1    **Genomic analysis of the tryptome reveals molecular mechanisms of gland cell evolution**

2    Leslie S. Babonis[1,3], Joseph F. Ryan[1,2], Camille Enjolras[1], and Mark Q. Martindale[1,2]

3

4    **Affiliations:**

5    [1]Whitney Laboratory for Marine Bioscience, University of Florida, St. Augustine, FL 32080

6    [2]Department of Biology, University of Florida, Gainesville, FL 32611, USA.

7    [3]Corresponding author: babonis@whitney.ufl.edu

8    **Short title: Molecular mechanisms of gland cell evolution**

9

10

11  **Abstract**

12  Understanding the drivers of morphological diversity is a persistent challenge in evolutionary

13  biology. Here, we investigate functional diversification of secretory cells in the sea anemone

14  *Nematostella vectensis* to understand the mechanisms promoting cellular specialization across

15  animals. We demonstrate regionalized expression of gland cell subtypes in the internal ectoderm

16  of *N. vectensis* and show that adult gland cell identity is acquired very early in development. A

17  phylogenetic survey of trypsins across animals suggests this gene family has undergone

18  numerous expansions. We reveal unexpected diversity in trypsin protein structure and show that

19  trypsin diversity arose through independent acquisitions of non-trypsin domains. Finally, we

20  show that trypsin diversification in *N. vectensis* was effected through a combination of tandem

21  duplication, exon shuffling, and retrotransposition. Together we reveal that numerous

22  evolutionary mechanisms drove trypsin duplication and divergence during the morphological

23  specialization of cell types and suggest the secretory cell phenotype is highly adaptable as a

24  vehicle for novel secretory products.

25

26  **Introduction**

27  The development of new tissue layers provides the opportunity to spatially segregate cell types

28  enabling the compartmentalization of different functions. Cnidarians are diploblasts, comprised

29  of an internal endodermal epithelium separated from an external ectodermal epithelium by a

30  largely acellular matrix called mesoglea. Anthozoans (corals, sea anemones, and their kin) are

31  unusual among cnidarians in their possession of internal tissues (pharynx and mesenteries) that

32  arise by secondary epithelial fold morphogenesis following completion of gastrulation (1).

33  Additional growth and differentiation of both internalized layers results in the morphogenesis of

2

34    the pharynx and mesenteries and in an adult form quite different from that of medusozoans. In

35    anthozoans, both layers (endoderm and ectoderm) are in contact with the gastric cavity, whereas

36    in medusozoans (and, indeed, most other animals), the gastrovascular cavity is lined only by

37    endoderm.

38

39    *Nematostella vectensis,* the starlet sea anemone, has become a valuable model for studies of

40    animal body plan evolution (2-6); yet, little is known about the extent of cell diversity in the

41    tissues that comprise the pharynx and mesenteries. The endodermal component of the

42    mesenteries houses the germ cell precursors and two types of muscle cells and the few recent

43    studies of the mesenteries in *Nematostella* have focused largely on these endodermal functions

44    (7-9). The ectodermal component of the mesenteries is known to be populated by cnidocytes and

45    gland cells (10) and two recent studies demonstrated the expression of multiple proteases in the

46    mesenteries of *N. vectensis* (11, 12). Trypsins are the largest family of proteases, and although

47    they have diverse functions, most trypsins are secreted to the extracellular environment and are

48    therefore expressed in zymogen-type gland cells (13). A previous study cataloging trypsin

49    diversity from prokaryotes and eukaryotes identified 75 trypsins in the genome of *N. vectensis*

50    (14),  suggesting that the few cell types identified anatomically as zymogen gland cells (10) may

51    belie the digestive capacity of the mesenteries.

52

53    We sought to understand the evolutionary mechanisms promoting functional diversification at

54    the cell and tissue levels in the mesenteries of *N. vectensis,* and to characterize the evolutionary

55    history of a large (super)family of proteases expressed abundantly in the mesenteries.  Building

56    on a previous study using RNA-seq to characterize the expression profile of the mesenteries in *N.*

57    *vectensis* (11), we show that the continuous epithelium comprising the internal ectoderm in *N.*

58    *vectensis* is functionally partitioned into different regions associated with distinct morphologies

59    and functions. Additionally, we show numerous lineage-specific expansions of trypsins and that

60    trypsin diversification arises through novel domain acquisition. Finally, we propose a model by

61    which the expansion of trypsins may have promoted specialization of gland cell subtype in

62    cnidarians.

63

64    **Results**

65    *Morphology and function of the internal ectoderm*

66    We examined the fine structure of the internal and external ectoderm in the region of the mouth

67    of *N. vectensis* during feeding for evidence of morphological and functional variation (Figure 1).

68    Cells in the external ectoderm around the mouth are organized into a low cuboidal type

69    epithelium that covers the closed mouth between feeding events (Figure 1A-D). In the presence

70    of prey, the pharynx is partially everted, exposing the tall columnar epithelium of the pharyngeal

71    ectoderm (Figure 1E-G). After passing through the pharynx (Figure 1H), ingested prey remains

72    in contact with the ectodermal portion of the mesenteries, which is populated by cnidocytes and

73    gland cells (Figure 1I-K).

74

75    The pharyngeal ectoderm is populated by numerous distinct electron dense (zymogen-secreting)

76    and electron lucent (mucus-secreting) gland cells (Figure 2A-F). The adjacent non-secretory cells

77    in this epithelium have distinctive apical electron-dense vesicles (Figure 2F). The proximal

78    region of the mesentery (adjacent to the body wall) is comprised of endoderm while the distal

79    portion (the free edge) is comprised of ectoderm (Figure 2G). The ectodermal region gives rise to

80    both the cnidoglandular tract at the most distal extent (Figure 2H-J, M-O) and the ciliated tract

81    more proximally (Figure 2K,L). Thin sections of the ectodermal mesentery in the oral region

82    (near the pharynx) show abundant zymogen gland cells (Figure 2H-J), some of which contain

83    secretory vesicles with heterogeneous contents (Figure 2J). Ciliated tracts are short and are

84    present only in the oral end of each mesentery. Cells of the ciliated tract are highly proliferative

85    and have apical motile cilia but do not have other distinguishing features (Figure 2K,L). The

86    aboral mesentery lacks a ciliated tract but the cnidoglandular tract still contains numerous

87    distinct zymogen gland cells, some with motile apical cilia (Figure 2M-O). Mucus-secreting cells

88    were found in the pharyngeal ectoderm (Figure 2D) and in the external ectoderm of the body

89    wall and tentacles (Supplemental file 1new), but never in the endoderm.

90

91    *Proteolytic enzymes are expressed in the developing mesenteries*

92    We previously identified numerous genes encoding different classes of proteases to be

93    upregulated in the adult mesentery of *N. vectensis* (11) . Using *in situ* hybridization we examined

94    the spatial and temporal expression of 10 proteases of various classes during early development

95    of the pharynx/mesenteries to understand the ontogeny of digestive function and the onset of

96    terminal gut cell differentiation. All genes examined were expressed in individual ectodermal

97    cells of the mesenteries at the primary polyp stage, just after metamorphosis (Figure 3A,B); two

98    protease genes (NVJ_82725 and NVJ_83864) were also expressed in the pharyngeal ectoderm of

99    the primary polyp. There was surprisingly little variation in the onset of protease expression,

100   although serine proteases (trypsins) consistently exhibited expression in the early planula stage

101   before differentiation of the presumptive pharynx and mesenteries (Figure 3B). Double

102   fluorescent *in situ* hybridization for two metalloprotease genes (NVJ_88668 and NVJ_2109)

5

103    indicates both co-expression of these two enzymes in few cells at the aboral end of the pharynx

104    and independent expression of the two genes in distinct cells of the ectodermal mesenteries in the

105    late tentacle bud stage (Figure 3C).

106

107    The surprising lack of any obvious spatial segregation in protease expression led us to

108    hypothesize that many proteases may be co-expressed together in the few anatomically

109    distinguishable gland cells identified above (Figure 2). Using the raw data from a single-cell

110    RNA-Seq study published previously (15), we show co-expression of six of the ten proteases we

111    studied by *in situ* hybridization in a single putative gland cell (Figure 3D). Using the raw data

112    from the same study and a very low cutoff for gene expression ($N \geq 1$ read), we examined more

113    fully the co-expression of the large superfamily of trypsin proteases and  found 6,727 cells

114    expressing at least one trypsin gene. Nearly 50% of the trypsin-expressing cells (3,282/6,727)

115    appear to express only a single trypsin, while the remaining cells exhibited co-expression of up

116    to 24 trypsins (Figure 3E). For each trypsin, we then examined the relationship between the

117    ubiquity of expression (the total number of cells in which that trypsin is expressed) and the

118    number of cells in which it is co-expressed with other trypsins and found a strong positive

119    correlation (Figure 3F), confirming that the trypsins with the broadest expression profiles were

120    most likely to be co-expressed with other trypsins.

121

122    *The tryptome is lineage-specific*

123    To characterize the tryptome (all proteins with a trypsin domain) of *N. vectensis*, we searched the

124    JGI gene models (https://genome.jgi.doe.gov/Nemve1/Nemve1.home.html) for all sequences

125    containing a significant Trypsin or Trypsin_2 domain using hmmsearch (HMMER 3.1b2;

6

126   http://hmmer.org) and constructed domain architecture diagrams for each protein (Figure 4). Of

127   the 72 trypsin gene models that remained after curation (see Methods), 28 encode a trypsin

128   domain but lack any other conserved domains and the other 44 encode a trypsin domain and at

129   least one additional conserved domain. In total, trypsin domains were found in association with

130   24 other domains in *N. vectensis*. To determine if any of these associated domains were

131   overrepresented in the tryptome, we compared the abundance of trypsin-associated domains in

132   the tryptome and in the proteins predicted from the JGI gene models (N = 27,273 protein

133   predictions). Six domains were found to be represented in high abundance (≥10%) in the

134   tryptome: DIM, ShK, Lustrin_cystein, Sushi, MAM and SRCR (Figure 4A). The DIM and

135   Lustrin_cysteine domains are present in low abundance throughout the predicted proteome (1

136   and 4 total domains, respectively), artificially inflating their perceived abundance in the

137   tryptome. For ShK, Sushi, MAM, and SRCR, ≥15% of the domains found in the proteome were

138   associated with trypsins.

139

140   To determine whether the makeup of the tryptome was unique to *N. vectensis*, we searched for

141   proteins with these same domain architectures in representatives from all domains of life (other

142   cnidarians, bilaterians, non-metazoan eukaryotes, and a selection of prokaryotes). Two domain

143   architectures were found to be present across taxa: those with only a trypsin domain, and those

144   with a trypsin and a PDZ domain (Figure 4B). Trypsins appear to have expanded considerably

145   after the origin of animals, as both choanoflagellate lineages had fewer than five trypsins but the

146   ctenophore *Mnemiopsis leidyi* and the placozoan *Trichoplax adhaerens* (representing two of the

147   earliest diverging animal lineages) both have at least twenty. Among the cnidarians, other

148   actinarians (sea anemones) shared more trypsin proteins in common with *N. vectensis* than any

7

149   other lineage; however, we identified three trypsins (NVJ_105548, NVJ_105271, and

150   NVJ_199428) specific to *N. vectensis* that were absent even from *Edwardsiella lineata*

151   (representing the genus sister to *Nematostella*).

152

153   *Trypsins diversified independently in cnidarians and bilaterians*

154   To characterize the diversification of animal trypsins, we built a phylogeny of trypsin domains

155   from taxa representing each of the five major animal lineages: bilaterians, cnidarians,

156   placozoans, sponges, and ctenophores. Using this tree, we identify six clades of trypsins and

157   classify them by their function in human: a non-catalytic group, the intracellular trypsins,

158   tryptases and transmembrane trypsins, trypsins involved in coagulation and immune response,

159   chymotrypsins, and the clade including granzymes, pancreatic trypsins, kallikreins, hepatocyte

160   growth factors, and elastases (Figure 5A). Each of these includes representatives from

161   bilaterians, cnidarians, and at least one placozoan, sponge, or ctenophore and likely represent the

162   suite of trypsin clades present in the last common ancestor of animals. The *N. vectensis* tryptome

163   includes representatives of 5 of 6 clades likely present in the common ancestor of animals. *N.*

164   *vectensis* may have lost representatives of the tryptase/transmembrane clade as this these trypsins

165   appear to be present in *M. leidyi*, *A. digitifera*, and bilaterians (Figure 5A). Proteins with only a

166   trypsin domain were distributed throughout the tree, rather than being clustered in a single clade.

167   Two proteins from *N. vectensis* (NVJ_203589 and NVJ_23745) had divergent trypsin domains

168   and were detected only by the Trypsin_2 HMM, both of which appear to be part of the clade that

169   includes trypsin-PDZ proteins (Supplemental file 2).

170

171    We compared the distribution of conserved domains from different clades of trypsins in *N.*

172    *vectensis* and *H. sapiens* (Figure 5B). In *N. vectensi*s, domain diversity is greatest among the

173    trypsins that group with human chymotrypsins, as trypsin domains from this clade co-occur with

174    14 associated domains. Chymotrypsins from *N. vectensis* share four domains in common with

175    the immune/coagulation group, which is represented by 10 domains. The "pancreatic" group

176    (including granzymes, kallikreins, HGF, and elastase) is characterized by only 5 domains, 3 of

177    which are shared with chymotrypsins. Trypsins from the non-catalytic clade lacks associated

178    domains and the intracellular clade uses unique domains (Death and PDZ). Four trypsin-

179    associated domains (Sushi, EGF_CA, CUB, and FXa_inhibition) were found in the

180    immune/coagulation clade of trypsins in both *N. vectensis* and *H. sapiens* and the PDZ domain is

181    restricted to the intracellular clade of trypsins in both taxa. Surprisingly, there were no other

182    domains found in common between *N. vectensis* and *H. sapiens* trypsins from the same clade.

183    (See supplemental file 3 for distribution of human trypsin domain architectures.)

184

185    These analyses of trypsin distribution across animals revealed a surprisingly large tryptome in

186    sea anemones (*N. vectensis* and *E. lineata*) relative to most other taxa (including another

187    cnidarian, *Hydra magnipapillata*). To determine whether the tryptome of *N. vectensis* is

188    reflective of other cnidarians, we built a phylogeny using representatives of each class within

189    Cnidaria (Figure 6). We identify 16 clades of trypsins that include representatives of at least two

190    lineages of anthozoans and two lineages of medusozoans, suggesting these clades may have been

191    present in the stem cnidarian. Two clades (the trypsin-MAM and trypsin-ShK clades) seem to

192    have undergone further expansion in anthozoans after their divergence from medusozoans.

193

194    *The Nematostella tryptome diversified through numerous mechanisms*

195    To understand the mechanisms generating trypsin diversity in *N. vectensis*, we examined the

196    evolutionary relationships of the 72 trypsin proteins in the tryptome (Figure 7A). Among the 72

197    predicted proteins, 85% (61/72) had all three conserved residues constituting the catalytic triad

198    and are likely to function as proteases, 79% (57/72) were predicted to have a signal peptide and

199    are presumably secreted, and 7% (5/72) were predicted to have a transmembrane domain (see

200    Supplemental File 4). Four of the five clades of trypsins from *N. vectensis* (excluding the

201    intracellular clade) include secreted trypsins, membrane-bound trypsins, and trypsins with

202    divergent sequence that have likely lost their catalytic function. While most (70/72) of the

203    trypsin domains were encoded across multiple exons (Supplemental file 5), two genes

204    (NVJ_128003 and NVJ_216003) lack introns completely.

205

206    Numerous trypsins from the "pancreatic" and chymotrypsin clades were associated with ShK

207    domains. Likewise, over 30% (26/82) of the ShK domains in *N. vectensis* are associated with

208    trypsins (Figure 4A). To determine if the combination of the trypsin and ShK domains may have

209    duplicated together, we built a phylogeny of all 108 ShK domains from the *N. vectensis* proteins

210    predicted from gene models (Figure 7B). Despite the abundance of trypsin-ShK associations, the

211    distribution of trypsin-ShK proteins suggests that this domain combination may have duplicated

212    together only once, giving rise to NVJ_218669 and NVJ_218670; Figure 7A).

213

214    *Trypsin diversity increases through new associations with old domains*

215    Gene age can be estimated using a phylostratigraphic approach; in such analyses, the minimum

216    age of a gene is inferred by identifying the last common ancestor in which the gene is present

10

217    (16, 17). We examined the age of the trypsins found in *N. vectensis* and the age of each

218    associated domain across all domains of life to understand the evolution of trypsin diversity.

219    Trypsin-PDZ and a subset of the trypsin-only proteins likely arose before bacteria/archaea split

220    from eukaryotes, over 2 billion years ago (Figure 8). While trypsin-only proteins are present in

221    every lineage examined, trypsin-PDZ proteins appear to have been lost in several taxa including

222    *C. owczarzaki*, *M. leidyi*, *A. vanhoeffeni*, and *C. cruxmelitensis* (Figure 4). All other associations

223    between trypsin and other conserved domains appear to have originated after the stem metazoan

224    diverged from the rest of life (~800 million years ago). Many of the trypsin-associated domains

225    originated long before they became associated with trypsin; for example, the astacin domain was

226    present in the ancestor of all life but the trypsin-astacin association likely did not arise until the

227    origin of Cnidaria (Figure 8A). By contrast, the SRCR domain and its association with trypsin

228    likely arose in the stem metazoan as trypsin-SRCR proteins were found in *M. leidyi*

229    (Supplemental file 6).

230

231    There does not appear to be a relationship between the age of the domain and the origin of its

232    association with trypsin (Figure 8B). Two trypsin associations were found only in *N. vectensis*:

233    trypsin-DIM (NVJ_199428) and trypsin-WSC (NVJ_105271), and one association was found

234    only in *Edwarsiidae* (*Nematostella + Edwardsiella)*: trypsin-Lustrin_cystein (NVJ_164017). The

235    WSC domain is present throughout eukaryotes (Figure 8A) but was associated with trypsin in

236    only in *N. vectensis*. The Lustrin_cystein domain seems to have arisen in the last common

237    ancestor of parahoxozoa (Placozoa + Cnidaria + Bilateria). These two associations represent

238    extreme cases whereby trypsin diversity in *N. vectensis* arose through acquisition of both young

239    (Lustrin_cystein) and old (WSC) domains.

240

## Discussion

*The not-so-simple cnidarian ectoderm*

Although cnidarian body plans develop from only two tissue layers, morphological diversity

varies widely across taxa. Similarly, only a dozen or so morphologically unique cell types have

been described (10, 18), but cnidarian genomic and functional diversity rival that of any

bilaterian lineage (4, 19). While the ectodermal layer comprising the external and pharyngeal

epithelia may be contiguous, these regions are morphologically and functionally distinct in *N.*

*vectensis* (Figure 1) (10, 12). In this study, we further demonstrate that the continuous layer of

internal ectoderm from the pharynx through the mesenteries is equally heterogeneous. The

pharyngeal ectoderm houses numerous zymogen and mucous cells while the ectoderm of the

mesenteries houses only the former (Figure 2). This anatomical heterogeneity is supported by

variable gene expression: some proteases are expressed throughout the pharyngeal and

mesentery ectoderm, while others are restricted only to the mesentery ectoderm (Figure 3A,B)

(also see (12, 20)). Furthermore, the combinatorial expression of only two proteases can result in

the development of at least three distinct cell types (Figure 3C). Together, the combination of a

diverse tryptome and extensive trypsin co-expression (Figure 3E) suggests cell functional

diversity in cnidarians may well exceed historical expectations.

We found no evidence of endodermal gland cells (zymogen type or mucous type) in our TEM or

*in situ* hybridization results (Figure 2,3, Supplemental file 1). Indeed, all non-neuronal secretory

cells (including mucous cells, zymogen cells, and cnidocytes), are restricted to the ectoderm in

*N. vectensis* but their distribution is heterogeneous. Zymogen cell diversity, for example, is much

12

263    higher in the internal than the external ectoderm (Figure 2, Supplemental file 1). This is

264    consistent with the histological analyses of Frank and Bleakney (10) but seems to be in contrast

265    with the distribution of gland cells in medusozoans. In *Hydra*, for example, zymogen gland cells

266    are found exclusively in the endoderm (21). These observations suggest that the internalization

267    of the ectoderm in anthozoans was a pivotal event in the diversification of specialized zymogen

268    cells. Cell products secreted from the tentacle ectoderm may quickly become diluted in the water

269    column, whereas the closed environment of the gastrovascular cavity limits the space over which

270    secreted products can diffuse; thus, internalization created distinct selective pressures in different

271    regions of the ectoderm. The selective pressure to secrete digestive enzymes into the enclosed

272    gastrovascular cavity may have driven the development of gland cells in the internal ectoderm of

273    anthozoans and the endoderm of medusozoans (and many bilaterians). As such, we see no reason

274    to homologize the ectoderm of anthozoan mesenteries and the endodermal lining of the

275    vertebrate midgut/pancreas (12). We consider it more likely that these tissues have converged on

276    similar morphologies and gene expression profiles in response to similar selection pressures

277    associated with extracellular digestion.

278

13

279 *N. vectensis trypsins have many putative functions*

280 The trypsin domain catalyzes the cleavage of polypeptides at internal amino acid residues and is

281 therefore essential for processing large proteins into smaller peptide chains. Digestive trypsins

282 are synthesized in secretory cells with zymogen type secretory granules where they are packaged

283 into vesicles for release into the gut. We show that there are at least 10 morphologically distinct

284 zymogen gland cell types in the pharyngeal and mesentery ectoderm of *N. vectensis* (Figure 2).

285 Further, we demonstrate that numerous proteases are expressed in the same tissues (Figure 3)

286 and that the vast majority of trypsins in *N. vectensis* encode a signal peptide (Figure 7A). Using

287 published single-cell expression data (15), we identified ten putative gland cells that express

288 trypsins, at least two of which also express synaptotagmin (Supplemental file 4), which

289 facilitates fusion of the vesicle with the cell membrane during regulated secretion. These

290 combined features point to the expression of trypsins in secretory cells of the internal ectoderm

291 and strongly support a role for trypsins in extracellular protein degradation in *N. vectensis*.

292

293 Numerous trypsins were expressed outside of the putative gland cells identified by Sebe-Pedros

294 et al (15). At least 20 cells categorized by these authors as neural cells exhibited trypsin

295 expression but unlike gland cells, the maximum number of trypsins expressed by any putative

296 neuron is three (Supplemental file 4). We show trypsin-expressing cells differentiating very early

297 in development, in the invaginating pharynx/mesenteries (Figure 3). Neurons expressing

298 RFamide and Elav are also undergoing terminal differentiation in this tissue at this

299 developmental stage (18, 22). Indeed, the trypsin protease NVJ_99932 (Figure 3) is co-expressed

300 with two other trypsins (NVJ_230861 and NVJ_130234) in a putative neuron expressing GABA

301 and dopamine receptors (Supplementary file 7). In vertebrates, secretion of neurotrypsin from the

14

302  pre-synaptic membrane facilitates degradation of the extracellular matrix during synaptic

303  plasticity and axon guidance (23). Although 17 different trypsins were expressed in putative

304  neurons, none of the trypsins from *N. vectensis* clustered with human neurotrypsin (Figure 5),

305  suggesting this function may have been acquired independently from different ancestral trypsins.

306

307  Trypsins are important regulators of tissue remodeling; as such, upregulation of trypsins and

308  other proteases occurs coincident with wound healing and tissue regeneration (24).  Recent

309  studies of regeneration in *N. vectensis* demonstrated that a new pharynx will regenerate from the

310  oral ends of the mesenteries after amputation (25) and that many proteases are expressed

311  abundantly during this process (26). Thus, the mesenteries appear to play an important role in

312  directing the tissue remodeling process in *N. vectensis*. In support of this, a study of wound

313  healing in response to a body wall injury demonstrated that the mesenteries come into direct

314  contact with damaged tissue during the healing process (27). This study also showed that two

315  trypsins (NVJ_107554 and NVJ_112683) are among the top genes undergoing upregulation

316  during wound healing in *N. vectensis*. While NVJ_112683 was not reported in the single-cell

317  dataset, NVJ_107554 is expressed in two putative gland cells (metacells C12 and C19,

318  Supplemental file 4). One of these (metacell C12) is also the site of expression of three proteases

319  examined by *in situ* hybridization in this study (Figure 3). Thus, mesentery-expressed trypsins

320  play important roles in the cell and tissue biology of *N. vectensis* during wound healing and

321  regeneration and these roles may vary through ontogeny.

322

323  Beyond their roles in digestion and tissue remodeling, trypsins are an important component of

324  the innate immune system. In vertebrates, immune trypsins play a role in blood coagulation and

325    are part of the complement system which recognizes foreign particles (28). In symbiotic

326    cnidarians, immune trypsins play a role in the beneficial interaction between the host and the

327    alga (29). While *N. vectensis* does not host symbiotic algae, a previous study aimed an

328    understanding the origin of the innate immune system reported the expression of three immune

329    system trypsins in *N. vectensis*: MASP (NVJ_138799) and two paralogs of Factor B

330    (NVJ_41116, NVJ_204186), each of which were expressed in the endoderm (gastrodermis) of

331    juvenile polyps (30). We found that the two factor B orthologs were also co-expressed in a single

332    putative gastrodermal cell (Supplemental file 4) further supporting a role for the endoderm in the

333    immune response of *N. vectensis*. One trypsin (NVJ_127465) was not reported in the single-cell

334    dataset (15) but was among the genes found to be significantly upregulated in the tissue-specific

335    transcriptome of nematosomes, which may also play a role in the immune system of *N. vectensis*

336    (11). This gene clustered with human chymotrypsin genes, not the immune system trypsins

337    (Figure 5); as such it may have acquired a role in the immune system secondarily.

338

339    *Trypsin functional diversity has undergone numerous expansions*

340    Two groups of trypsins (intracellular and non-catalytic) are found in all domains of life (Figure

341    4), both of which form monophyletic groups in animals (Figure 5). Trypsin diversity expanded

342    rapidly with the origin of animals, as representatives of two of the earliest diverging animal

343    lineages (ctenophores and placozoans) have at least 20 trypsins (Figure 4). The phylogeny of

344    animal trypsins suggests the last common ancestor of animals may have had at least six major

345    groups of trypsins (Figure 5). Sponges are unusual among animals in that they have only three

346    trypsins – two trypsin-PDZ paralogs and a trypsin-Sushi protein (Supplemental file 6). This

347    suggests either extensive loss of trypsins in Porifera or independent diversification of trypsins in

16

348  ctenophores and in the stem of parahoxozoa. The evolutionary history of trypsin-Sushi, trypsin-

349  SRCR, and trypsin-ShK proteins sheds little light on this topic; while all three associations are

350  found in ctenophores, trypsin-Sushi proteins are missing from placozoans and trypsin-SRCR and

351  trypsin-ShK proteins are missing from both sponges and placozoans (Figure 4,8, Supplemental

352  file 6). Presently, it is difficult to determine if these associations were either lost in sponges and

353  placozoans or arose independently in ctenophores and the stem of planulozoa. Considering the

354  association between the ShK and trypsin domains also seems to have been lost in the bilaterian

355  lineage (Figure 4, Supplemental file 6), we think independent gain of ShK domains in

356  ctenophores and cnidarians is likely.

357

358  The ancestral cnidarian may have had a far more diverse suite of trypsins than the ancestral

359  animal. Indeed, our data suggest there were at least 17 lineages of trypsins present in the last

360  common cnidarian ancestor (Figure 6) and 12 of the associations between trypsin and another

361  conserved domain in *N. vectensis* are specific to cnidarian lineages (Figure 8). Furthermore,

362  while the diversity of trypsins in *N. vectensis* rivals that of *H. sapiens* (Figure 4), there is little

363  conservation in the associated domains in these taxa (Figure 5B). Furthermore, within each clade

364  of trypsins, sequences from cnidarians and bilaterians form distinct groups, suggesting that

365  secretory cell function expanded independently in the cnidarian and bilaterian stem lineages.

366  There was extensive divergence in the trypsin gene superfamily during the diversification of

367  cnidarians but anthozoans seem to have undergone additional radiations in at least two trypsin

368  clades. Anthozoans are the most speciose group of cnidarians and are largely sessile; thus,

369  selection for trophic specialization and sympatric niche diversification may be stronger among

370  anthozoans than medusozoans.

17

371

372    *Numerous evolutionary mechanisms contributed to the rise of trypsin diversity*

373    Multidomain proteins are more common than proteins with only a single domain as domain

374    recombination increases versatility in protein function (31). Selection to maintain the catalytic

375    activity of the trypsin domain while allowing the context in which this domain is expressed to

376    vary was a primary driver of diversification of this gene family. Surprisingly, there was little

377    conservation in trypsin-associated domains across animals, even among cnidarians (Figure 4,

378    Supplemental file 6), suggesting the trypsin domain underwent significant duplication before the

379    diversification of cnidarians and that the associated domains have been continuously gained and

380    lost in each cnidarian lineage. Furthermore, nearly 40% (28/72) of the proteins comprising the *N.*

381    *vectensis* tryptome have only a trypsin domain (Figure 7) yet these trypsin-only proteins did not

382    form a monophyletic group (Figure 5,6). These results suggest that trypsin domains themselves

383    may be rapidly gained and lost from evolutionarily unrelated proteins, further underscoring the

384    selective advantage of a trypsin domain.

385

386    Trypsins were associated with 24 other conserved protein domains in the *N. vectensis* tryptome,

387    only few of which were over-represented in the tryptome (Figure 4A). The ShK domain is a

388    short peptide found in a K-channel inhibitor originally isolated from the sea anemone

389    *Stichodactyla helianthus* (32). The ShK phylogeny (Figure 7B) further suggests this domain is

390    gained and lost easily, as ShK domains from sister trypsins were almost never monophyletic.

391    Consistent with this, every ShK domain in the tryptome of *N. vectensis* was encoded by only a

392    single exon (Supplemental file 5), supporting the possibility of rapid evolution through exon

393    shuffling. Two trypsin-ShK proteins (NVJ_218669 and NVJ_218670) were found to be sister in

18

394    both phylogenies, suggesting they arose by duplication of the combined domains. These two

395    genes are encoded on the same scaffold and are separated by approximately 1000bp of genomic

396    DNA; thus, they are likely the result of a recent tandem duplication event. What role the ShK

397    domain plays when it is paired with the trypsin domain is not known but the overabundance of

398    these two combined domains in cnidarian tryptomes (Supplemental file 6) combined with the

399    multiple independent origins of this domain combination in *N. vectensis* (Figure 7B) suggest the

400    pairing provides a strong selective advantage in the biology of cnidarians.

401

402    We also found evidence of trypsin diversification independent of the acquisition of associated

403    domains. One gene (NVJ_127465) encodes three trypsin domains, all of which form a

404    monophyletic group suggesting this gene structure arose through tandem duplication of the

405    trypsin domain. Likewise, we identify four cases where sister trypsins are found on the same

406    scaffold, suggesting tandem gene duplication. Additionally, two trypsins were found to lack

407    introns (NVJ_128003 and NVJ_216003), suggesting these two arose through recent

408    retrotransposition. These two genes are also on the same scaffold, suggesting retrotransposition

409    may have been followed by tandem gene duplication. The tryptome from *H. sapiens* also

410    includes two proteins with three trypsin domains each (Supplemental file 6). All six of these

411    trypsin domains from *H. sapiens* are found in the tryptase/transmembrane clade (Supplemental

412    file 3), whereas the three domains in NVJ_127465 group with chymotrypsins. Despite their

413    similar domain architecture, triple-trypsin domain proteins appear to have evolved multiple

414    times.

415

19

416    Diversification of the trypsin superfamily through gene duplication and divergence has been

417    continuous, suggesting an important role for trypsins in the evolutionary success of animal

418    lineages. Early expansions, before the diversification of animals may have promoted the origins

419    of the innate immune system and extracellular digestion, facilitating the evolution of large body

420    size and increased longevity in this lineage. Lineage-specific expansions (between genera, for

421    example) in the digestive trypsins enable the development of specialized of taxon-specific

422    tryptomes, which can support niche specialization in otherwise similar taxa (Figure 4,6,8). The

423    expansion of the tryptome of *N. vectensis* was facilitated by gene duplication followed by the

424    acquisition of additional domains (Figure 8). In some cases, the acquisition of an associated

425    domain by a trypsin protein occurred soon after the origin of the acquired domain (e.g., trypsin-

426    SRCR), while in other cases this acquisition occurred long after the associated domain first

427    appeared (e.g., trypsin-astacin). Indeed, we found no relationship between the age of the domain

428    and the age of the association with trypsin (Figure 8B), further supporting the idea that trypsin

429    domain architectures diversify continuously and are not dependent on the origin of novel

430    domains.

431

432    New genes are thought to arise rapidly *via* duplication and divergence (33) but the proportion of

433    novel genes that become functionally integrated into signaling networks may be very small. The

434    factors that promote retention of new genes following duplication are not well understood but

435    may involve complimentary degenerative mutations promoting subfunctionalization in the

436    duplicates (34). Under this framework, sister trypsins that result from a recent duplication event

437    should be expressed in differing contexts, for example in different developmental stages or in

438    different cell types. Consistent with this, analysis of the single-cell expression of NVJ_218669

20

439    and NVJ_218670, which appear to have undergone recent duplication, suggests that

440    NVJ_218669 is expressed in a subset of the cells that express NVJ_218670 (Supplemental file

441    4). Only 51 of the 72 trypsins in *N. vectensis* were reported in the single-cell study (15), making

442    a detailed analysis of shared expression patterns in sister trypsins premature; however, a targeted

443    study of this protein superfamily in the future may reveal novel cis-regulatory regions, further

444    enlightening the processes involved in the diversification of this protein family.

445

446    *Secretory cells and the evolution of cnidarian body plans*

447    Resolving the embryological origin of cnidarian gland cells will be important for understanding

448    the evolution of life history in Cnidaria. If the anthozoan polyp body plan is ancestral to all

449    cnidarians (35), then the origin of strobilation (medusa formation) and its associated tissue

450    remodeling in the stem medusozoan may have necessitated the sacrifice of the internalized tissue

451    layers of the ancestral pharynx and mesenteries. In this case, the stem medusozoan may have

452    overcome this loss by shifting the development of their gland cell population to the endoderm

453    without sacrificing the selective advantage of secreting their products into the gastrovascular

454    cavity. In support of this hypothesis, gland cells in *Hydra* are known to undergo differentiation in

455    a location-specific manner, suggesting the identity of this cell lineage is highly sensitive to

456    positional cues from other cells in their environment (36). Furthermore, a recent study of single-

457    cell dynamics in *Hydra* demonstrated that gland cells acquire their identity in the endoderm only

458    after their precursor migrates out of the ectoderm and across the mesoglea (37). Both of these

459    studies point to the highly plastic nature of gland cell identity in *Hydra* but similar analyses in

460    more medusozoans are needed to understand the relationship between gland cell development

461    and cnidarian life history evolution.

21

462

**Conclusions**

463

464 The transition from unicellular to multicellular life was marked by many transitions that enabled

465 functional specialization. Unicellular taxa used trypsins for intracellular protein regulation but

466 the origin of the regulated secretion system created new opportunities for protease activity in

467 multiple tissue compartments. Secretion of molecules to the extracellular space enabled the

468 development of the nervous, endocrine, immune, and digestive systems, and permitted spatial

469 and temporal separation of multiple functions performed by a single cell. The diversification of

470 animals was associated with a large expansion of trypsins. Trypsins with transmembrane

471 domains first appear in the choanoflagellates but trypsins with signal peptides did not appear

472 until the origin of animals. Subsequent duplication and divergence (e.g., through exon shuffling

473 and retrotransposition) of genes encoding secreted proteases enabled nuanced variation in the

474 function of these secretory cells before the increase in anatomical diversity (Figure 9).

475

**Methods**

476

*Electron microscopy, cell proliferation assay, and* in situ *hybridization*

477

478 Adult polyps were immobilized for 10 mins in 7.5% $MgCl_2$ and processed for transmission

479 electron microscopy as described previously (38). Samples were imaged on a Hitachi HT7700 at

480 the University of Hawaii's Biological Electron Microscopy facility. To identify proliferating

481 nuclei, live polyps were incubated in 100uM EdU (in 1/3X seawater) for 30 mins at room

482 temperature. Animals were then immobilized and fixed briefly (1.5 min) at 25C in 4%

483 paraformaldehyde with 0.2% glutaraldehyde in phosphate buffered saline with 0.1% Tween-20

484 (PTw) followed by a long fixation (60 min) in 4% paraformaldehyde in PTw at 4C. Fixed tissues

22

485    were analyzed using the Click-IT EdU kit (#C10340, Invitrogen, USA) following the

486    manufacturer's protocol. Nuclei were counter stained in a 30 min incubation in DAPI at room

487    temperature and samples were imaged on a Zeiss 710 confocal microscope at the Whitney Lab

488    for Marine Bioscience. To characterize the localization of target genes, we performed *in situ*

489    hybridization following a standard protocol for *N. vectensis* (39).

490

491    *Protein domain analysis*

492    To identify trypsin-domain proteins from *N. vectensis* we first searched the JGI protein models

493    using the default settings with hmmsearch (HMMER 3.1b2; http://hmmer.org/) and two target

494    HMMs: Trypsin (PF00089) and Trypsin_2 (PF13365). This approach yielded 99 putative

495    trypsin-domain containing proteins (hereafter referred to as trypsins) with an E-value ≤ 1e-05

496    (40). Where multiple partial non-overlapping trypsin domains were identified from the same

497    protein, we assumed these represented one single contiguous domain (41). Based on a reciprocal

498    BLAST comparison with transcriptome data available publicly (11), we found 68/99 of the JGI

499    gene models coding for trypsin proteins were incomplete. We manually corrected these

500    sequences using the transcriptome data and used these corrected sequences for downstream

501    analyses. We then used the transcriptome data to search protein models for evidence of

502    pseudogenes (with premature stop codons) using the translation and alignment features in

503    Geneious v 7.1.8 (https://www.geneious.com) and manually examined models for duplicate

504    predictions using the JGI genome viewer. Based on these analyses, we removed 27 sequences,

505    resulting in a final set of 72 curated trypsin protein models (FASTA file available at:

506    https://github.com/josephryan/2019-Babonis_et_al_trypsins).

507

23

508    We examined the domain architecture of trypsin proteins from *N. vectensis* by searching for non-

509    Trypsin domains in the amino acid sequences using hmmscan (HMMER 3.1b2) and the complete

510    Pfam-A database (downloaded Oct 27, 2017). Hmmscan identifies regions of similarity between

511    protein queries and domain models (protein profiles) derived from numerous proteins within the

512    family from a range of animals (Bateman et al 2004). Following the protocol of Koch et al (40),

513    we ran hmmscan using the default parameters and report only those domains with an

514    independent (domain-specific) E-value ≤ 0.05 that were found in a protein containing a

515    significant Trypsin (or Trypsin_2) domain. Domains that overlapped by ≤ 20% were both

516    retained; when the overlap was >20% the domain with the lower E-value was retained. In

517    addition to domain analysis, we manually searched an alignment of the corrected set of trypsin

518    protein models from *N. vectensis* for the conserved residues that comprise the trypsin catalytic

519    triad (necessary for inferring protease activity): H-57, D-102, or S-195.  Finally, we searched the

520    corrected amino acid sequences for signal peptides and transmembrane domains using SignalP

521    v4.1 (42) and TMHMM v2.0 (43), respectively.

522

523    To characterize the origin of trypsin domain architecture, we hmmscan with the same approach

524    described above to identify and characterize trypsins from representatives across all domains of

525    life. We sampled three bilaterians (*Capitella teleta*, *Branchiostoma floridae*, *Homo sapiens*), ten

526    cnidarians (*N. vectensis*, *Edwardsiella lineata*, *Aiptasia pallida*, *Anthopleura elegantissima*,

527    *Acropora digitifera*, *Renilla renilla, Hydra magnipapillata, Calvadosia cruxmelitensis, Atolla*

528    *vanhoeffeni, Alatina alata*), three non-planulozoan animals (*Mnemiopsis leidyi*, *Amphimedon*

529    *queenslandica, Trichoplax adhaerens)*, five non-metazoan eukaryotes (*Dictyostelium discoidum,*

530    *Schizosaccharomyces pombe, Capsaspora owczarzaki, Monosiga brevicolis, Salpingoeca*

531   *rosetta*) and a combined database of representative archeaea and bacteria (*Candidatus aquiluna,*

532   *Candidatus nitrosopumilus, Candidatus pelagibacter, Glaciecola pallidula*, *Marinobacter*

533   *adhaerens,* a marine gamma proteobacterium, and a marine group I thaumarchaeote). Protein

534   models were predicted from transcriptome data previously for *N. vectensis, E. lineata, A. pallida,*

535   *A. elegantissima, A. alatina, A. vanhoeffeni*, and *P. carnea (11)*. Proteomes for *R. renilla* and *C.*

536   *cruxmelitensis* were predicted from the transcriptome data reported by Kayal et al., (35) using the

537   same methods. For all other taxa, protein models were downloaded directly (commands available

538   at: https://github.com/josephryan/2019-Babonis_et_al_trypsins).

539

540   *Phylotocol (phylogenetic transparency)*

541   All phylogenetic investigations were planned prior to running any analyses and all are reported

542   in this manuscript. In most cases, these analyses were outlined beforehand in a phylotocol (44)

543   that is posted on our GitHub site: https://github.com/josephryan/2019-Babonis_et_al_trypsins.

544   Any analyses performed prior to being added to our phylotocol were later added to the document

545   and justified.

546

547   *Phylogenetics*

548   To understand the diversification of animal trypsins, we built a phylogeny using predicted

549   proteins from *M. leidyi, A. queenslandica, T. adhaerens, N. vectensis, E. lineata, H.*

550   *magnipapillata, C. teleta, B. floridae,* and *H. sapiens*. First, we used a custom script to generate

551   alignments from these protein files using the Trypsin HMM (commands available at:

552   https://github.com/josephryan/2019-Babonis_et_al_trypsins). All trees were constructed using a

553   maximum likelihood framework with RAxML and IQ-TREE (45-47). We used the model finder

554   function with IQ-TREE (-m MF) to determine the best substitution model for the alignment and

555   then ran three approaches in parallel: RAxML with 25 parsimony starting trees, RAxML with 25

556   random starting trees, a single run with IQ-TREE (which, by default, uses a broad sampling of

557   initial starting trees). We selected the best tree by comparing the maximum likelihood scores of

558   all three approaches.

559

560   Using the Trypsin HMM we recovered 97% (70/72) of the curated trypsin proteins from *N.*

561   *vectensis*. The two remaining trypsin proteins (NVJ_23745 and NVJ_203589) were recovered

562   using the Trypsin_2 HMM. (Note: the Trypsin_2 HMM recovered only 89% (64/72) of the

563   curated trypsins.) To understand the evolutionary relationships of these two proteins to the rest of

564   the trypsin family, we generated another phylogeny using the same procedure as above and an

565   alignment built using the Trypsin_2 HMM. This best tree recovered using the Trypsin_2 HMM

566   is provided in Supplemental file 2. After inspecting both trees, we removed sequences from *B.*

567   *floridae* for ease of viewing and re-ran the full analyses. All tree files and alignment files are

568   available on our Github site (https://github.com/josephryan/2019-Babonis_et_al_trypsins).

569

570   To evaluate whether *N. vectensis* has undergone lineage-specific expansion of trypsins or if the

571   common ancestor of all cnidarians had an equally diverse trypsin protein repertoire, we built a

572   phylogeny of trypsin proteins from cnidarians only using a subset of the proteomes listed above.

573   Specifically, we used four species of anthozoans (*N. vectensis, E. lineata, R. renilla, A.*

574   *digitifera*) and four medusozoans (*H. magnipapillata, C. cruxmelitensis, A. vanhoeffeni, A.*

575   *alata*). We then pruned all non-*Nematostella* taxa from this tree using Phyutility v.2.2.6 (48) to

576   generate a tree for *N. vectensis* trypsins only. To examine the evolutionary history of ShK

577    domains from *N. vectensis*, we used hmmsearch with the ShK HMM and a custom script (as

578    above) to identify and align all ShK domains from the predicted proteome. We then used the

579    approach described above to produce a phylogeny of ShK domains.

580

### Acknowledgements

586

### Authors' contributions

588    Study design/concept: LSB, MQM, JFR; animal/tissue methods: LSB, CE; phylogenetics: LSB,

589    JFR; other analyses: LSB; writing: LSB; review and editing: MQM, JFR, CE. All authors read

590    and approved the final manuscript.

591

### Competing interests

593    The authors have no competing interests.

594

### Figures

596    **Figure 1.** Morphology of the pharynx and mesenteries. (A) Adult polyp; the pharynx/mesentery

597    transition is denoted by the dotted line. (B-D) Polyp at rest; the pharyngeal ectoderm (green) is

598    retracted inside the oral ectoderm (yellow). (E-G) Partial eversion of pharynx occurs during

599    capture/handling of prey (*Artemia* sp., indicated by *). (H-J) Ingested prey passes through the

27

600     mouth and pharynx and remains in contact with the ectoderm (white arrows) of the mesenteries

601     during digestion; colored arrow indicates endoderm of mesenteries (pigmentation from

602     consumption of *Artemia*). (K) Cnidocytes (black arrow) and gland cells (green arrow) are

603     restricted to the ectoderm of the mesenteries. C,D,F,G,K are DIC micrographs. D,G are false

604     colored. E,H are oral views, the remaining images are lateral views. Dotted lines in D,G,K

605     denote position of mesoglea. White arrowheads point to the mouth, black arrowheads denote

606     transition from external to internal ectoderm.

607

608     **Figure 2.** Fine structure of the pharyngeal ectoderm and ectodermal mesenteries. (A)

609     Diagrammatic representation of regions where thin sections were examined through the pharynx

610     (I) and at two positions in the mesenteries (II, III). Top – sagittal view of a primary polyp,

611     middle – cross-section at position I, bottom – partial cross sections at positions II (left) and III

612     (right). (B-F) TEMs of a region of the pharyngeal ectoderm (I) indicated by the box in A. (B)

613     Two zymogen gland cells are false colored for emphasis. (C) Cross sections of cilia emerging

614     from the pharyngeal ectoderm into the pharyngeal canal. (D) Mucus-secreting cell, false colored.

615     (E) Ten zymogen-type gland cells and two cnidocytes. (F) Electron-dense vesicles (black

616     arrowheads in E,F,N,O) in the apex of non-glandular ectodermal cells. (G) SEM of a mesentery

617     from position II in panel A. The ectodermal (EC) portion has two parts: the ciliated tract (black

618     line) and cnidoglandular tract (white line); the endodermal (EN) portion is false colored. (H-J)

619     TEMs of the cnidoglandular tract of a mesentery from position II in A. Sections correspond to

620     the position of the dotted line in G. (I) One zymogen type gland cell is false colored. (J) Some

621     zymogen vesicles have heterogeneous contents (black arrow). (K) 3D reconstruction of a

622     confocal z-stack through the ciliated tract of a mesentery from position II in A; pink – nuclei of

28

623    proliferating cells (EdU), blue – quiescent nuclei (DAPI). (L) TEM section through the ciliated

624    tract at the position indicated by the dotted line in K and by the box in H. Nuclei corresponding

625    to panel K are false colored pink. (M-O) TEMs of a mesentery lacking a ciliated tract, position

626    III in panel A. (M) Endoderm is false colored. (N) Zymogen gland cells are false colored. (O)

627    Ciliary rootlet (black arrow) in the apex of a zymogen gland cell. Scale bars: white line – 10 um,

628    black line – 500 nm, grey line (panels G, K) – 50 um. Panels B and H are composites of multiple

629    micrographs. White arrowheads indicate cnidocytes throughout.

630

631    **Figure 3.** Expression of proteases in *N. vectensis*. (A-C) *In situ* hybridization of proteases in the

632    (A) M family – metallopeptidases, astacins, and carboxypeptidases and (B) the S1 family –

633    trypsins during early development. The oral pole is to the left in all images; black arrows indicate

634    expression in the ectoderm of the pharynx. (C) Double fluorescent *in situ* hybridization showing

635    the co-expression (yellow arrow) of two metalloprotease genes (Nv_88668 and Nv_2109) in the

636    aboral ectoderm of the pharynx and independent expression of both genes in individual cells of

637    the ectodermal mesentery (red and green arrows). Nuclei are labeled with DAPI. (D) Single-cell

638    expression of genes from panels A and B in metacells identified by (15). Metacells C9-C12 are

639    part of the "gland" cell cluster and C26 is part of the "neuron" cell cluster. Two genes

640    (NVJ_200868 and NVJ_2109) are clearly expressed in the ectodermal mesenteries but were not

641    reported in the single-cell study (indicated by dotted lines). (E) Histogram of the number of

642    expressed trypsins per cell (from (15)); red font highlights a small second mode in the

643    distribution. (F) Scatterplot of the number of cells exhibiting co-expression of multiple trypsins

644    as a function of the number of total cells in which a particular trypsin is expressed.

645

646    **Figure 4.** Characterization of the 72 proteins comprising the *N. vectensis* tryptome. (A)

647    Abundance of the 24 trypsin-associated domains in the tryptome (black) and the remaining

648    proteins in the predicted proteome (white). Numbers inside the bars reflect domain counts in

649    both datasets; grey numbers to the right of the bars indicate total proteins containing the

650    indicated domain. Pfam domains are listed on the left; see Supplemental file 4 for pfam domains

651    IDs. (B) Phylogenetic distribution of trypsin domain architectures. Domain architecture models

652    from the *N. vectensis* tryptome are shown on the right; proteins that differ only in the number of

653    copies of associated domains were collapsed into a single row and these protein IDs are

654    underlined. The number of proteins with identical domain composition are indicated for each

655    taxon in the table. Actinarians (sea anemones) are highlighted in blue. [1]Maximum number of

656    trypsins from any one taxon in the bacteria/archaea database. [2]El,Ae,Ap,Av have trypsin-TSP_1-

657    ShK proteins which were included in this count. Bac – bacteria/archaea, Dd – *Dictyostelium*

658    *discoidum*, Sp – *Schizosaccharomyces pombe*, Co – *Capsaspora owczarzaki*, Sr – *Salpingoeca*

659    *rosetta*, Mb – *Monosiga brevicollis*, Ml – *Mnemiopsis leidyi*, Aq – *Amphimedon queenslandica*,

660    Ta – *Trichoplax adhaerens*, El – *Edwardsiella lineata*,  Ap – *Aiptasia pallida*, Ae – *Anthopleura*

661    *elegantissima*, Ad – *Acropora digitifera*, Rr – *Renilla renilla*, Hm – *Hydra magnipapillata*, Pc –

662    *Podocoryna carnea*, Av – *Atolla vanhoeffeni*, Aa – *Alatina alata*, Cc – *Calvadosia*

663    *cruxmelitensis*, Ct – *Capitella teleta*, Hs – *Homo sapiens*

664

665    **Figure 5.** Evolutionary history of animal trypsins. (A) Phylogeny of animal trypsins; clades are

666    named after the orthologs from *H. sapiens*. Domain architectures for *N. vectensis* proteins are

667    shown. Proteins with multiple trypsin domains are polyphyletic; in such cases, the domain

668    architecture diagram points to a single trypsin domain and the position of the other trypsin

30

669    domains is indicated by open/closed circles. Open stars indicate clades that descended from

670    single genes present in the last common ancestor of animals. N indicates the position of human

671    neurotrypsin (NP_003610.2). The trypsin-death protein groups with trypsin-PDZ proteins but is

672    not found on this tree. *N. vectensis* proteins characterized by *in situ* hybridization are indicated:

673    arrows (this study), † (30), or * (12).  (B) Comparison of the domain content of trypsins from *N.*

674    *vectensis* and *H. sapiens*. The non-catalytic trypsins do not have associated domains and are not

675    included. The tryptase/transmembrane trypsins are not represented in *N. vectensis* and are also

676    not included. The "pancreatic" group includes granzyme, kallikrein, HGF, and elastase.

677

678    **Figure 6.** Gene phylogeny of cnidarian trypsin domains. Domain architectures for *N. vectensis*

679    proteins are shown; domains are colored as shown in Figure 4A. Proteins with multiple trypsin

680    domains are polyphyletic; in such cases, the domain architecture diagram points to a single

681    trypsin domain and the position of the other trypsin domains is indicated by open/closed circles.

682    Open stars indicate clades that descended from single genes present in the last common ancestor

683    of cnidarians; closed stars represent possible anthozoan expansions. Clades are labeled as in

684    Figure 5A; unshaded clades do not include any taxa from Figure 5 and cannot be resolved. The

685    "pancreatic" group includes proteins that are sister to human granzyme, kallikrein, HGF, and

686    elastase. Chymotrypsins are not monophyletic on this tree but are shaded together to facilitate

687    comparison with Figure 5A.

688

689    **Figure 7**. Relationships among trypsin and ShK proteins in *N. vectensis*. (A) Phylogeny of *N.*

690    *vectensis* trypsin domains; domain architecture diagrams are shown to the right. Proteins with

691    multiple trypsin domains appear multiple times in the tree; in these cases, the domain represented

31

692    at each position in the tree is bolded. SP – signal peptide, TM - transmembrane domain. Trypsin

693    domains lacking H-57, D-102, or S-195 are designated non-catalytic (see Supplemental file 4).

694    Hypothesized ancestral states (models with dotted lines) are indicated in two positions for

695    trypsin-ShK proteins; a revised hypothesis is indicated by the trypsin-ShK model with solid

696    lines. (B) Phylogeny of *N. vectensis* ShK domains from the predicted proteome. Proteins with

697    multiple ShK domains appear multiple times in the tree; in these cases, the domain represented at

698    each position is indicated with an arrowhead. Proteins that have both trypsin and ShK domains

699    are colored green. Colored symbols to the right of the domain architectures are provided to

700    facilitate comparisons of identical proteins in panels A and B. Domains are colored as shown in

701    Figure 4A.

702

703    **Figure 8.** Evolutionary history of the tryptome. (A) Domain architectures at ancestral nodes

704    indicate the origin of each domain combination; domains are reported at nodes where they

705    appear for the first time with an independent (domain-specific) E-value ≤ 0.05. The origin of the

706    SRCR domain and the origin of its association with trypsin are indicated by white arrows; black

707    arrows indicate the origin of the astacin domain and the association between astacin and trypsin.

708    E – Edwardsiidae only (*Edwardsiella + Nematostella*), N – *Nematostella* only. WAP domains (*)

709    either evolved twice (in *D. discoideum* and the ancestor of animals) or this domain was lost in

710    the intervening lineages. The DIM domain (†) has a complex distribution in bacteria, fungi,

711    cnidarians, and insects. (B) There is no relationship between the phlyostratigraphic age of the

712    domain and the age of its association with trypsin.

713

714 **Figure 9.** Secretory vesicles permit functional expansion without anatomical variation. Blue and

715 white cells reflect the intracellular expression of blue and white gene products. The origin of the

716 signal peptide directing proteins to the regulated secretion system/secretory vesicle (white arrow)

717 permitted segregation of gene products into two distinct compartments: intracellular and

718 intravesicular. The subsequent duplication of only few genes (black arrow) could result in the

719 acquisition of numerous new cell types through unique and combinatorial expression. Green –

720 co-expression of blue and yellow, purple – co-expression of blue and red, brown – co-expression

721 of blue, yellow, and red.

722

723 **Supplemental Material**

724 **Supplemental file 1:** Gland cells of the external ectoderm. (A-C) Mucus cells (false colored

725 yellow) in the tentacle ectoderm. (D) Zymogen cell (false colored green) in the body wall

726 ectoderm. White arrows point to cnidocytes, black arrowheads point to electron dense apical

727 vesicles in cells adjacent to gland cells. A, B – SEM, C,D – TEM. Scale bars: black – 5um, white

728 – 10um.

729

730 **Supplemental file 2**: Phylogeny of Trypsin_2 domains across animals. Protein models for *N.*

731 *vectensis* are shown. NVJ_203589 and NVJ_23745 were not detected by the trypsin HMM and

732 do not appear in Figure 5. Colors as in Figure 4A.

733

734 **Supplemental file 3:** Human trypsin domain architecture mapped on the animal trypsin

735 phylogeny. Proteins with multiple trypsin domains are polyphyletic; in such cases, the diagram

736 points to a single trypsin domain and the position of the other trypsin domains is indicated by

33

737    symbols. *N. vectensis* proteins characterized by *in situ* hybridization are indicated by arrows (this

738    study), † (30), or * (12).

739

740    **Supplemental file 4**: Excel file tabulating: presence/absence of signal peptides and

741    transmembrane domains in *N. vectensis* trypsins, amino acid sequences for trypsin catalytic

742    domains for all taxa, pfam IDs for all domains, and a summary of single-cell expression of

743    trypsins published previously (15).

744

745    **Supplemental file 5:** All trypsin domains are encoded by multiple exons in *N. vectensis*

746    (excluding the NVJ_128003 and NVJ_216003) but many of the associated domains are encoded

747    by a single exon (indicated by triangle). Domains that span intron/exon boundaries by 10 or

748    fewer nucleotides were considered to be encoded by a single exon.

749

750    **Supplemental file 6:** Trypsin protein domain architectures from all taxa examined in this study.

751

752    **Supplemental file 7:** Trypsin protein IDs from all taxa examined in this study.

753

754

755    **References**

756    1.    Magie CR, Daly M, Martindale MQ. Gastrulation in the cnidarian *Nematostella vectensis*

757    occurs via invagination not ingression. Dev Biol. 2007;305(2):483-97.

758    2.    Fritzenwanker JH, Genikhovich G, Kraus Y, Technau U. Early development and axis

759    specification in the sea anemone *Nematostella vectensis*. Dev Biol. 2007;310(2):264-79.

760  3.    Finnerty JR, Pang K, Burton P, Paulson D, Martindale MQ. Origins of bilateral symmetry:

761      Hox and dpp expression in a sea anemone. Science. 2004;304(5675):1335-7.

762  4.    Kusserow A, Pang K, Sturm C, Hrouda M, Lentfer J, Schmidt HA, et al. Unexpected

763      complexity of the Wnt gene family in a sea anemone. Nature. 2005;433(7022):156-60.

764  5.    Wijesena N, Simmons DK, Martindale MQ. Antagonistic BMP-cWNT signaling in the

765      cnidarian *Nematostella vectensis* reveals insight into the evolution of mesoderm. P Natl Acad Sci

766      USA. 2017;114(28):E5608-E15.

767  6.    Leclere L, Rentzsch F. RGM Regulates BMP-Mediated Secondary Axis Formation in the

768      Sea Anemone *Nematostella vectensis*. Cell Rep. 2014;9(5):1921-30.

769  7.    Moiseeva E, Rabinowitz C, Paz G, Rinkevich B. Histological study on maturation,

770      fertilization and the state of gonadal region following spawning in the model sea anemone,

771      *Nematostella vectensis*. Plos One. 2017;12(8).

772  8.    Jahnel SM, Walzl M, Technau U. Development and epithelial organisation of muscle cells

773      in the sea anemone *Nematostella vectensis*. Front Zool. 2014;11.

774  9.    Renfer E, Amon-Hassenzahl A, Steinmetz PRH, Technau U. A muscle-specific transgenic

775      reporter line of the sea anemone, *Nematostella vectensis*. P Natl Acad Sci USA.

776      2010;107(1):104-8.

777  10.  Frank P, Bleakney J. Histology and sexual reproduction of the anemone *Nematostella

778      vectensis* Stephenson 1935. Journal of Natural History. 1976;10(4):441-9.

779  11.  Babonis LS, Martindale MQ, Ryan JF. Do novel genes drive morphological novelty? An

780      investigation of the nematosomes in the sea anemone *Nematostella vectensis*. Bmc Evol Biol.

781      2016;16.

782    12.  Steinmetz PRH, Aman A, Kraus JEM, Technau U. Gut-like ectodermal tissue in a sea

783    anemone challenges germ layer homology. Nat Ecol Evol. 2017;1(10):1535-42.

784    13.  Page MJ, Di Cera E. Serine peptidases: Classification, structure and function. Cell Mol Life

785    Sci. 2008;65(7-8):1220-36.

786    14.  Page MJ, Di Cera E. Evolution of Peptidase Diversity. J Biol Chem. 2008;283(44):30010-4.

787    15.  Sebe-Pedros A, Saudemont B, Chomsky E, Plessier F, Mailhe MP, Renno J, et al. Cnidarian

788    Cell Type Diversity and Regulation Revealed by Whole-Organism Single-Cell RNA-Seq. Cell.

789    2018;173(6):1520-+.

790    16.  Domazet-Loso T, Brajkovic J, Tautz D. A phylostratigraphy approach to uncover the

791    genomic history of major adaptations in metazoan lineages. Trends Genet. 2007;23(11):533-9.

792    17.  Long M, Betran E, Thornton K, Wang W. The origin of new genes: Glimpses from the

793    young and old. Nat Rev Genet. 2003;4(11):865-75.

794    18.  Marlow HQ, Srivastava M, Matus DQ, Rokhsar D, Martindale MQ. Anatomy and

795    Development of the Nervous System of *Nematostella vectensis*, an Anthozoan Cnidarian. Dev

796    Neurobiol. 2009;69(4):235-54.

797    19.  Putnam NH, Srivastava M, Hellsten U, Dirks B, Chapman J, Salamov A, et al. Sea anemone

798    genome reveals ancestral eumetazoan gene repertoire and genomic organization. Science.

799    2007;317(5834):86-94.

800    20.  Moran Y, Praher D, Schlesinger A, Ayalon A, Tal Y, Technau U. Analysis of Soluble

801    Protein Contents from the Nematocysts of a Model Sea Anemone Sheds Light on Venom

802    Evolution. Mar Biotechnol. 2013;15(3):329-39.

803    21.  Haynes JF, Davis LE. Ultrastructure of Zymogen Cells in Hydra Viridis. Z Zellforsch Mik

804    Ana. 1969;100(2):316-&.

805    22.  Nakanishi N, Renfer E, Technau U, Rentzsch F. Nervous systems of the sea anemone

806    *Nematostella vectensis* are generated by ectoderm and endoderm and shaped by distinct

807    mechanisms. Development. 2012;139(2):347-57.

808    23.  Almonte AG, Sweatt JD. Serine proteases, serine protease inhibitors, and protease-activated

809    receptors: Roles in synaptic function and behavior. Brain Res. 2011;1407:107-22.

810    24.  Daley WP, Peters SB, Larsen M. Extracellular matrix dynamics in development and

811    regenerative medicine. J Cell Sci. 2008;121(3):255-64.

812    25.  Amiel AR, Johnston HT, Nedoncelle K, Warner JF, Ferreira S, Rottinger E.

813    Characterization of Morphological and Cellular Events Underlying Oral Regeneration in the Sea

814    Anemone, *Nematostella vectensis*. Int J Mol Sci. 2015;16(12):28449-71.

815    26.  Schaffer AA, Bazarsky M, Levy K, Chalifa-Caspi V, Gat U. A transcriptional time-course

816    analysis of oral vs. aboral whole-body regeneration in the Sea anemone *Nematostella vectensis*.

817    Bmc Genomics. 2016;17.

818    27.  Dubuc TQ, Traylor-Knowles N, Martindale MQ. Initiating a regenerative response; cellular

819    and molecular features of wound healing in the cnidarian *Nematostella vectensis*. Bmc Biol.

820    2014;12.

821    28.  Heutinck KM, ten Berge IJM, Hack CE, Hamann J, Rowshani AT. Serine proteases of the

822    human immune system in health and disease. Mol Immunol. 2010;47(11-12):1943-55.

823    29.  Poole AZ, Kitchen SA, Weis VM. The Role of Complement in Cnidarian-Dinoflagellate

824    Symbiosis and Immune Challenge in the Sea Anemone Aiptasia pallida. Front Microbiol.

825    2016;7.

826   30.  Kimura A, Sakaguchi E, Nonaka M. Multi-component complement system of Cnidaria: C3,

827   Bf, and MASP genes expressed in the endodermal tissues of a sea anemone, *Nematostella*

828   *vectensis*. Immunobiology. 2009;214(3):165-78.

829   31.  Vogel C, Teichmann SA, Pereira-Leal J. The relationship between domain duplication and

830   recombination. J Mol Biol. 2005;346(1):355-65.

831   32.  Castaneda O, Harvey AL. Discovery and characterization of cnidarian peptide toxins that

832   affect neuronal potassium ion channels. Toxicon. 2009;54(8):1119-24.

833   33.  Tautz D, Domazet-Loso T. The evolutionary origin of orphan genes. Nat Rev Genet.

834   2011;12(10):692-702.

835   34.  Force A, Lynch M, Pickett FB, Amores A, Yan YL, Postlethwait J. Preservation of duplicate

836   genes by complementary, degenerative mutations. Genetics. 1999;151(4):1531-45.

837   35.  Kayal E, Bentlage B, Pankey MS, Ohdera AH, Medina M, Plachetzki DC, et al.

838   Phylogenomics provides a robust topology of the major cnidarian lineages and insights on the

839   origins of key organismal traits. Bmc Evol Biol. 2018;18.

840   36.  Siebert S, Anton-Erxleben F, Bosch TCG. Cell type complexity in the basal metazoan Hydra

841   is maintained by both fly stem cell based mechanisms and transdifferentiation. Dev Biol.

842   2008;313(1):13-24.

843   37.  Siebert S, Farrell JA, Cazet JF, Abeykoon YL, Primack AS, Schnitzler CE, et al. Stem cell

844   differentiation trajectories in Hydra resolved at single-cell resolution. bioRxiv. 2018:460154.

845   38.  Dubuc TQ, Dattoli AA, Babonis LS, Salinas-Saavedra M, Rottinger E, Martindale MQ, et

846   al. In vivo imaging of *Nematostella vectensis* embryogenesis and late development using

847   fluorescent probes. Bmc Cell Biol. 2014;15.

848    39.    Wolenski FS, Layden MJ, Martindale MQ, Gilmore TD, Finnerty JR. Characterizing the

849    spatiotemporal expression of RNAs and proteins in the starlet sea anemone, *Nematostella*

850    *vectensis*. Nat Protoc. 2013;8(5):900-15.

851    40.    Koch BJ, Ryan JF, Baxevanis AD. The Diversification of the LIM Superclass at the Base of

852    the Metazoa Increased Subcellular Complexity and Promoted Multicellular Specialization. Plos

853    One. 2012;7(3).

854    41.    Triant DA, Pearson WR. Most partial domains in proteins are alignment and annotation

855    artifacts. Genome Biol. 2015;16.

856    42.    Petersen TN, Brunak S, von Heijne G, Nielsen H. SignalP 4.0: discriminating signal

857    peptides from transmembrane regions. Nat Methods. 2011;8(10):785-6.

858    43.    Krogh A, Larsson B, von Heijne G, Sonnhammer ELL. Predicting transmembrane protein

859    topology with a hidden Markov model: Application to complete genomes. J Mol Biol.

860    2001;305(3):567-80.

861    44.    Ryan JF, Debiasse MB. Phylotocol: Promoting transparency and overcoming bias through

862    publicly posted, a priori methodological protocols in phylogenetics. Integr Comp Biol.

863    2018;58:E410-E.

864    45.    Kalyaanamoorthy S, Minh BQ, Wong TKF, von Haeseler A, Jermiin LS. ModelFinder: fast

865    model selection for accurate phylogenetic estimates. Nat Methods. 2017;14(6):587-+.

866    46.    Nguyen LT, Schmidt HA, von Haeseler A, Minh BQ. IQ-TREE: A Fast and Effective

867    Stochastic Algorithm for Estimating Maximum-Likelihood Phylogenies. Mol Biol Evol.

868    2015;32(1):268-74.

869    47.    Stamatakis A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large

870    phylogenies. Bioinformatics. 2014;30(9):1312-3.

871    48.  Smith SA, Dunn CW. Phyutility: a phyloinformatics tool for trees, alignments and molecular

872    data. Bioinformatics. 2008;24(5):715-6.
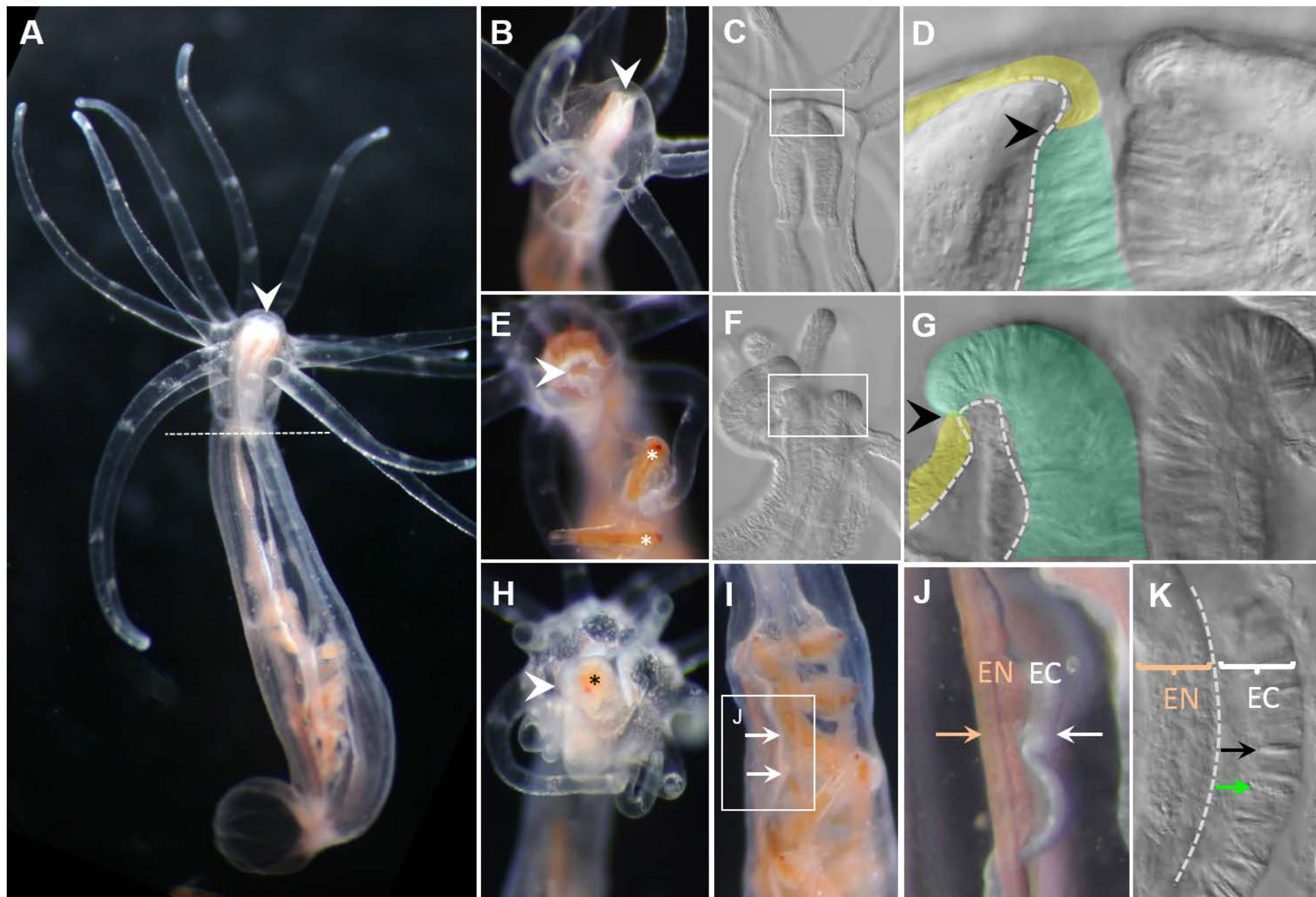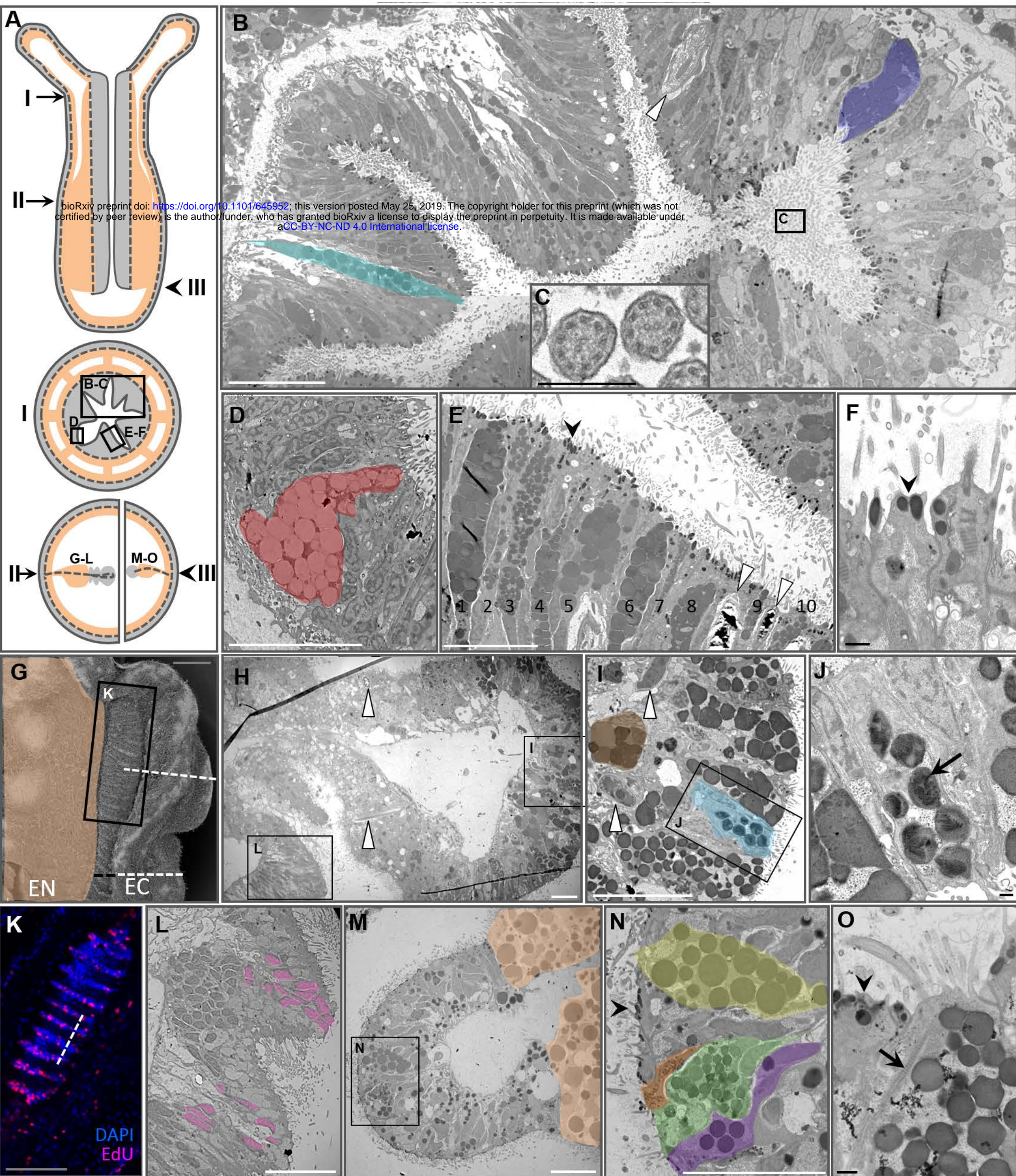
873

874
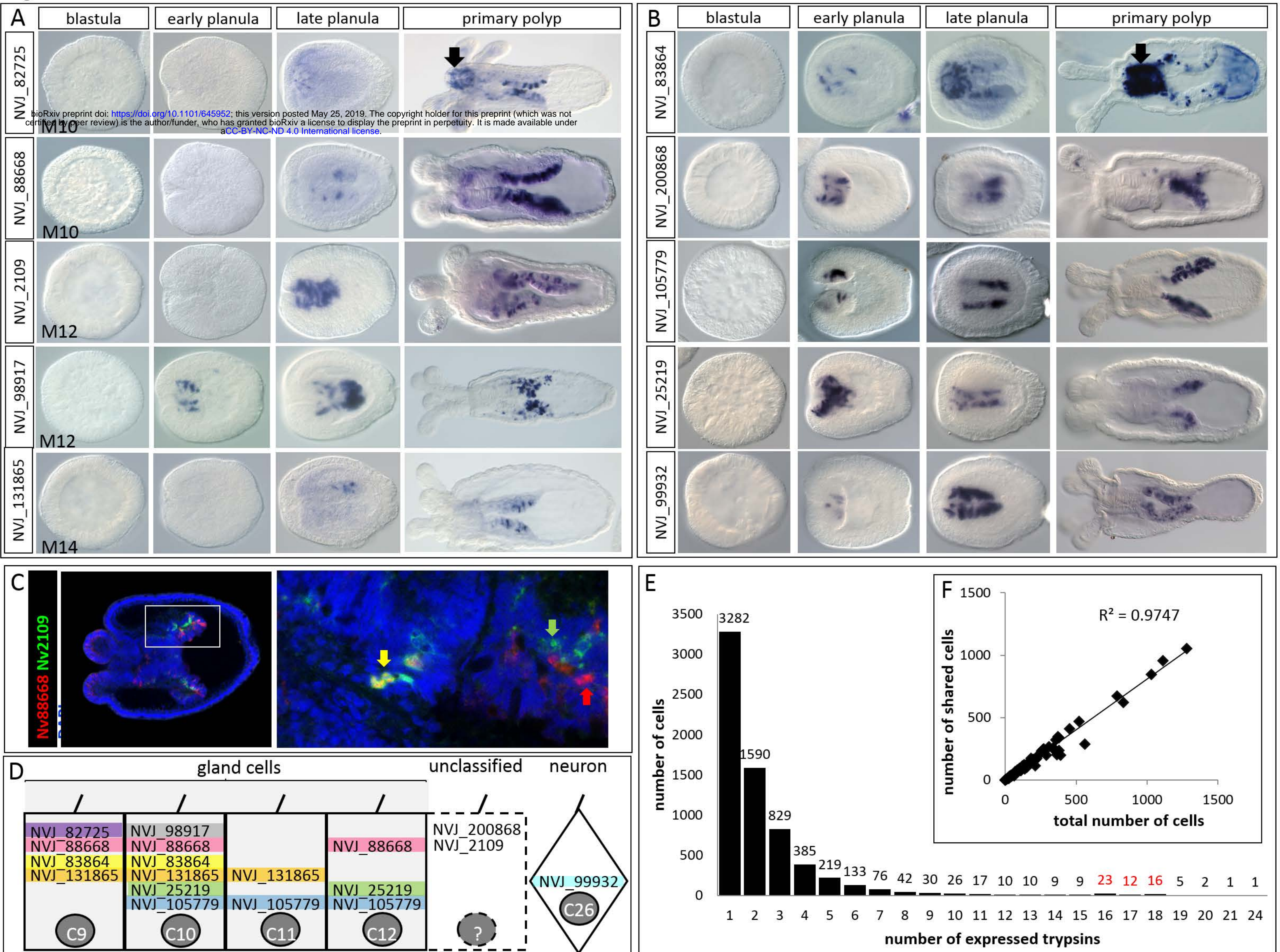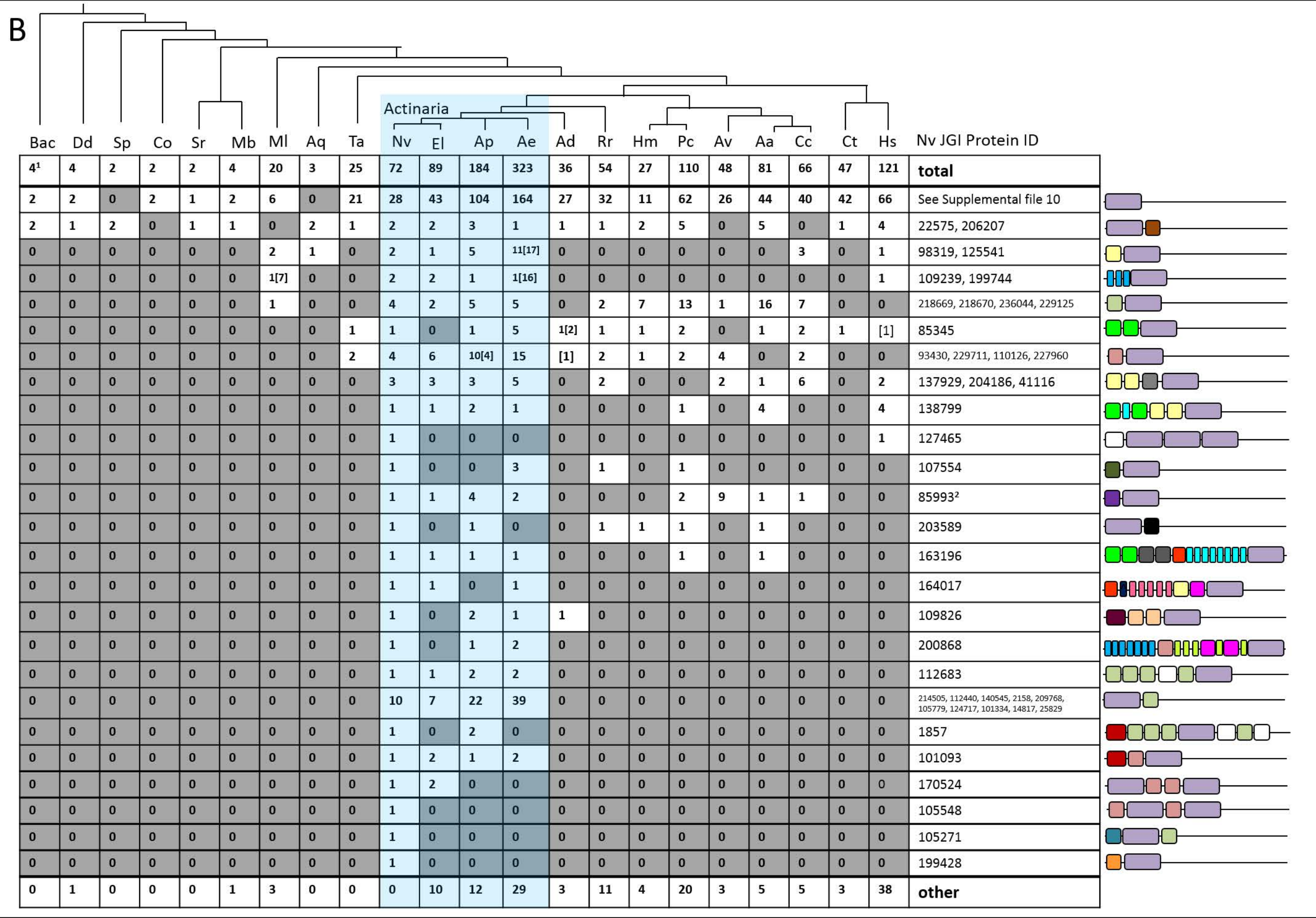
# Figure 1

# Figure 2

# Figure 3

Figure 4

## A

| Pfam Domain | in 'tryptome' | in predicted proteome | |
|---|---|---|---|
| Trypsin | 76 | 0 | 72 |
| DIM | 1 | 0 | 1 |
| ShK | 26 | 82 | 70 |
| Lustrin_cystein | 1 | 4 | 5 |
| Sushi | 21 | 112 | 39 |
| SRCR | 20 | 122 | 69 |
| MAM | 16 | 100 | 55 |
| Ldl_recept_b | 6 | 127 | 75 |
| Thyroglobulin | 4 | 89 | 23 |
| CUB | 2 | 49 | 20 |
| Death | 1 | 26 | 21 |
| WAP | 1 | 30 | 19 |
| Ig_2 | 10 | 402 | 186 |
| PDZ_2 | 2 | 84 | 69 |
| Astacin | 2 | 94 | 92 |
| WSC | 1 | 70 | 55 |
| PLAT | 1 | 71 | 49 |
| VWA | 3 | 214 | 125 |
| Ldl_recept_a | 2 | 149 | 36 |
| I-set | 5 | 405 | 187 |
| EGF_CA | 10 | 862 | 241 |
| F5_F8_typeC | 5 | 532 | 366 |
| FXa_inhib | 2 | 402 | 155 |
| TSP_1 | 1 | 279 | 122 |
| cEGF | 1 | 581 | 190 |

■ in 'tryptome' (N = 72)   □ in predicted proteome (N = 27,273)

## B

Actinaria spans Nv, El, Ap, Ae.

| Bac | Dd | Sp | Co | Sr | Mb | Ml | Aq | Ta | Nv | El | Ap | Ae | Ad | Rr | Hm | Pc | Av | Aa | Cc | Ct | Hs | Nv JGI Protein ID |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 4[1] | 4 | 2 | 2 | 2 | 4 | 20 | 3 | 25 | 72 | 89 | 184 | 323 | 36 | 54 | 27 | 110 | 48 | 81 | 66 | 47 | 121 | total |
| 2 | 2 | 0 | 2 | 1 | 2 | 6 | 0 | 21 | 28 | 43 | 104 | 164 | 27 | 32 | 11 | 62 | 26 | 44 | 40 | 42 | 66 | See Supplemental file 10 |
| 2 | 1 | 2 | 0 | 1 | 1 | 0 | 2 | 1 | 2 | 2 | 3 | 1 | 1 | 1 | 2 | 5 | 0 | 5 | 0 | 1 | 4 | 22575, 206207 |
| 0 | 0 | 0 | 0 | 0 | 0 | 2 | 1 | 0 | 2 | 1 | 5 | 11[17] | 0 | 0 | 0 | 0 | 0 | 0 | 3 | 0 | 1 | 98319, 125541 |
| 0 | 0 | 0 | 0 | 0 | 0 | 1[7] | 0 | 0 | 2 | 2 | 1 | 1[16] | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 109239, 199744 |
| 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 4 | 2 | 5 | 5 | 0 | 2 | 7 | 13 | 1 | 16 | 7 | 0 | 0 | 218669, 218670, 236044, 229125 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 1 | 5 | 1[2] | 1 | 1 | 2 | 0 | 1 | 2 | 1 | [1] | 85345 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 4 | 6 | 10[4] | 15 | [1] | 2 | 1 | 2 | 4 | 0 | 2 | 0 | 0 | 93430, 229711, 110126, 227960 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 3 | 3 | 3 | 5 | 0 | 2 | 0 | 0 | 2 | 1 | 6 | 0 | 2 | 137929, 204186, 41116 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 2 | 1 | 0 | 0 | 0 | 1 | 0 | 4 | 0 | 0 | 4 | 138799 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 127465 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 3 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 107554 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 4 | 2 | 0 | 0 | 0 | 2 | 9 | 1 | 1 | 0 | 0 | 85993² |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 203589 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 163196 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 164017 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 2 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 109826 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 200868 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 2 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 112683 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 10 | 7 | 22 | 39 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 214505, 112440, 140545, 2158, 209768, 105779, 124717, 101334, 14817, 25829 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1857 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 2 | 1 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 101093 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 170524 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 105548 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 105271 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 199428 |
| 0 | 1 | 0 | 0 | 0 | 1 | 3 | 0 | 0 | 0 | 10 | 12 | 29 | 3 | 11 | 4 | 20 | 3 | 5 | 5 | 3 | 38 | other |

# Figure 5

Figure 6

**Anthozoa**
*N. vectensis*
*E. lineata*
*R. renilla*
*A. digitifera*

**Medusozoa**
*H. magnipapillata*
*C. cruxmelitensis*
*A. vanhoeffeni*
*A. alatina*

chymotrypsinogen

non-catalytic

intracellular

immune/coagulation

"pancreatic"

Figure 7

# Figure 8

# Figure 9