

1 **Calcisponges have a ParaHox gene and dynamic expression of**
2 **dispersed NK homeobox genes**

3

4 Sofia A.V. Fortunato^{1,2}, Marcin Adamski¹, Olivia Mendivil Ramos^{3,†}, Sven Leininger^{1,◇},

5 Jing Liu¹, David E.K. Ferrier³ and Maja Adamska^{1§}

6 ¹Sars International Centre for Marine Molecular Biology, University of Bergen,

7 Thormøhlensgate 55, 5008 Bergen, Norway

8 ²Department of Biology, University of Bergen, Thormøhlensgate 55, 5008 Bergen,

9 Norway

10 ³ The Scottish Oceans Institute, Gatty Marine Laboratory, School of Biology, University of

11 St Andrews, East Sands, St Andrews, Fife KY16 8LB, UK.

12 † Current address: Stanley Institute for Cognitive Genomics, Cold Spring Harbor

13 Laboratory, 1 Bungtown Road, Cold Spring Harbor, NY 11724, USA.

14 ◇ Current address: Institute of Marine Research, Nordnesgaten 50, 5005 Bergen, Norway

15 §Corresponding author

16 Email address: maja.adamska@sars.uib.no

17

18 **Summary**

19 **Sponges are simple animals with few cell types, but their genomes paradoxically**
20 **contain a wide variety of developmental transcription factors¹⁻⁴, including**
21 **homeobox genes belonging to the *Antennapedia* (ANTP) class^{5,6}, which in**
22 **bilaterians encompass *Hox*, *ParaHox* and *NK* genes. In the genome of the**
23 **demosponge *Amphimedon queenslandica*, no *Hox* or *ParaHox* genes are present,**
24 **but *NK* genes are linked in a tight cluster similar to the *NK* clusters of bilaterians⁵.**
25 **It has been proposed that *Hox* and *ParaHox* genes originated from *NK* cluster**
26 **genes after divergence of sponges from the lineage leading to cnidarians and**
27 **bilaterians^{5,7}. On the other hand, synteny analysis gives support to the notion that**
28 **absence of *Hox* and *ParaHox* genes in *Amphimedon* is a result of secondary loss**
29 **(the ghost locus hypothesis)⁸. In this study, we analyzed complete suites of *Antp***
30 **class homeoboxes in two calcareous sponges, *Sycon ciliatum* and *Leucosolenia***
31 ***complicata*. Our phylogenetic analyses demonstrate that these calcisponges**
32 **possess orthologues of bilaterian *NK* genes (*Hex*, *Hmx* and *Msx*), a varying number**
33 **of additional *NK* genes and one *ParaHox* gene, *Cdx*. Despite generation of scaffolds**
34 **spanning multiple genes, we find no evidence of clustering of *Sycon* *NK* genes. All**
35 ***Sycon* *Antp*-class genes are developmentally expressed, with patterns suggesting**
36 **involvement in cell type specification in embryos and adults, metamorphosis and**
37 **body plan patterning. The present study demonstrates that *ParaHox* genes**
38 **predate the origin of sponges, thus confirming the ghost locus hypothesis⁸, and**
39 **highlights the need to analyze genomes of multiple sponge lineages in order to**
40 **obtain a complete picture of the ancestral composition of the first animal genome.**

41 Sponges (Porifera) are strong candidates for being the earliest extant lineage(s) of
42 animals⁹. The genome sequence of the demosponge *Amphimedon queenslandica*
43 provided rich material for comparative studies on the origins of metazoan
44 developmental genes, cell types and body plan¹. Among others, it has fuelled hypotheses
45 about the origin of one of the most widely studied groups of developmental genes: the
46 *Antennapedia* (ANTP) class homeoboxes, including *Hox*, *ParaHox* and *NK* genes^{5,7,8,10,11}.
47 *Antp* genes have been found in all animals and are involved in multiple developmental
48 processes, including body plan patterning and neurogenesis¹². *Hox*, *ParaHox* and *NK*
49 genes are often found in clusters^{5,13}, and in some animals their expression is temporally
50 or spatially correlated to their position within the cluster (this being known as
51 colinearity)¹². The *Amphimedon* genome contains eight *NK* genes, but neither *Hox* nor
52 *ParaHox* genes are present⁵. Six *NK* genes are linked in a tight cluster, and their simple
53 embryonic and larval expression patterns are not consistent with colinearity^{5,6}. Lack of
54 *Hox* and *ParaHox* genes in *Amphimedon*, and also in the ctenophore *Mnemiopsis leidyi*⁷,
55 has previously been interpreted as reflecting the ancestral condition, and gave rise to
56 the ParaHoxozoa hypothesis, in which all animal lineages apart from poriferans and
57 ctenophores are collectively known as the ParaHoxozoa. Others¹⁰ have interpreted the
58 phylogenetic evidence differently, suggesting that both *Hox* and *ParaHox* genes were
59 originally present in sponges, but have subsequently been lost. This view has been
60 recently revived by identification of *Hox* and *ParaHox* “ghost loci” (regions that display
61 synteny to bilaterian *Hox* and *ParaHox* loci, but lack the *Hox/ParaHox* genes themselves)
62 in the genome of *Amphimedon*⁸.

63 We expected that expanding the range of sequenced sponge genomes would provide
64 new information about the evolutionary history of genes important for the origin and

65 evolution of the animal kingdom. Calcisponges form a poriferan lineage which has been
66 separated from the demosponges for at least 600 million years⁹. We have recently
67 started analysis of the developmental toolkits of two calcisponges, *Sycon ciliatum* and
68 *Leucosolenia complicata*^{2,4,14}. Here, we have searched for *Antp*-class homeobox genes in
69 transcriptomic and genomic assemblies of these species.

70 We retrieved ten *Antp*-class homeodomains in *Sycon*, constituting nine transcripts (one
71 with two homeoboxes), and twelve *Antp*-class homeodomains in *Leucosolenia*,
72 constituting nine transcripts (one with four homeoboxes) (Supplementary dataset 1).
73 Our phylogenetic analyses demonstrate that the repertoire of *Antp*-class genes is similar
74 between the two calcisponges, but strikingly different than in the demosponge
75 *Amphimedon*. Calcisponges and demosponges have clear orthologues of the bilaterian
76 genes *Hex* and *Msx*; calcisponges also have a clear *NK5 (Hmx)* orthologue, which seems
77 to be lacking in *Amphimedon*. In contrast, this demosponge has possible *Bsh*, *BarH/BarI*
78 and *Tlx* genes, which are not recognizable in calcisponges. While in *Amphimedon* there is
79 a single gene associated with the bilaterian NK2/3/4 clade¹⁵, several paralogues are
80 present in the two calcisponges. They contain multi-homeobox genes with non-
81 orthologous relationships between *Sycon* and *Leucosolenia*, and other genes containing
82 single homeoboxes in the NK2/3/4 clade. Affiliation of *Sycon* and *Leucosolenia* *NKB* and
83 *NKG* and the *LcoNKF* genes with a particular bilaterian NK family is not clear. No Hox
84 genes were found; however, a pair of the calcisponge genes showing affinities with the
85 ParaHox Cdx subfamily given the concordance of Neighbour-Joining (NJ) and Maximum
86 Likelihood (ML) analyses (Fig. 1, Extended Data Fig. 1).

87 Given the importance of this potential assignment, we performed further phylogenetic
88 analyses of these putative Cdx orthologues. In addition to the ELEKEF motif which is

89 shared by many Hox and ParaHox, but not NK-type homeodomains, the Cdx family has
90 some distinctive residues in its homeodomain, most notably the YIT motif present only
91 in a small number of other ANTP class homeodomains (Supplementary note 1 and
92 Extended Data Fig. 2). Phylogenetic analyses focused on these few families in addition to
93 families represented in sponges, based on greater taxon sampling than in the overall
94 classification, produced a significantly supported clustering of *SciCdx* and *LcoCdx* with
95 *Cdx* genes from other species in NJ, ML and Bayesian analyses (Extended Data Figs. 3-5).

96 We have also investigated the genomic neighbourhood of *SciCdx* to help resolve the
97 identity of this homeobox gene (Fig. 2; Supplementary note 2 and Supplementary Table
98 1). From the 14 genes on the *SciCdx* scaffold that have clear human orthologues
99 (Supplementary Table 1), four are orthologues of genes linked to ParaHox loci in
100 humans. One of these, SAR1A/B also has a conserved neighbouring relationship with the
101 ParaHox cluster in the cnidarian, *Nematostella vectensis* (Fig. 2a, b). Although these gene
102 numbers are insufficient to reach statistical significance, the neighbour relationships are
103 consistent with the identification of *SciCdx* as a ParaHox gene. Furthermore, as one
104 would expect from the ghost locus hypothesis and the identification of *SciCdx* as a bona
105 fide *ParaHox* gene, we also find clustering of *Sycon* orthologues of ParaHox and Hox
106 neighbour genes into two distinct groups in the *Sycon* genome to statistically significant
107 levels (Fig. 2c-e). Altogether the evidence is consistent with the identification of *SciCdx*
108 and *LcoCdx* as the first examples of sponge *ParaHox* genes.

109 All of the *Sycon* NK genes are found on separate scaffolds (Extended data fig. 6) with
110 multiple additional genes surrounding the homeobox genes. We interpret this as the
111 ancient NK cluster having been broken apart in the *Sycon* genome. Alternatively, our
112 current assembly is not sufficient to provide evidence of a cluster with multiple genes

113 inserted between the *NK* genes. It has been previously shown that arrangements of *NK*
114 genes are variable between different species, ranging from intact and conserved *NK*
115 clusters^{5,16,17} to clusters that are partially broken^{15,18}.

116 We studied the expression of *Antp*-class genes in *Sycon* using a combination of *in situ*
117 hybridization with quantitative transcriptome analysis (Fig. 3, Supplementary note 3
118 and Extended Data Figs 7-9). For all *Antp*-class genes, except *SciHex*, expression can be
119 detected in oocytes and during cleavage (Fig. 3a and Extended Data Fig. 7a-g). During
120 embryogenesis, the most striking expression domain of the majority of the identified
121 genes is the cruciform cells, which are putative larval sensory cells^{2,14}. Beginning at the
122 four-cell stage, stronger expression of *SciNKA* marks the cytoplasm destined to become
123 partitioned into the cruciform cells (Fig. 3d and Extended Data Fig. 7h-q) and expression
124 of *SciHmx* is also markedly elevated in these cells (Fig. 3e). *SciNKC* and *NKD* are uniquely
125 and strongly expressed in the cruciform cells of more advanced (pre-inversion stage)
126 embryos (Fig. 3f-g). *SciNKA* is additionally detected in macromeres of embryos and
127 larvae, and along with *SciNKG* and *SciNKB* domains, forms a set of adjacent stripes along
128 the larval anterior-posterior axis (Fig. 3h-p). This pattern is reminiscent of “striped”
129 patterns reported for *NK* genes in bilaterians¹⁹, and might be indicative of roles for the
130 calcisponge *NK* genes in axial patterning of the larval body plan or in cell type
131 determination with cells destined for specific fates distributed along the larval axis. For
132 example, the macromeres give rise to the pinacocytes of the outer cell layer²⁰, and the
133 *SciNKG*-positive micromeres are good candidates for future sclerocytes (spicule
134 producing cells) given co-expression of *SciNKG* and sclerocyte-specific carbonic
135 anhydrases (*scl-CA1* and *scl-CA2*)²¹. All of the *Antp*-class genes except *SciNKC* and *SciNKD*
136 are expressed during metamorphosis (Fig. 3b, Extended Data Fig. 9) in sub-populations

137 of cells in all three cell layers (Extended Data Fig. 7). The clear expression of *SciCdx* in
138 the inner cell mass during formation of the choanocyte chamber (Fig. 3q) is particularly
139 striking in light of the recently revived notion of homology of the sponge choanoderm
140 with bilaterian endoderm¹⁴, as *ParaHox* expression in bilaterians is often associated
141 with the developing gut. In adults, most of the *Antp*-class genes display differential
142 expression along the body axis (Fig. 3b and Extended Data Table 1). *SciNKG* and *SciNKA*
143 are strongly expressed in sclerocytes, while *SciMsx* and *SciHmx* transcripts are
144 predominantly detected within and around the oscular sphincter (Fig. 3r-w).

145 In summary, analysis of *Antp*-class genes in a previously understudied lineage of
146 sponges allowed us to demonstrate pre-poriferan ancestry of *ParaHox* genes, thus
147 confirming the ghost locus hypothesis and rejecting the *ParaHoxozoa* hypothesis of
148 *Hox/ParaHox* gene origins. Expression patterns of the identified genes indicate that
149 developmental functions of *Antp*-class genes also predate poriferans, with probable
150 involvement in specification of potentially homologous structures
151 (choanoderm/endoderm and cross cells/sensory cells) as well as morphological
152 novelties (calcareous spicules). Differences in *Antp*-class gene repertoires between the
153 demosponge *Amphimedon* and the two calcisponges, *Sycon* and *Leucosolenia*, are
154 striking, and the fact that both classes of sponges share a subset of genes with
155 bilaterians indicates independent gene loss events in the two poriferan lineages.

156 **Methods**

157 Genome and transcriptome assemblies will be described in detail elsewhere (Adamski,
158 Leininger and Adamska, unpublished results). Briefly, the high quality draft genome
159 assembly of *S. ciliatum* was generated using two (360bp and 530bp) paired-end libraries
160 and several mate-pair libraries ranging from 2.0 to 9.0 kb and the preliminary draft

161 assembly of *L. complicata* was generated from a single 295 bp paired-end library, all
162 prepared and sequenced by Illumina technology. Assembly was performed using
163 SOAPdenovo2²² and scaffolding using SSPACE v2.2²³, and resulted in N50 = 150kbp and
164 450bp for *S. ciliatum* and *L. complicata*, respectively. Transcriptomes were assembled
165 using Trinity²⁴. For *S. ciliatum*, genomic scaffolds and transcripts of sponge origin (as
166 opposed to those derived from associated organisms) were identified by aligning the
167 resulting assembly to reads from an Illumina sequenced library obtained from
168 laboratory grown, eukaryotic- contamination free juveniles. The calcisponge *Antp*-class
169 sequences were retrieved from these assemblies using TBLASTN with representative
170 query homeodomain sequences from *Amphimedon queenslandica*, *Mus musculus*,
171 *Tribolium castaneum* and *Branchiostoma floridae*. For phylogenetic analysis, we selected
172 *B. floridae* and *T. castaneum* to provide a framework for the classification of the sponge
173 sequences, as these species have been shown to collectively contain homologues of all
174 major bilaterian *Antp*-class genes²⁵. Their homeodomain sequences were extracted from
175 HomeoDB²⁶. Prottest3.0²⁷ and Modelgenerator v0.85²⁸ were used to determine the best
176 suitable model of sequence evolution (LG+G). Phylogenetic analyses were based on
177 Neighbour-Joining (Phylip v3.69), Maximum Likelihood (PhyML v3.0) and Bayesian
178 inference (MrBayes v3.1.2) methods. Gene expression was studied using available
179 packages^{29,30}. *S. ciliatum* *Antp*-class gene amplification, cloning, sequencing, probe
180 production and single *in situ* hybridization were performed as described previously². In
181 the double *in situ* experiment, samples were hybridized simultaneously with
182 digoxigenin-labelled *SciNKB* probe and fluorescein-labelled *SciNKG* probe. After
183 detection of the digoxigenin-labelled probe with NBT/BCIP substrate, the anti-
184 digoxigenin antibody was removed by two 5-minute washes in 0.1 M glycine/HCl, pH
185 2.2/0.1% Tween 20 followed by three additional maleic acid buffer washes. A second

186 round of pre-blocking, antibody incubation and post-antibody washes were as in the
187 single probe protocol with the exception that anti-Fluo-AP antibody was used and the
188 colour developed using Fast Red tablets (Roche) according to manufacturer's
189 instructions. Photographs demonstrating gene expression are representative of multiple
190 individual specimens, with following replicates: oocytes and embryos: 3-4 small pieces
191 of adult sponge, each containing tens to hundreds of oocytes or embryos of a given
192 developmental stage; young syconoid sponges: at least 5 individual specimens;
193 juveniles: small petri dishes or wells of multi-well plates containing at least 10 juveniles.
194 At least two independent experiments were carried for each probe.

195 **References:**

- 196 1 Srivastava, M. *et al.* The *Amphimedon queenslandica* genome and the evolution of animal
197 complexity. *Nature* **466**, 720-726,
198 doi:[http://www.nature.com/nature/journal/v466/n7307/abs/nature09201.html#supplemen](http://www.nature.com/nature/journal/v466/n7307/abs/nature09201.html#supplementary-information)
199 [tary-information](http://www.nature.com/nature/journal/v466/n7307/abs/nature09201.html#supplementary-information) (2010).
- 200 2 Fortunato, S. *et al.* Genome-wide analysis of the sox family in the calcareous sponge *Sycon*
201 *ciliatum*: multiple genes with unique expression patterns. *EvoDevo* **3**, 14 (2012).
- 202 3 Larroux, C. *et al.* Genesis and expansion of metazoan transcription factor gene classes. *Mol*
203 *Biol Evol* **25**, 980 - 996 (2008).
- 204 4 Fortunato, S., Leininger, S. & Adamska, M. Evolution of the Pax-Six-Eya-Dach network: the
205 calcisponge case study. *EvoDevo* **5**, 23 (2014).
- 206 5 Larroux, C. *et al.* The NK Homeobox Gene Cluster Predates the Origin of Hox Genes. *Current*
207 *biology : CB* **17**, 706-710 (2007).
- 208 6 Fahey, B., Larroux, C., Woodcroft, B. J. & Degnan, B. M. Does the High Gene Density in the
209 Sponge NK Homeobox Gene Cluster Reflect Limited Regulatory Capacity? *The Biological*
210 *Bulletin* **214**, 205-217 (2008).
- 211 7 Ryan, J. *et al.* The homeodomain complement of the ctenophore *Mnemiopsis leidyi* suggests
212 that Ctenophora and Porifera diverged prior to the ParaHoxozoa. *EvoDevo* **1**, 9,
213 doi:10.1186/2041-9139-1-9 (2010).
- 214 8 Mendivil Ramos, O., Barker, D. & Ferrier, David E. K. Ghost Loci Imply Hox and ParaHox
215 Existence in the Last Common Ancestor of Animals. *Current biology : CB* **22**, 1951-1956
216 (2012).
- 217 9 Wörheide, G. *et al.* in *Advances in Marine Biology* Vol. Volume 61 (eds Maria J. Uriz Manuel
218 Maldonado Mikel A. Becerro & Turon Xavier) 1-78 (Academic Press, 2012).
- 219 10 Peterson, K. J. & Sperling, E. A. Poriferan ANTP genes: primitively simple or secondarily
220 reduced? *Evolution & Development* **9**, 405-408, doi:10.1111/j.1525-142X.2007.00179.x
221 (2007).
- 222 11 Ferrier, D. E. K. in *Hox Genes: Studies from the 20th to the 21st Century* Vol. 689 *Advances in*
223 *Experimental Medicine and Biology* (ed J. S. Deutsch) 91-100 (Springer-Verlag Berlin, 2010).

- 224 12 Garcia-Fernandez, J. The genesis and evolution of homeobox gene clusters. *Nat Rev Genet* **6**,
225 881-892 (2005).
- 226 13 Chourrout, D. *et al.* Minimal ProtoHox cluster inferred from bilaterian and cnidarian Hox
227 complements. *Nature* **442**, 684-687,
228 doi:http://www.nature.com/nature/journal/v442/n7103/suppinfo/nature04863_S1.html
229 (2006).
- 230 14 Leininger, S. *et al.* Developmental gene expression provides clues to relationships between
231 sponge and eumetazoan body plans. *Nat Commun* **5**, doi:10.1038/ncomms4905 (2014).
- 232 15 Luke, G. N. *et al.* Dispersal of NK homeobox gene clusters in amphioxus and humans. *Proc.*
233 *Natl. Acad. Sci. U. S. A.* **100**, 5292-5295, doi:10.1073/pnas.0836141100 (2003).
- 234 16 Hui, J. H. L. *et al.* Extensive Chordate and Annelid Macrosynteny Reveals Ancestral Homeobox
235 Gene Organization. *Mol. Biol. Evol.* **29**, 157-165, doi:10.1093/molbev/msr175 (2012).
- 236 17 Hui, J. H. L., Holland, P. W. H. & Ferrier, D. E. K. Do cnidarians have a ParaHox cluster?
237 Analysis of synteny around a Nematostella homeobox gene cluster. *Evolution & Development*
238 **10**, 725-730, doi:10.1111/j.1525-142X.2008.00286.x (2008).
- 239 18 Seo, H.-C. *et al.* Hox cluster disintegration with persistent anteroposterior order of
240 expression in *Oikopleura dioica*. *Nature* **431**, 67-71,
241 doi:http://www.nature.com/nature/journal/v431/n7004/suppinfo/nature02709_S1.html
242 (2004).
- 243 19 Saudemont, A. *et al.* Complementary striped expression patterns of NK homeobox genes
244 during segment formation in the annelid *Platynereis*. *Dev. Biol.* **317**, 430-443,
245 doi:<http://dx.doi.org/10.1016/j.ydbio.2008.02.013> (2008).
- 246 20 Amano, S. & Hori, I. Metamorphosis of Calcareous Sponges .2. Cell Rearrangement and
247 Differentiation in Metamorphosis. *Invertebrate Reproduction & Development* **24**, 13-26
248 (1993).
- 249 21 Voigt, C., Adamski, M., Sluzek, M. & Adamska, M. Genome-wide screening identifies two
250 carbonic anhydrases involved in biomineralization in calcareous sponges. *BMC Evol. Biol.*
251 *pending minor revisions* (2014).
- 252 22 Luo, R. *et al.* SOAPdenovo2: an empirically improved memory-efficient short-read de novo
253 assembler. *GigaScience* **1**, 18, doi:10.1186/2047-217X-1-18 (2012).
- 254 23 Boetzer, M., Henkel, C., Jansen, H., Butler, D. & Pirovano, W. Scaffolding pre-assembled
255 contigs using SSPACE. *Bioinformatics (Oxford, England)* **27**, 578-579,
256 doi:10.1093/bioinformatics/btq683 (2011).
- 257 24 Grabherr, M. *et al.* Full-length transcriptome assembly from RNA-Seq data without a
258 reference genome. *Nature biotechnology* **29**, 644-652, doi:10.1038/nbt.1883 (2011).
- 259 25 Takatori, N. *et al.* Comprehensive survey and classification of homeobox genes in the
260 genome of amphioxus, *Branchiostoma floridae*. *Dev. Genes Evol.* **218**, 579-590,
261 doi:10.1007/s00427-008-0245-9 (2008).
- 262 26 Zhong, Y.-F. & Holland, P. HomeoDB2: functional expansion of a comparative homeobox
263 gene database for evolutionary developmental biology. *Evolution & development* **13**, 567-
264 568, doi:10.1111/j.1525-142X.2011.00513.x (2011).
- 265 27 Darriba, D., Taboada, G., Doallo, R. & Posada, D. ProtTest 3: fast selection of best-fit models
266 of protein evolution. *Bioinformatics (Oxford, England)* **27**, 1164-1165,
267 doi:10.1093/bioinformatics/btr088 (2011).
- 268 28 Keane, T., Creevey, C., Pentony, M., Naughton, T. & McInerney, J. Assessment of methods for
269 amino acid matrix selection and their use on empirical data shows that ad hoc assumptions
270 for choice of matrix are not justified. *BMC Evolutionary Biology* **6**, 29 (2006).
- 271 29 Anders, S. & Huber, W. Differential expression analysis for sequence count data. *Genome*
272 *Biology* **11**, R106 (2010).
- 273 30 Li, B. & Dewey, C. RSEM: accurate transcript quantification from RNA-Seq data with or
274 without a reference genome. *BMC Bioinformatics* **12**, 323 (2011).

275 **Acknowledgements**

276 This study was funded by the Sars Centre core budget to MA. Sequencing was performed
277 at Norwegian High Throughput Sequencing Centre funded by the Norwegian Research
278 Council. OMR and DEKF acknowledge support from the BBSRC and the School of
279 Biology, University of St Andrews. We thank Brith Bergum for collecting samples in the
280 2011 season and Rita Holdhus from the Genomics Core Facility at the University of
281 Bergen for help with DNA shearing.

282 **Author Contributions**

283 SAVF carried out gene identification and cloning, analysed gene expression by in situ
284 hybridization, and participated in the phylogenetic analyses and manuscript writing.
285 Mar.A performed sequence assembly, annotation, quantification of gene expression and
286 participated in sample collection, phylogenetic analyses and manuscript writing. OMR
287 performed the synteny analyses, participated in phylogenetic analyses, manuscript
288 writing and design of the research approach for synteny and phylogenetic analyses. SL
289 isolated samples for sequencing of genomes, generated MP libraries, and participated in
290 sample collection. JL generated samples for sequencing of *S. ciliatum* metamorphosis
291 stages. DEKF participated in design of the research approach for synteny and
292 phylogenetic analyses and writing of the manuscript. Maj.A conceived the study,
293 participated in data analysis, sample collection and writing of the manuscript.

294 **Author Information**

295 Genome assembly of *Sycon ciliatum* and the coding sequences and their translations
296 from transcriptome assemblies of *S. ciliatum* and *Leucosolenia complicata* used in this
297 study can be accessed through <http://compagen.zoologie.uni-kiel.de/> and are also
298 deposited at <http://datadryad.org/> (doi:10.5061/dryad.tn0f3). The RNA-Seq data are

299 deposited at www.ebi.ac.uk/arrayexpress (E-MTAB-2430, 2431, 2890), and the cloned
300 coding sequences of *S. ciliatum* *Antp*-class genes are deposited at European Nucleotide
301 Archive under accession codes HGXXXXXX to YYYYYY.

302 Reprints and permissions information is available at
303 www.nature.com/reprints. The authors declare no competing financial interests.
304 Readers are welcome to comment on the online version of the paper. Correspondence
305 and requests for materials should be addressed to M. Adamska
306 (maja.adamska@sars.uib.no).

307

308 **Figure Legends**

309 **Figure 1. Phylogenetic tree of the *Antp*-class homeodomains.** A neighbour joining
310 (NJ) tree is displayed. Three support values are shown: left/black value is NJ bootstrap
311 support, middle/blue is Maximum Likelihood bootstrap support and right/red is
312 posterior probability from Bayesian analysis. Bootstrap values below 10% and
313 posterior probability values below 0.5 are not shown except for associations of
314 calcisponge sequences. The root was determined by using selected *Prd*-class genes as an
315 outgroup. Acronyms of the species are: *Amphimedon queenslandica*, Aqu; *Leucosolenia*
316 *complicata*, Lco; *Sycon ciliatum*, Sci; *Branchiostoma floridae*, Bfl; and *Tribolium*
317 *castaneum*, Tca. Scale bar indicates number of aminoacid substitutions per site.

318 **Figure 2. *SciCdx* synteny and ghost loci simulations.** (a, b) Genomic neighborhoods of
319 *SciCdx* gene and *N. vectensis* ParaHox cluster¹⁷, colors indicate orthologous relationship
320 to human genes with following chromosomal location: yellow – ParaHox neighbours,
321 orange – not linked to Hox/ParaHox loci, yellow-orange – mix of ParaHox and non-
322 Hox/ParaHox neighbours, grey – no orthology. Green lines highlight the conserved *Sar1*-

323 *Cdx* linkage. **(c-e)** Monte Carlo simulations of human Hox and ParaHox neighbour
324 orthologue distributions and their overlap across *S. ciliatum* scaffolds. Arrows indicate
325 numbers of scaffolds with Hox and ParaHox neighbour orthologues and their co-
326 localization in *S. ciliatum*; the observed distributions imply distinct Hox and ParaHox
327 loci.

328 **Figure 3. Expression of *S.ciliatum* Antp-class genes.** **a**, Adult specimen. **b**, expression
329 “heat map”, * indicate statistically higher^{29,30} expression in the apical/top region in
330 comparison to middle (>M) or basal (>B) parts. **c**, oocytes. **d, e**, cleavage and **f-k**, pre-
331 inversion stage embryos: arrows/cc/rainbow colouring – forming cross cells expressing
332 multiple *Antp*-class genes, */mac/blue – macromeres expressing *SciNKA*; mic/red –
333 equatorial micromeres expressing *SciNKG* and *SciNKB*. **l-p**, post-inversion embryos and
334 larvae; **q**, post-larva; **r-w**, top parts of sponges: *SciNKG* and *SciNKA* in sclerocytes, *SciMsx*
335 and *SciHmx* in cells of the oscular sphincter (arrows). Scale bars represent 10µm, except
336 **l-o**: 25µm, **q, r, t, v, w** – 50µm.

337 **Extended Data Figure 1. Phylogenetic tree of the ANTP class homeodomains**
338 **including representative bilaterian and non- bilaterian sequences.** A Neighbour
339 Joining (NJ) tree using the JTT+G (0.5) (1000 bootstraps) model of protein evolution is
340 displayed. A combination of three support values obtained for three phylogenetic
341 methods is shown: left (black) value is bootstrap (BT) support from NJ, middle (blue) is
342 bootstrap support from Maximum Likelihood (LG+G 0.5) and right (red) is posterior
343 probability from Bayesian analysis (LG+G 0.5). BT values below 10% and PP values
344 below 0.5 are not shown except for associations of calcisponge sequences. The root was
345 determined by using selected *Prd*-class genes as an outgroup. Acronyms of the species
346 used are: *Amphimedon queenslandica*, Aqu (Porifera/demosponges); *Leucosolenia*

347 *complicata*, Lco; *Sycon ciliatum*, Sci (Porifera/calcsponges); *Nematostella vectensis*, Nve
348 (Cnidaria); *Trichoplax adhaerens*, Tad (Placozoa); *Mnemiopsis leidyi*, Mle (Ctenophora);
349 *Branchiostoma floridae*, Bfl (Chordata); and *Tribolium castaneum*, Tca (Arthropoda).

350 Scale bar indicates number of aminoacid substitutions per site.

351 **Extended Data Figure 2. Variability of the 'YIT/YIS' homeodomain motif within the**
352 **Cdx/Cad, En and Dbx families in bilaterians, cnidarians, a placozoan and sponges.**

353 Acronyms of the species used are Hsa (*Homo sapiens*), Bfl (*Branchiostoma floridae*), Cte
354 (*Capitella teleta*), Lgi (*Lottia gigantea*), Nve (*Nematostella vectensis*), Tad (*Trichoplax*
355 *adhaerens*), Tca (*Tribolium castaneum*), Sci (*Sycon ciliatum*), Lco (*Leucosolenia*
356 *complicata*), Edi (*Eleutheria dichotoma*), Nv (*Nereis virens*), Pdu (*Platynereis dumerilii*),
357 Pfl (*Ptychodera flava*), Dre (*Danio rerio*), Dme (*Drosophila melanogaster*), Ame (*Apis*
358 *mellifera*), Gga (*Gallus gallus*), Xla (*Xenopus laevis*), Mmu (*Mus musculus*) and Aqu
359 (*Amphimedon queenslandica*).

360 **Extended Data Figure 3. Phylogenetic tree including ANTP class homeodomain**
361 **subfamilies represented in sponges and two additional subfamilies characterized**
362 **by presence of YIT motif (Cdx and En), but excluding divergent *A. queenslandica***
363 **sequences (NK5/6/7a/b and BarH). NJ (JTT, 1000) bootstrap support values are in**
364 **black, ML (LG+G 0.4, 1000 replicates) bootstrap support values are in blue and BY (LG+G**
365 **0.4) posterior probabilities values in red. Only bootstrap support values equal to or**
366 **above 500 are shown. All subfamilies except Cdx are collapsed for clarity. Acronyms of**
367 **the species used are Hsa (*Homo sapiens*), Bfl (*Branchiostoma floridae*), Cte (*Capitella***
368 ***teleta*), Lgi (*Lottia gigantea*), Nve (*Nematostella vectensis*), Tad (*Trichoplax adhaerens*),**
369 **Tca (*Tribolium castaneum*), Sci (*Sycon ciliatum*), Lco (*Leucosolenia complicata*), Edi**
370 **(*Eleutheria dichotoma*), Nv (*Nereis virens*), Pdu (*Platynereis dumerilii*), Pfl (*Ptychodera***

371 *flava*), Dre (*Danio rerio*), Dme (*Drosophila melanogaster*), Ame (*Apis mellifera*), Gga
372 (*Gallus gallus*), Xla (*Xenopus laevis*), Mmu (*Mus musculus*) and Aqu (*Amphimedon*
373 *queenslandica*). Scale bar indicates number of aminoacid substitutions per site.

374 **Extended Data Figure 4. Phylogenetic tree including ANTP class homeodomain**
375 **subfamilies represented in sponges and three additional subfamilies**
376 **characterized by presence of YIT/YIS motifs (Cdx, En and Dbx), but excluding**
377 **some of divergent *A. queenslandica* sequences (NK5/6/7a/b).** NJ (JTT, 1000
378 replicates) bootstrap support values are in black, ML (LG+G 0.4, 1000 replicates)
379 bootstrap support values are in blue and BY (LG+G 0.4) posterior probabilities values in
380 red. Only bootstrap support values equal to or above 500 are shown. All subfamilies
381 except Cdx are collapsed for clarity. Acronyms of the species used are Hsa (*Homo*
382 *sapiens*), Bfl. (*Branchiostoma floridae*), Cte (*Capitella teleta*), Lgi (*Lottia gigantea*), Nve
383 (*Nematostella vectensis*), Tad (*Trichoplax adhaerens*), Tca (*Tribolium castaneum*), Sci
384 (*Sycon ciliatum*), Lco (*Leucosolenia complicata*), Edi (*Eleutheria dichotoma*), Nv (*Nereis*
385 *virens*), Pdu (*Platynereis dumerilii*), Pfl (*Ptychodera flava*), Dre (*Danio rerio*), Dme
386 (*Drosophila melanogaster*), Ame (*Apis mellifera*), Gga (*Gallus gallus*), Xla (*Xenopus laevis*),
387 Mmu (*Mus musculus*) and Aqu (*Amphimedon queenslandica*). Scale bar indicates number
388 of aminoacid substitutions per site.

389 **Extended Data Figure 5. Phylogenetic tree including *Antp*-class homeodomain**
390 **subfamilies represented in sponges and three additional subfamilies**
391 **characterized by presence of YIT/YIS motifs (Cdx, En and Dbx).** NJ (JTT, 1000
392 replicates) bootstrap support values are in black, ML (LG+G 0.4, 1000 replicates)
393 bootstrap support values are in blue and BY (LG+G 0.4) posterior probabilities values in
394 red. Only bootstrap support values equal to or above 500 are shown. All subfamilies

395 except Cdx are collapsed for clarity. Acronyms of the species used are Hsa (*Homo*
396 *sapiens*), Bfl (*Branchiostoma floridae*), Cte (*Capitella teleta*), Lgi (*Lottia gigantea*), Nve
397 (*Nematostella vectensis*), Tad (*Trichoplax adhaerens*), Tca (*Tribolium castaneum*), Sci
398 (*Sycon ciliatum*), Lco (*Leucosolenia complicata*), Edi (*Eleutheria dichotoma*), Nv (*Nereis*
399 *virens*), Pdu (*Platynereis dumerilii*), Pfl (*Ptychodera flava*), Dre (*Danio rerio*), Dme
400 (*Drosophila melanogaster*), Ame (*Apis mellifera*), Gga (*Gallus gallus*), Xla (*Xenopus laevis*),
401 Mmu (*Mus musculus*) and Aqu (*Amphimedon queenslandica*). Scale bar indicates number
402 of aminoacid substitutions per site.

403 **Extended Data Figure 6. *Sycon ciliatum* scaffolds containing NK genes (blue) and**
404 ***Amphimedon queenslandica* scaffold containing cluster of NK genes (modified**
405 **after⁶; green).** Annotation of the neighboring genes (genes within 50 kbp from the NK
406 gene) in *S. ciliatum* was performed using blastp searches against refseq database.

407 **Extended Data Figure 7. Additional expression patterns of ANTP class homeobox**
408 **genes in embryonic development and during metamorphosis.** All of the investigated
409 genes (except Hex, not shown) are expressed in oocytes (a-g). The expression of *SciNKA*
410 is detectable in all blastomeres of the cleavage stage embryos, but the transcripts are
411 concentrated in the corner-most cytoplasm which becomes gradually partitioned to the
412 cross cells (arrows). This subcellular localization of cross-cells enriched transcripts is
413 also observed for *SciNanos*, expression of which, similarly to *SciNKA*, becomes ultimately
414 restricted to cross cells and macromeres in preinversion stage embryos (**l-q**). In
415 metamorphosing postlarvae, *SciNKA* is expressed in the cells of the outer layer (**r**),
416 *SciNKB* and *SciNKG* in (possibly non-overlapping) fractions of cells in the inner cell mass
417 (**s, t**); *SciHex* is weakly expressed throughout the inner cell mass (**u**) and *SciNKC* (**v**) and

418 NKD (not shown) are not detectable in the juveniles. Scale bars represent 10 μ m, except
419 **r-v**: 25 μ m.

420 **Extended Data Figure 8. Samples used for quantification of expression. a-f,**
421 metamorphosis in *S.ciliatum*, stages are based on²⁰ with modifications: Stage I,
422 approximately 12 hours post settlement: large flat cells derived from larval macromeres
423 envelop the inner cell mass composed of former micromeres (**a**); Stage II, approximately
424 24 hours post settlement: single-axis spicules (monaxons) are produced by sclerocytes,
425 which have differentiated from the inner cell mass cells (**b**). Stage III, 2-3 days after
426 settlement: choanocytes which have differentiated from the inner cells mass cells form a
427 single internal chamber (**c**). Stage IV, approximately 4 days after settlement: osculum
428 (exhalant opening) forms at the apical end of the spherical juvenile; first tri-radial
429 spicules become apparent (**d**). Stage V, approximately 10 days after settlement: the
430 juvenile is elongated along the apical-basal axis, long straight spicules form a crown
431 around the osculum (**e**). Young syconoid sponges, approximately 8 weeks after
432 settlement (**f**). (**a-e**) are photographs of live specimens in culture; photographs **a-d** are
433 top (apical) views, cartoon representations of sections and photograph **e** are side views.
434 Scale bars represent 100 μ m, except **f**: 1mm. (**g**) Details of replicates used for the
435 analysis. Several hundred juveniles were used in each sample. (**h**) Plot demonstrating
436 results of principal component analysis of the metamorphosis series and axial dissection
437 series of non-reproductive adults calculated according to²⁹ and utilizing information of
438 the top 500 differentially expressed genes as in default parameters. Metamorphosis
439 stages and parts of sponges are colour-coded, with the ovals added manually for easier
440 visualization of similarities and differences between the samples. Progress of
441 development, starting from freshly released larvae and until emergence of adult, but not

442 yet reproductive sponges, is indicated by arrows. Note similarities of samples within
443 replicates and with neighbouring stages of the metamorphosis series, and
444 distinctiveness of the top (apical) samples from the basal and middle samples of the
445 adults. (i) Heatmap representation of sample-to-sample distances among all samples
446 used in this study, calculated according to²⁹ and based on expression of all coding genes
447 in *S. ciliatum* (approximately 18K sequences). Note that replicates and neighbouring
448 stages group together, as indicated by highlighting.

449 **Extended Data Figure 9.** Heat-map representation of expression profiles demonstrated
450 in Fig. 3B in the main text, but with data from individual libraries presented separately.

451 **Extended Data Table 1. Quantification of differences in expression levels between**
452 **top, middle and bottom parts of non-reproductive adult specimens of *S. ciliatum*.**

453 'expression level' was calculated as sum of the posterior probability of each read coming
454 from a given gene over all reads³⁰ scaled using size factors of the libraries²⁹; 'fold
455 change' was calculated between expression levels in middle (middle-top) or bottom
456 (bottom-top) and top part of the sponge; 'adj. p-value' are p-value adjusted for multiple
457 testing with the Benjamini-Hochberg procedure²⁹. Values with statistical significance
458 not less than 90% (adj. p-value ≤ 0.1) and apical expression level higher than in the
459 middle or bottom part of the sponge are indicated by *.





