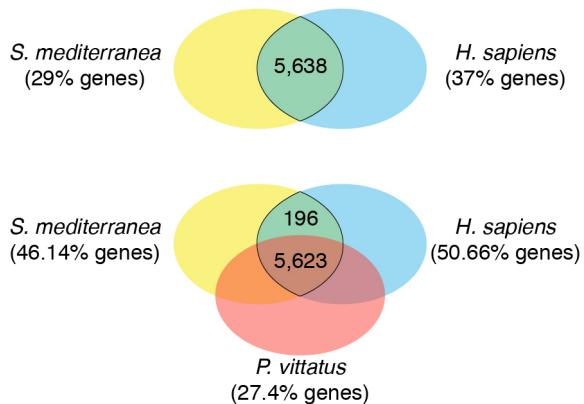
Smed-SDCCAG8 and Smed-CEP192 primer sequences

The following gene specific primers were used to clone *Smed-SDCCAG8* and *Smed-CEP192*: *Smed-SDCCAG8-F* 5' TTGACCAGGAACCGTACGAA 3', *Smed-SDCCAG8-R* 5' CACATTGCTCGAATCTGGCA 3'; *Smed-CEP192-F* 5' CAACAGGTCCGATTCCAACC 3', *Smed-CEP192-R* 5' AACGCAACAGGAACCAGAAC 3'.

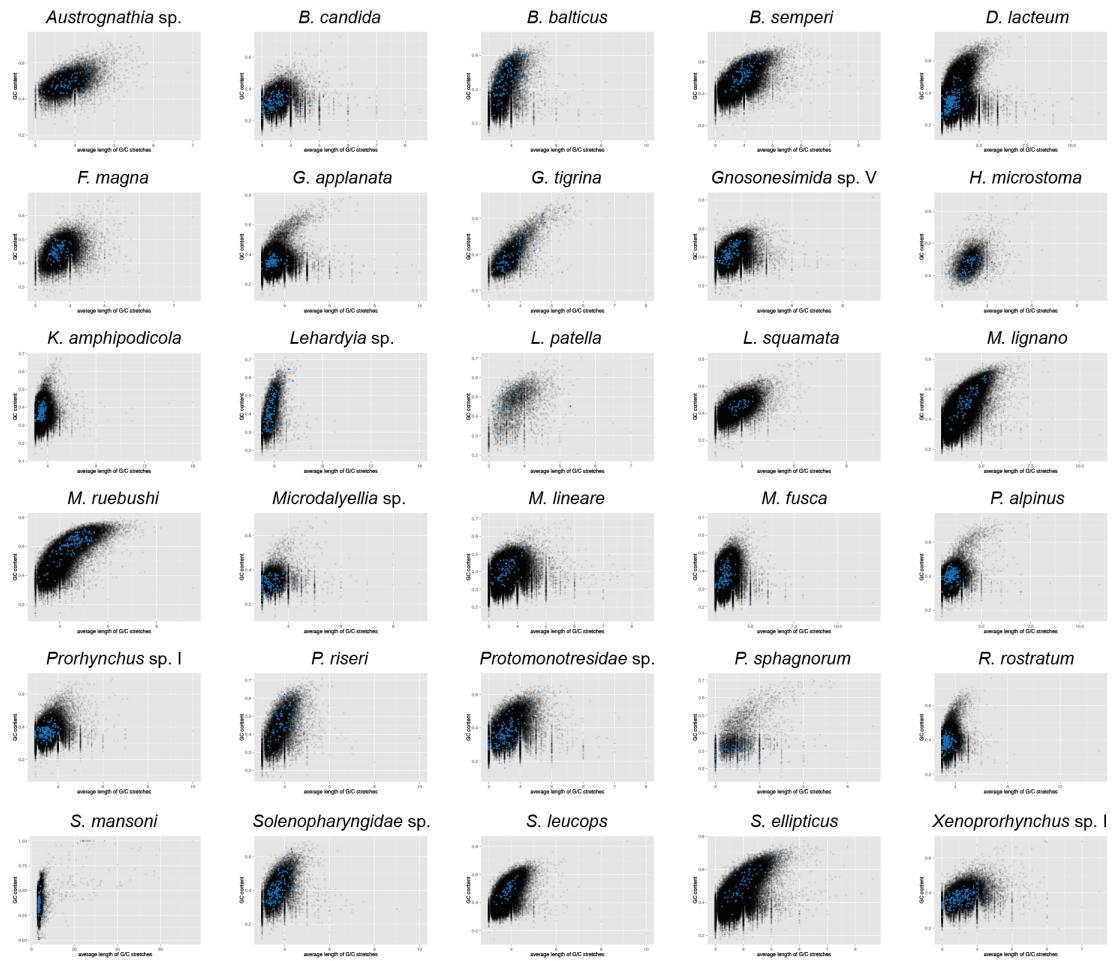
Supplemental Figures



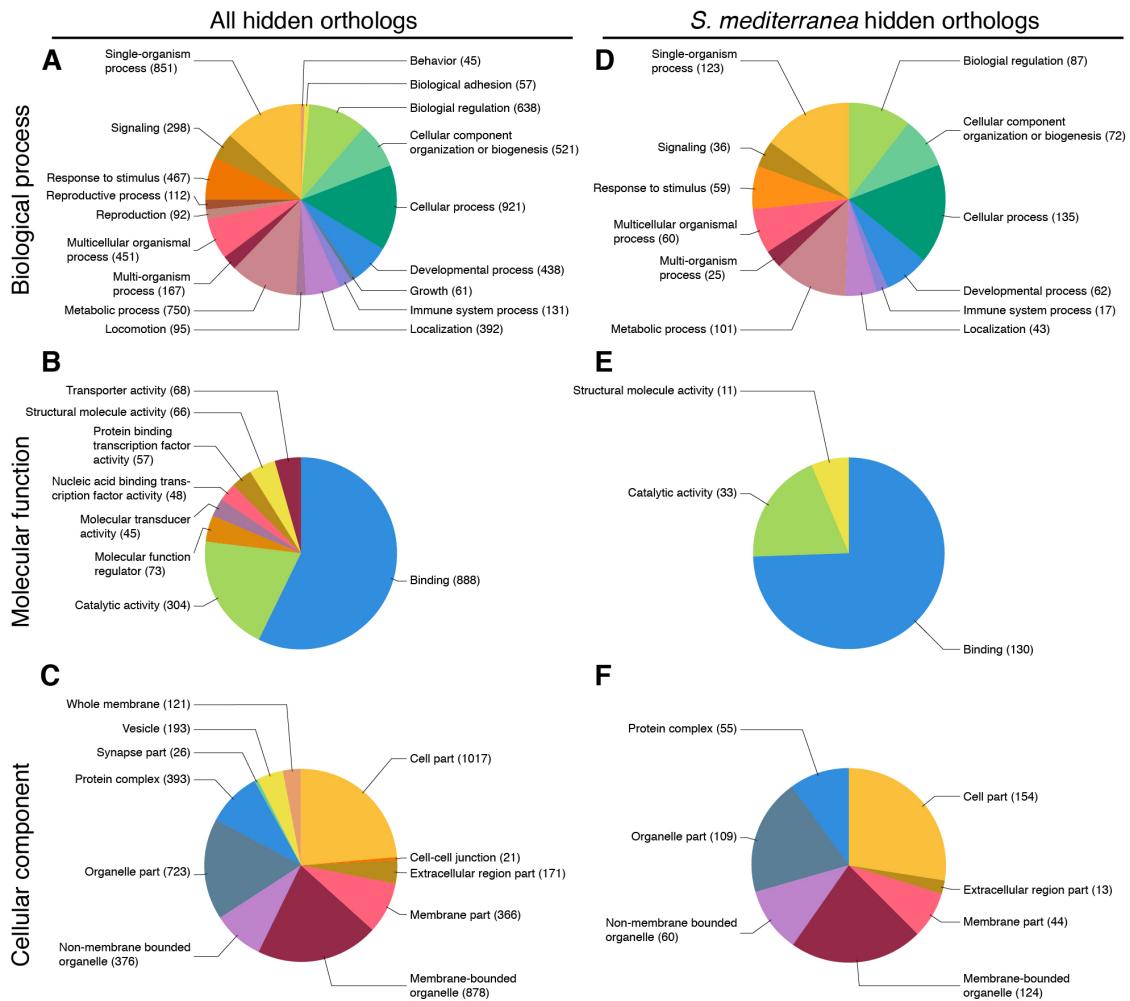
Supplemental Figure S1. Distribution of hidden orthologs in the analyzed flatworm transcriptomes. The figure shows the total number of hidden orthologs in the analyzed transcriptomes in a phylogenetic context and with respect to their completeness (percentage of recovered core eukaryote genes, CEGs). The quality of the transcriptomes seems to be a limitation for the recovery of hidden orthologs in some flatworm lineages (e.g. *Provortex cf. sphagnorum*). However, the number of hidden orthologs is very species-specific.



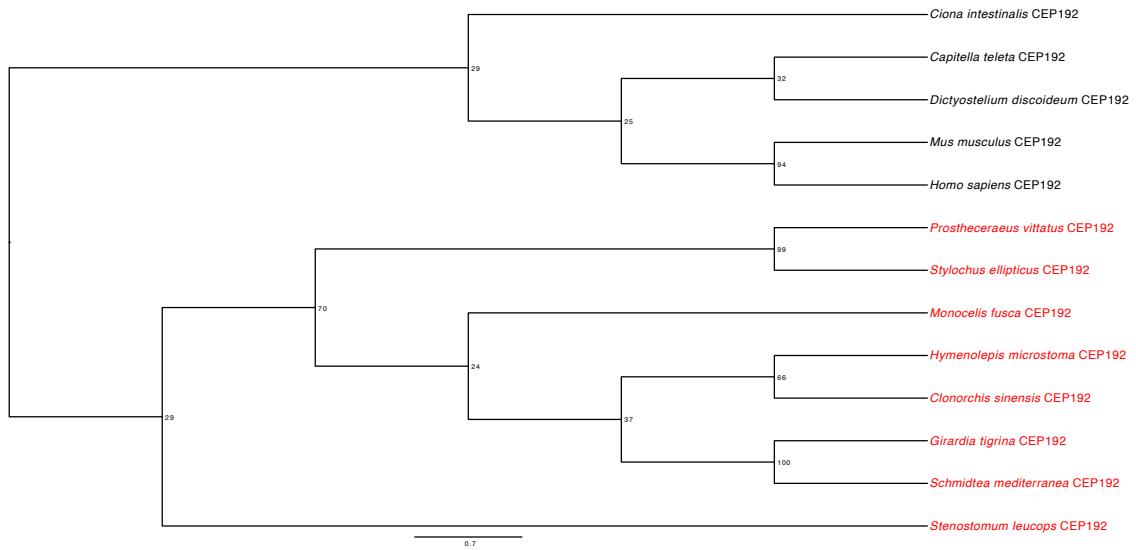
Supplemental Figure S2. Impact of a ‘bridge’ species in OrthoFinder. Number of orthogroups identified by OrthoFinder when only *S. mediterranea* and human are considered (top) and when a ‘bridge’ transcriptome is added (bottom). The inclusion of an intermediate species increases the number of sequences assigned to orthogroups from 29% to 46.14% of the transcriptome, and the number of orthogroups between *S. mediterranea* and human from 5,638 to 5,819. The inclusion of the *P. vittatus* ‘bridge’ transcriptome in an OrthoFinder analysis recovered only 62 out of the 75 hidden orthologs (82.7%) identified by the Leapfrog pipeline in the same *S. mediterranea* transcriptome with similar E-value cutoffs. For *S. mediterranea*, we used the transcriptome sequenced by J. Rink’s lab at the MPI in Dresden (Brandl et al. 2016).



Supplemental Figure S3. GC content in flatworm transcriptomes. GC content of each transcript plotted against its average length of G/C stretches for each flatworm species under study. The transcripts corresponding to hidden orthologs are in blue. Hidden orthologs do not differentiate from the majority of transcripts.



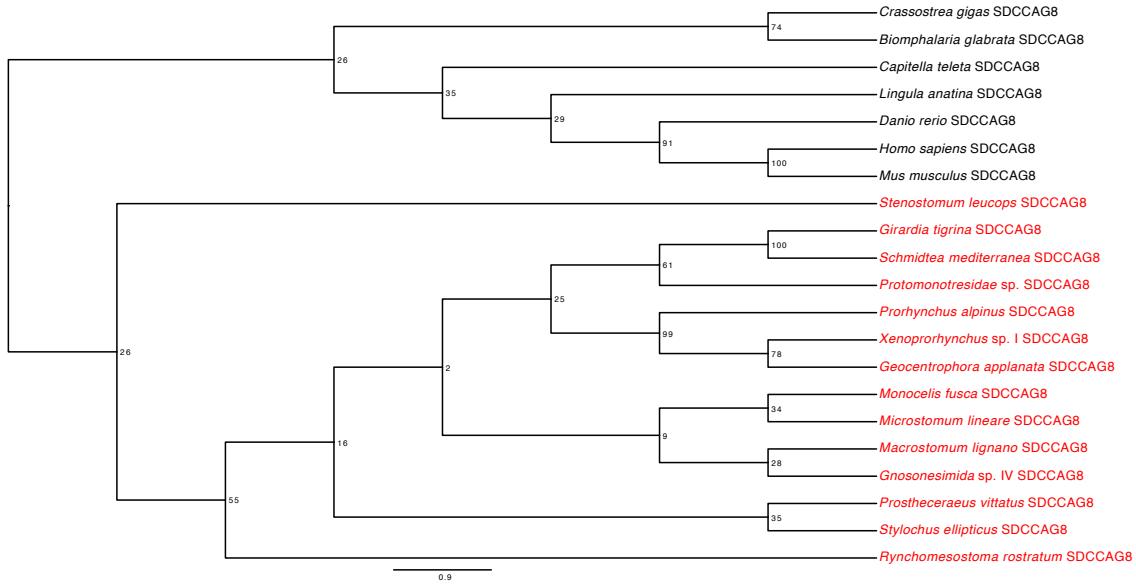
Supplemental Figure S4. Gene Ontology (GO) characterization of hidden orthologs. Distribution of GO terms for all recovered hidden orthologs (A–C) and for the hidden orthologs identified in *S. mediterranea* (D–F). Hidden orthologs include a great diversity of GO categories, with a big proportion of binding and catalytic activity. The number of GO nodes in each category is indicated in parentheses.



Supplemental Figure S5. Orthology analysis of the centrosomal CEP192 protein.

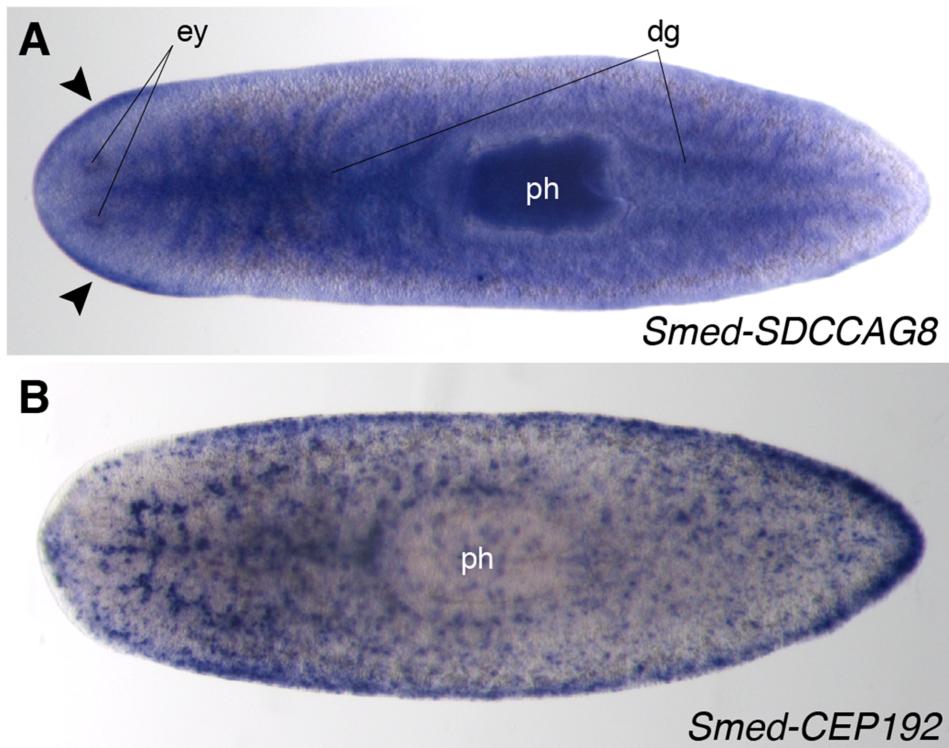
CEP192 proteins do not contain any identifiable protein domain, and there is no known related protein that can help root the tree. Flatworm sequences are highlighted in red.

Model of protein evolution: RtRev+I+G+F.

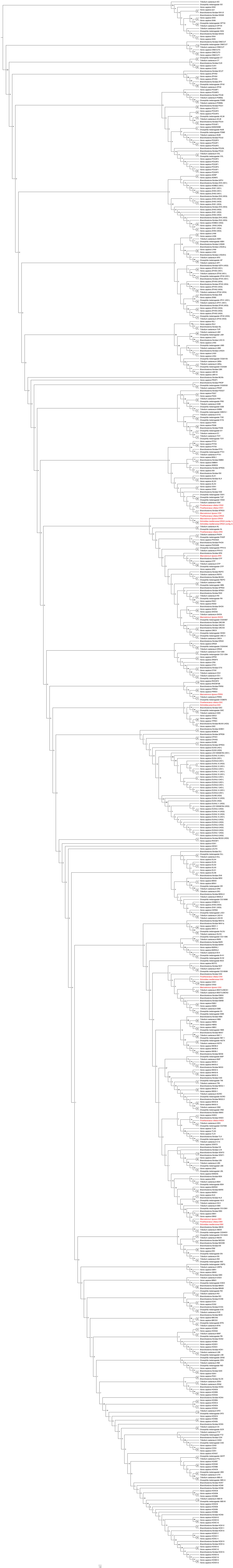


Supplemental Figure S6. Orthology analysis of the centrosomal SDCCAG8

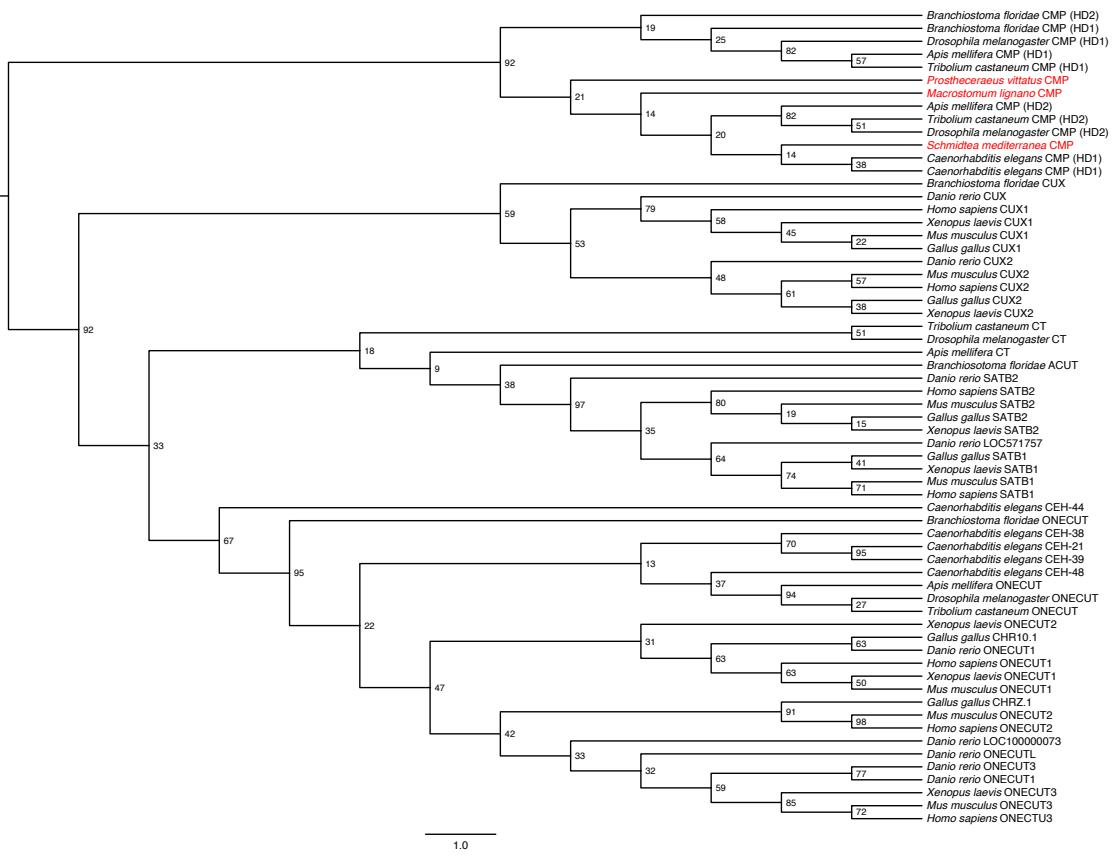
protein. SDCCAG8 proteins contain a SDCCAG8 domain (PFAM: PF15964), which is exclusive of these proteins. The domain is clearly recognizable in all flatworm sequences except *P. alpinus* (fragment too short) and the triclad *G. tigrina* and *S. mediterranea* (too divergent). Flatworm sequences are highlighted in red. Model of protein evolution: JTT+G+F.



Supplemental Figure S7. Expression of *Smed- SDCCAG8* and *Smed-CEP192* in adult planarians. (A–B) Whole mount *in situ* hybridization of *Smed-SDCCAG8* (n=6) and *Smed-CEP192* (n=6) in adult intact *S. mediterranea*. (A) *Smed-SDCCAG8* is expressed in the epidermis, anterior margin (black arrowheads), in the digestive system (dg) and pharynx, and faintly in the parenchyma. (B) *Smed-CEP192* is expressed in a salt-and-pepper manner in isolated cells of the epidermis, posterior margin of the animal, and parenchyma.



Supplemental Figure S8. Orthology analysis of the ANTP homeodomain class. The newly identify sequences in the macrostomid *M. lignano*, the polyclad *P. vittatus* and the triclad *S. mediterranea* are highlighted in red. Model of protein evolution: LG+G.



Supplemental Figure S9. Orthology analysis of the CUT homeodomain class. The newly identify sequences in the macrostomid *M. lignano*, the polyclad *P. vittatus* and the triclad *S. mediterranea* are highlighted in red. Model of protein evolution: LG+G.

Supplemental Tables

Supplemental Table 1. Transcriptomes analyzed in this study.

Supplemental Table 2. Recovered hidden orthologs. Hidden orthologs (as in HumRef2015) recovered in each transcriptome after running Leapfrog with the transcriptome of the polyclad *P. vittatus* used as the ‘bridge’.

Supplemental Table 3. List of hidden orthologs with reciprocal best BLAST hit against multiple human proteins.

Supplemental Table 4. PFAM domain composition of 130 human queries and their ortholog ‘bridge’ proteins in *P. vittatus*.

Supplemental Table 5. Number of hidden orthologs recovered with iterated ‘bridge’ transcriptomes.

Supplemental Table 6. Data set used for principal component analysis.

Supplemental Table 7. Length of hidden orthologs and ORFs in flatworm transcriptomes.

Supplemental Table 8. PFAM domains identified in the hidden orthologs.

Supplemental Table 9. Significantly enriched GO terms in the hidden orthologs recovered in *S. mediterranea*.

Supplemental References

- Bolger AM, Lohse M, Usadel B. 2014. Trimmomatic: a flexible trimmer for Illumina sequence data. *Bioinformatics* **30**: 2114-2120.
- Brandl H, Moon H, Vila-Farre M, Liu SY, Henry I, Rink JC. 2016. PlanMine - a mineable resource of planarian biology and biodiversity. *Nucleic Acids Res* **44**: D764-773.
- Grabherr M, Haas BJ, Yassour M, Levin J, Thompson D, Amit I, Adiconis X, Fan L, Raychowdhury R, Zeng Q et al. 2011. Full-length transcriptome assembly from RNA-seq data without a reference genome. *Nat Biotechnol* **29**: 644-652.
- Laumer CE, Hejnol A, Giribet G. 2015. Nuclear genomic signals of the 'microturbellarian' roots of platyhelminth evolutionary innovation. *Elife* **4**.
- Parra G, Bradnam K, Korf I. 2007. CEGMA: a pipeline to accurately annotate core genes in eukaryotic genomes. *Bioinformatics* **23**: 1061-1067.
- Parra G, Bradnam K, Ning Z, Keane T, Korf I. 2009. Assessing the gene space in draft genomes. *Nucleic Acids Res* **37**: 289-297.
- R Core Team. 2015. R: A Language and Environment for Statistical Computing. R Foundation for Statistical Computing, Vienna, Austria.
- Rieger R, Gehlen M, Haszprunar G, Holmlund M, Legniti A, Salvenmoser W, Tyler S. 1988. Laboratory cultures of marine Macrostomida (Turbellaria). *Progr Zool* **36**: 523.
- Smith-Unna RD, Boursnell C, Patro R, Hibberd JM, Kelly S. 2015. TransRate: reference free quality assessment of de-novo transcriptome assemblies. *BioRxiv* **021626**.
- Wickham H. 2009. *ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag, New York.