**Supplementary Material online**
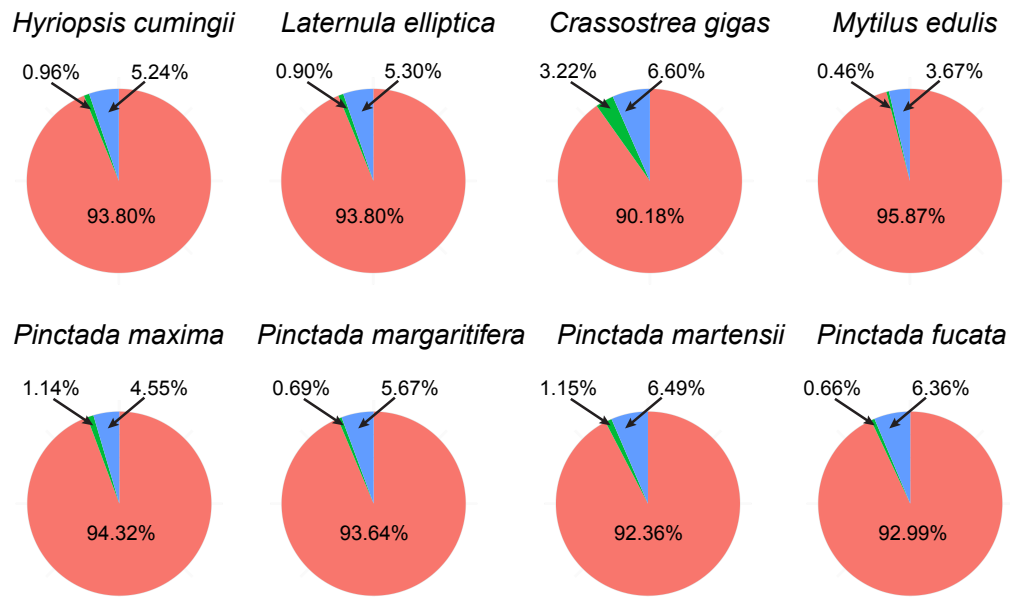
**Co-option and *de novo* gene evolution underlie molluscan shell diviersity**

Felipe Aguilera[1,2], Carmel McDougall[1] and Bernard M. Degnan[1*]

[1]Centre for Marine Sciences, School of Biological Sciences, The University of Queensland, Brisbane 4072, Australia
[2]Current address: Sars International Centre for Marine Molecular Biology, University of Bergen, Thormøhlensgate 55, Bergen 5008, Norway
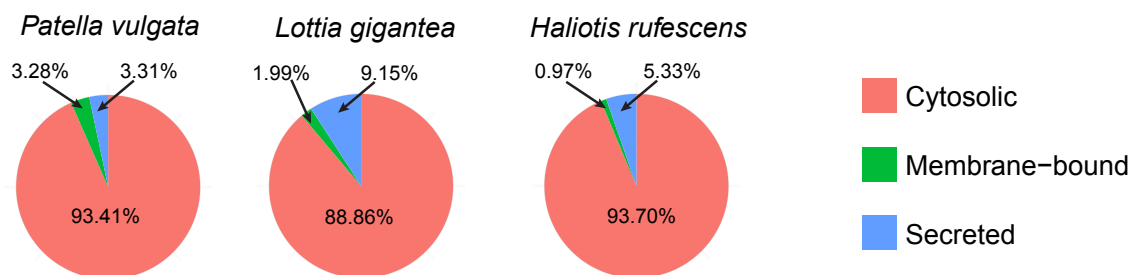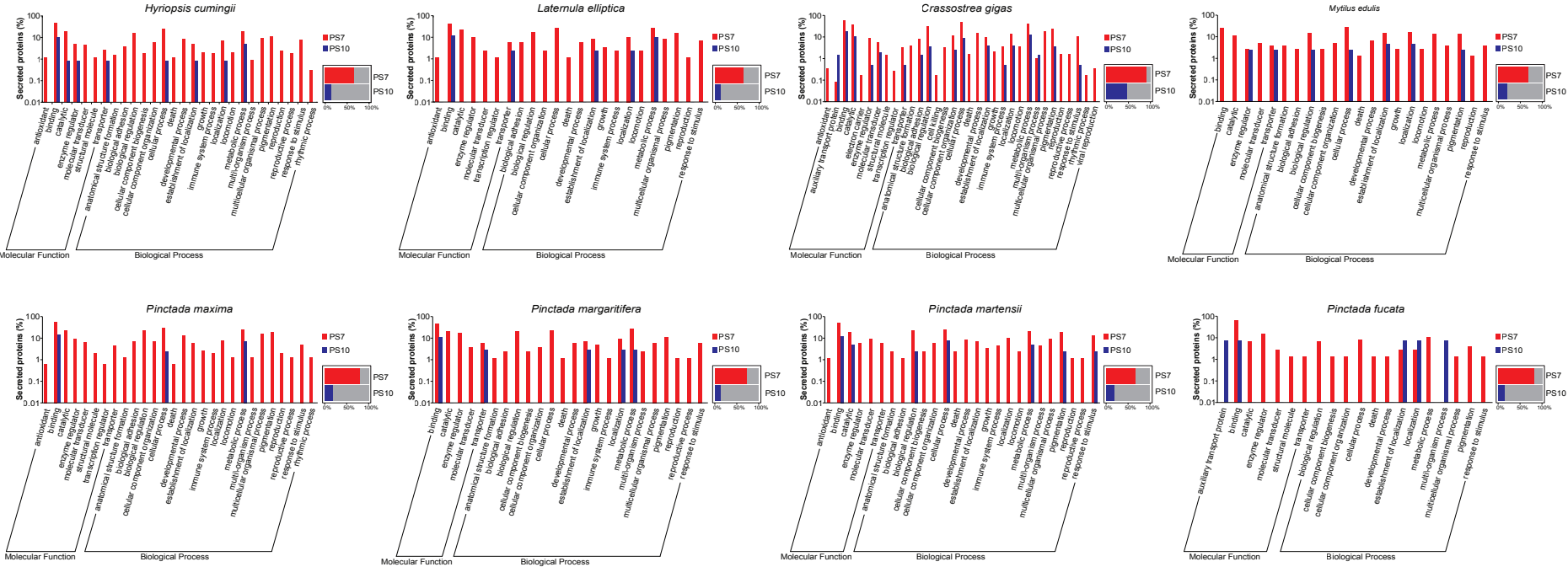
**Bivalves**



Figure S1. **Computational predictions of cytosolic, membrane-bound and secreted proteins for each species.** Pie charts represent the percentage of proteins predicted in each mantle transcriptome. Red color indicates cytosolic proteins; green color indicates membrane-bound proteins; and blue color indicates secreted proteins.

## BIVALVES



*Hyriopsis cumingii*

*Laternula elliptica*

*Crassostrea gigas*

*Mytilus edulis*

*Pinctada maxima*

*Pinctada margaritifera*

*Pinctada martensii*

*Pinctada fucata*

## GASTROPODS

*Patella vulgata*

*Lottia gigantea*
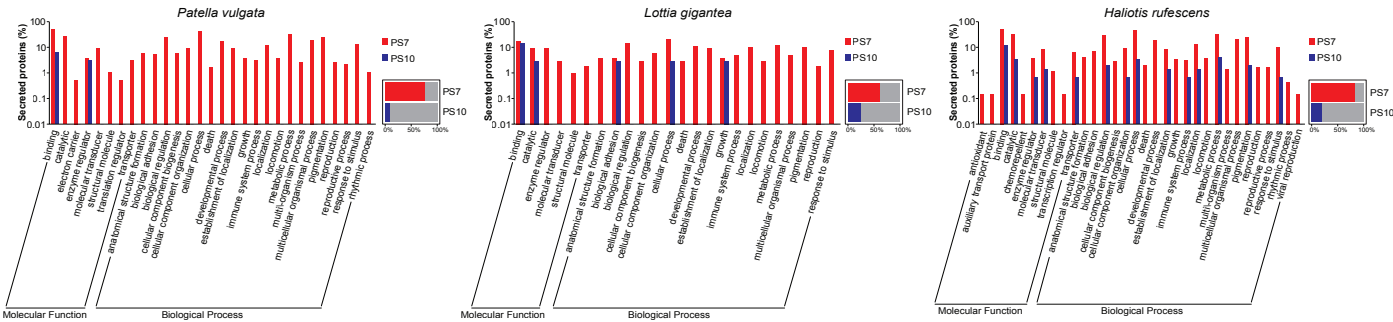
*Haliotis rufescens*

**Figure S2. GO annotations of genes that arose along the stems leading to bilaterian (PS7, red) and molluscan (PS10, blue) last common ancestors.** Horizontal bars - in small squares - represent the percentage of genes that were annotated (red and blue colors) and non-annotated (gray color) in each species.

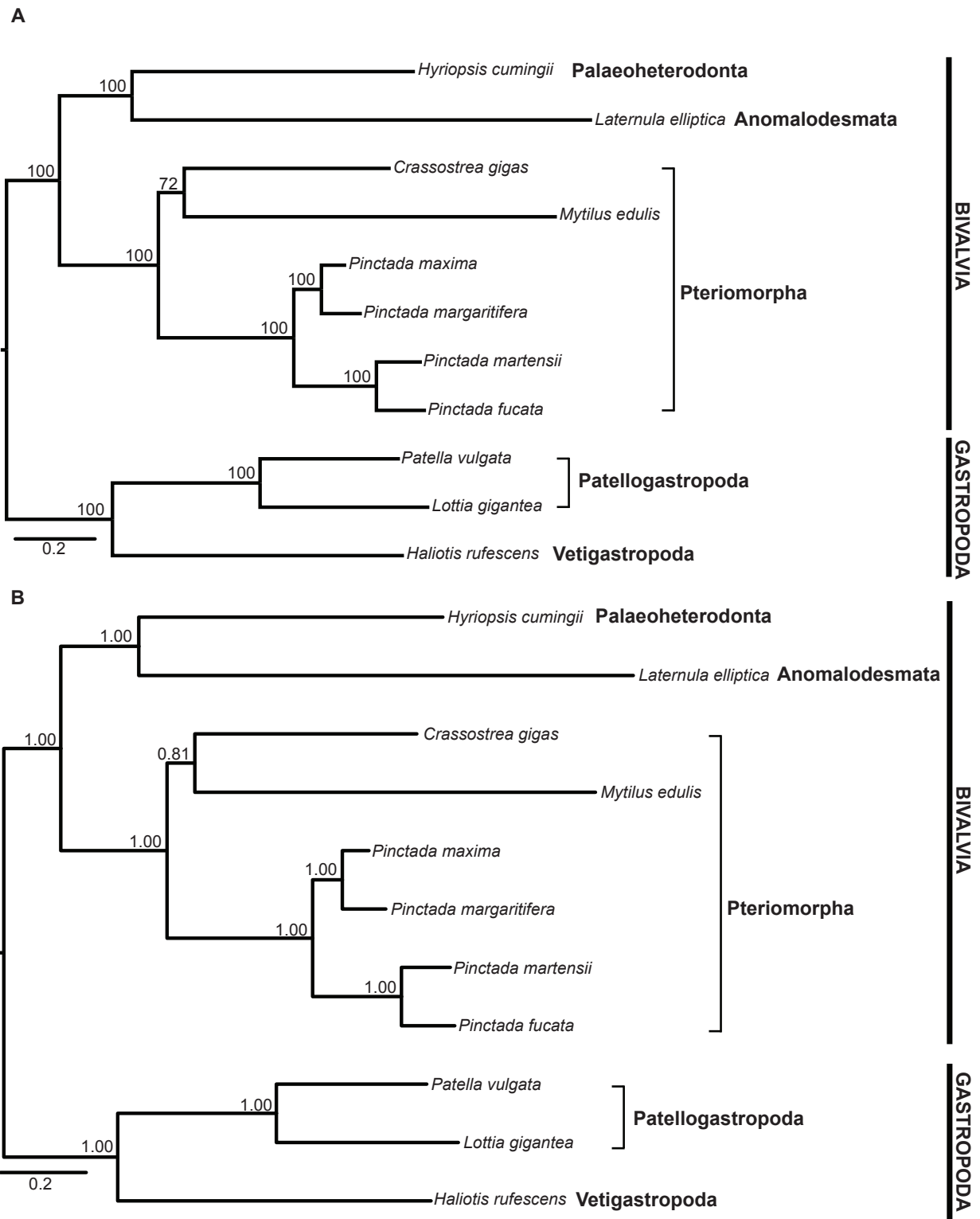**Figure S3. Phylogenetic relationships among bivalve and gastropod lineages based on 122 gene families that encode secreted proteins. A.** Phylogenetic tree based on Maximum Likelihood is shown. The most-likely tree (log likelihood = -138,848.94) sampled in RAxML using the PROTGAMMAWAGF substitution model. ML bootstrap (BS) support values are listed at each node. **B.** A 50% majority-rule consensus Bayesian tree is shown.

Posterior probabilities (PP) are indicated at each node. The length of the matrix is 13,604 amino acids.



**Figure S4. GO term (biological process and molecular function) enrichment in newly gained secreted gene families across conchiferan evolution.** Only significantly enriched (P<0.05, Fisher's exact test) GO terms at least present in two phylogenetic branches are shown. GO terms (biological process and molecular function) categories of the over-represented secreted gene families at each evolutionary time point are shown at the right. For a comprehensive list of enriched GO terms across molluscan evolution, see Supplementary Table S4, Supplementary Material online.

**Figure S5. Expression profiles of genes encoding secreted proteins – the mantle secretome – of eight bivalve and two gastropod species.** The boxplots show the median and interquartile ranges of the distribution of gene expression across phylostrata for bivalves and gastropods. Phylostratum (PS) is described and illustrated as per Figure 1.

**Figure S6. Expression levels of co-opted, lineage-specific and species-specific genes encoding secreted proteins.** Boxplots display the distribution of relative gene expression of genes classified as co-opted, lineage-specific and species-specific for each eight bivalves and two gastropods. Lines in the boxplots represent the median and interquartile range. Dots correspond to outlier genes.

# Carbonic anhydrase domain-containing proteins

Turbo marmoratus (BAB91157)

Haliotis gigantea (BAH58349)
Haliotis tuberculata (AEL22200)
Haliotis tuberculata (AEL22201)
Haliotis gigantea (BAH58350)

1

0.93

0.74

NG (68.00X)

(NG)2DG (2.17X)

L. gigantea (Contig275) ★
L. gigantea (374110482)
L.gigantea (238082) ★
L. gigantea (239188) ★

0.61

0.82

GNG(D)2N(G)2(R)2 (2.00X)

(D)2(Y)2(D)3YSN (2.60X)

(R)5FNG (2.00X)

(R)5FNG (2.00X)

GRGNGDNRD (2.00X)

(D)4Y (2.00X)

(D)4Y (2.00X)

L. gigantea (374110483)* ★

(D)2(Y)2(D)3YSN (2.60X)

GNG(D)2N(G)2(R)2 (2.00X)

G(D)2K (2.75X)

C. gigas (CGI10028495)

C. gigas (Contig2602)
C. gigas (CGI10014170) ★
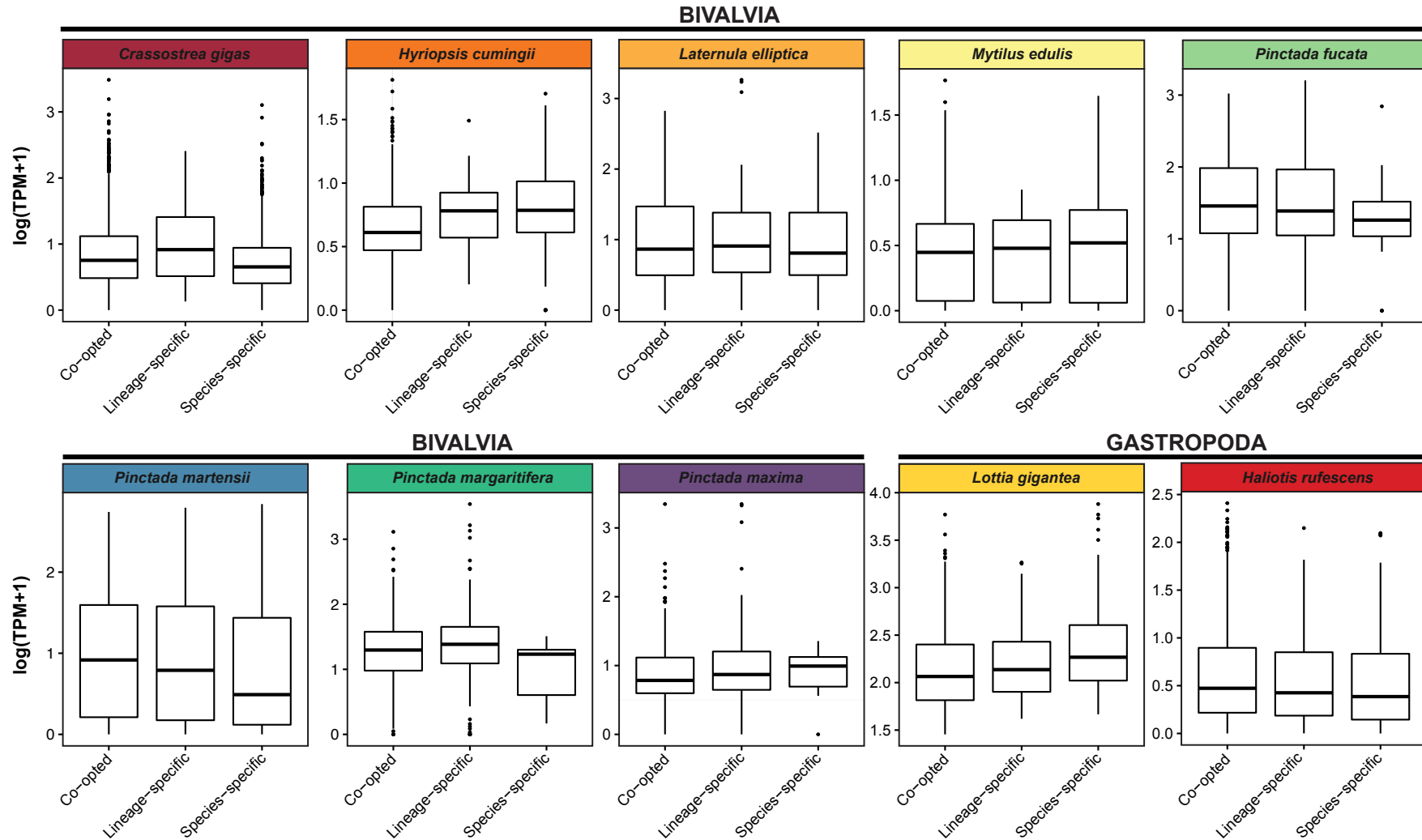C. gigas (comp46573_c0_seq1)*
H. cumingii (SRR530843_10829)*

(E)4DK(E)2NE (2.00X)

0.6

C. nemoralis (contig104312)* ★

NGDNGN (11.83X)

**Molluscan nacrein and nacrein-like cluster**

P. maxima (BAF42330)*
M. yessoensis (BAF42332)*
C. nippona (BAF42334)*
M. yessoensis (BAF42331)*
P. martensii (isotig00329)*
P. fucata (BAA11940)
P. martensii (isotig00326)*
C. nippona (BAF42333)*

0.76
0.54

1

0.56

NGD(NGN)2 (8.56X)

NGDNGN (10.83X)

NGDNGN (11.83X)

NGDNG(N)2GE (6.56X)

NGDNG(N)2GY (8.89X)

NGDNG(N)2GY (5.89X)

P. margaritifera (AEC03973)
P. margaritifera (AEC03971)
P. margaritifera (ADY69618)
P. margaritifera (AEC03972)
P. margaritifera (Contig542)*
P. margaritifera (PmargNacrein) ★
P. margaritifera (AEC03970)
P. maxima (BAA90540) ★
P. maxima (PmaxN66) ★
P. maxima (ACT55367)
P. maxima (Contig373)*

0.98
0.98 0.95
0.87
0.98

0.56

NGD(NGN)2 (8.56X)

(N)2GDNG(N)2(G)2 (4.00X)

NGN(NG)2(N)2(G)2N(NG)2N (9.00X)

(NGN)4NG (11.64X)

(N)2G(N)2(G)2YNG (15.14X)

D(NGN)3NG (2.50X)

(N)2G (65.33X)

(N)2G (65.33X)

(NGN)3GN (15.82X)

(NGN)3GN (15.82X)

NGDNG(N)2CDNG (2.09X)

(N)2GDNG(N)2CDNGN (2.00X)

S. ciliatum (lcl_24492S)
M. musculus (ENSMUSP00000103493)
G. gallus (ENSGALP00000004733)
M. musculus (ENSMUSP00000103490)
H. sapiens (ENSP00000390666)
H. sapiens (ENSP00000405388)
H. sapiens (ENSP00000285273)
H. sapiens (ENSP00000460238)
M. musculus (ENSMUSP00000035585)
M. musculus (ENSMUSP00000103495)
M. musculus (ENSMUSP00000003360)
H. sapiens (ENSP00000084798)

1

0.96

0.77

1

1

C. gigas (comp71517_c0_seq1)
H. rufescens (contig15309)
C. gigas (comp52387_c1_seq1)*
C. gigas (CGI10007458) ★

0.99
0.71

M. musculus (ENSMUSP00000121268)

0.97

M. musculus (ENSMUSP00000121268)
G. gallus (ENSGALP00000014565)

M. musculus (ENSMUSP00000099483)
H. sapiens (ENSP00000467465)
P. vulgata (PvulgataC11947)
G. gallus (ENSGALP00000009475)
M. musculus (ENSMUSP00000113400)

0.92
0.91
0.73
0.95

H. cumingii (SRR530843_38450)*

A. digitifera (adi_v1_16634, adi_v1_19067, adi_v1_1906600)
S. purpuratus (SPU_013459, SPU_103458)
S. purpuratus (SPU_001138, SPU_008894, SPU_022346)
G. gallus (ENSGALP00000003755)
M. musculus (ENSMUSP00000030817)
H. sapiens (ENSP00000366662, ENSP00000366654, ENSP00000447108, ENSP00000366661)
H. sapiens (ENSP00000435280)

1
1
0.52
0.69
0.88

S. pistillata (ACA53457)

H. cumingii (SRR530843_20497)*

S. purpuratus (SPU_012518)
P. maxima (comp_29068_c0_seq1)*
C. gigas (Contig7312)
C. gigas (comp70022_c1_seq1)
H. rufescens (contig13550)*
P. vulgata (PvulgataC6462)
C. gigas (comp79270_c0_seq1)*
P. maxima (comp_502_c0_seq1)
C. nemoralis (contig25891) ★
L. gigantea (66515)* ★
P. vulgata (PvulgataC16356)*
H. rufescens (contig4350)*
L. elliptica (comp118878_c0_seq1)*
H. rufescens (contig139881)*
H. rufescens (contig88255)*
C. gigas (comp52277_c1_seq1)*
C. gigas (comp51003_c0_seq1)*
C. gigas (Contig2947)
L. elliptica (SRR039932_115737)*

0.65
0.71
0.52
0.73
0.99
0.54
0.96
0.62
0.6
0.99

(Q)4AYP (12.57X)

QP(Q)3 (2.20X)

TVSEP (2.00X)

QS(Q)3HQA (2.00X)

H. sapiens (ENSP00000309649)
C. gigas (Contig1162)
L. gigantea (205401)* ★
H. cumingii (Contig20694)

0.98
0.96

S. pistillata (ACE95141)

L. elliptica (Contig3022)
H. cumingii (Contig20199)*
H. cumingii (Contig20059)*
H. cumingii (Contig22307)
H. cumingii (SRR530843_23095)*

0.66
0.98
1

S. ciliatum (lcl_16738, lcl43358, lcl14361, lcl32441)
S. ciliatum (lcl_48195, lcl36517, lcl44439)

1
0.68
1

1

0.57

0.58

1

0.3

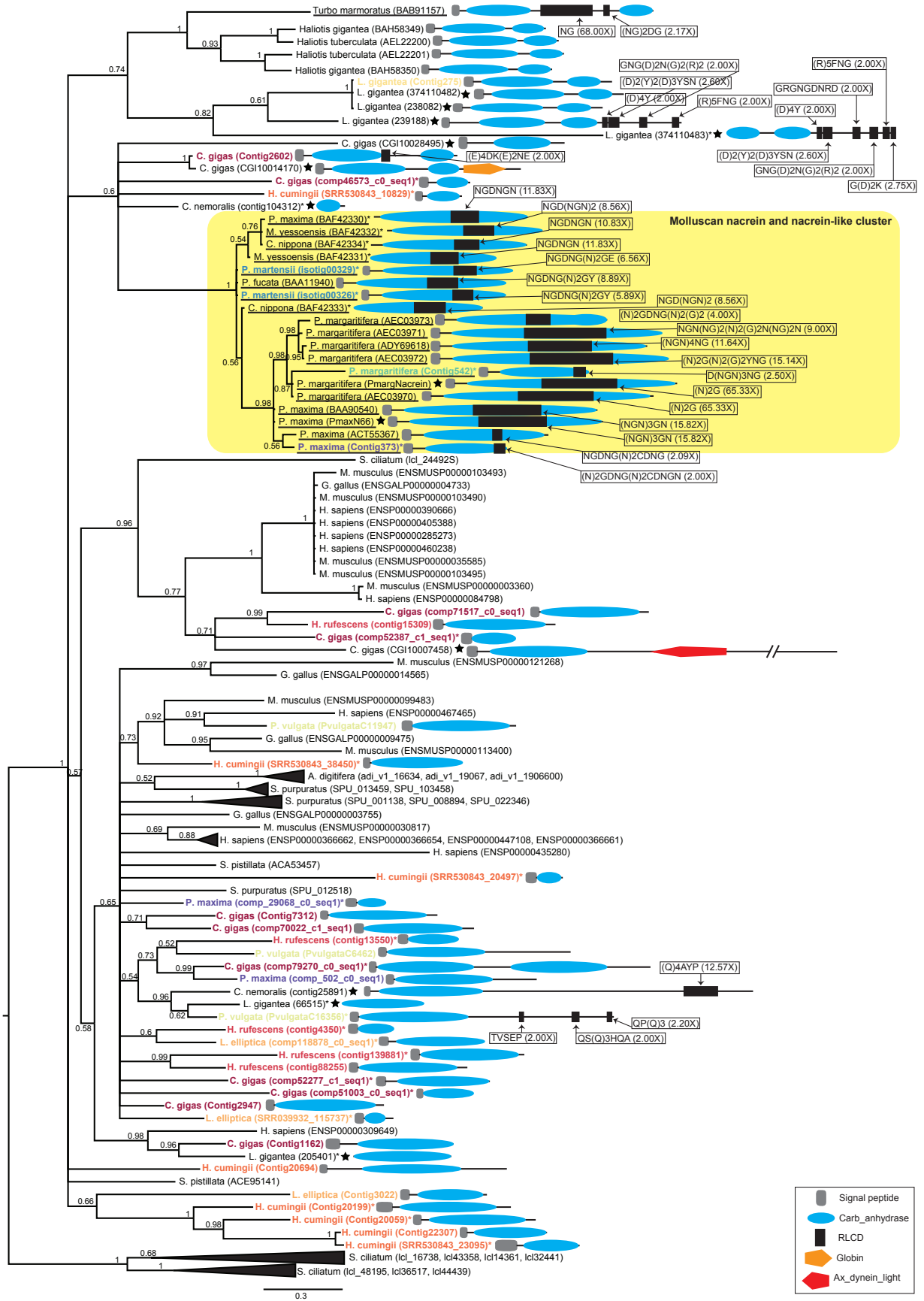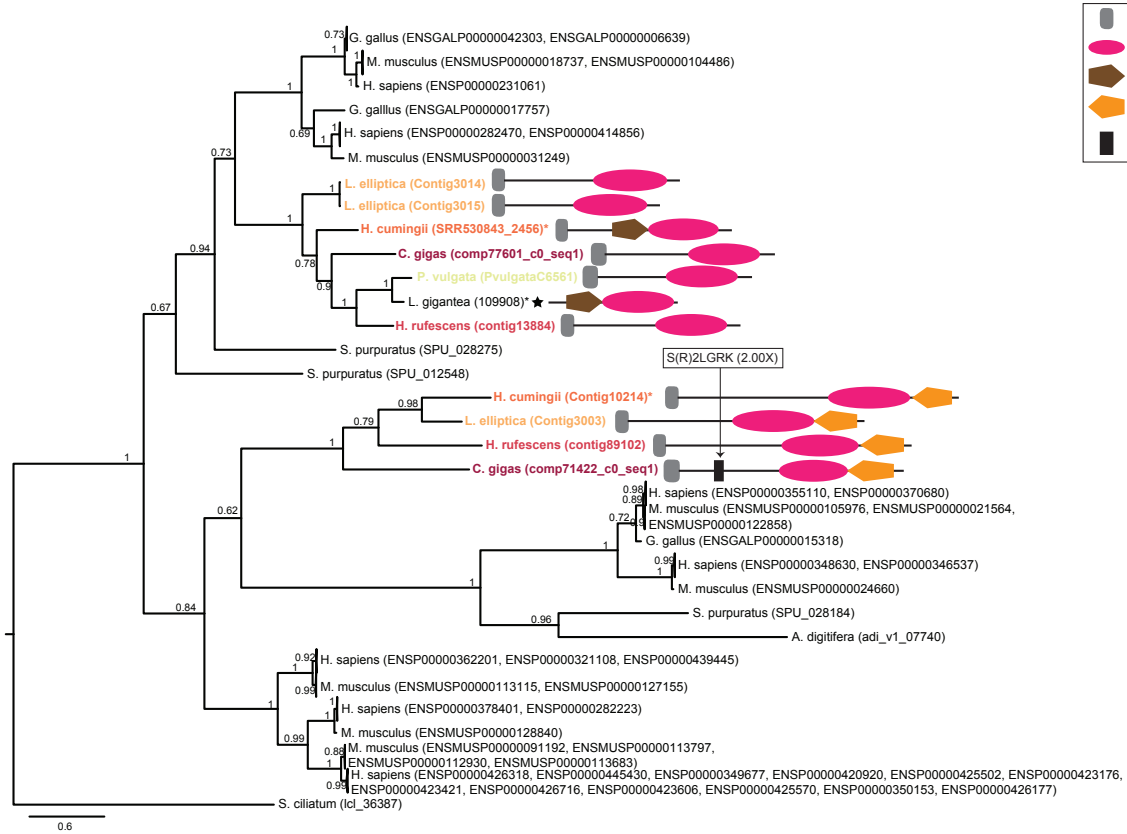| | Legend |
|---|---|
| ■ | Signal peptide |
| ⬭ | Carb_anhydrase |
| ■ | RLCD |
| ⬠ | Globin |
| ⬖ | Ax_dynein_light |

**Figure S7. Phylogenetic analysis of carbonic anhydrase domain-containing proteins.** Bayesian phylogenetic tree was obtained under the LG+G substitution model. The tree is rooted using the calcareous sponge sequences as outgroup. Statistical support for each node is indicated as posterior probabilities after 20,000,000 generations. Domain architectures of molluscan secreted proteins are shown on the right of each sequence, and domain nomenclature is shown in the rectangle. *Nacrein* and *nacrein-like* proteins are in bold/underline. Molluscan nacreins are highlighted in the yellow box. Proteins previously extracted from shells and/or known to be involved in shell formation are indicated by a black star. For each RLCD, the consensus repeat sequence is shown, followed by the copy number (in brackets).
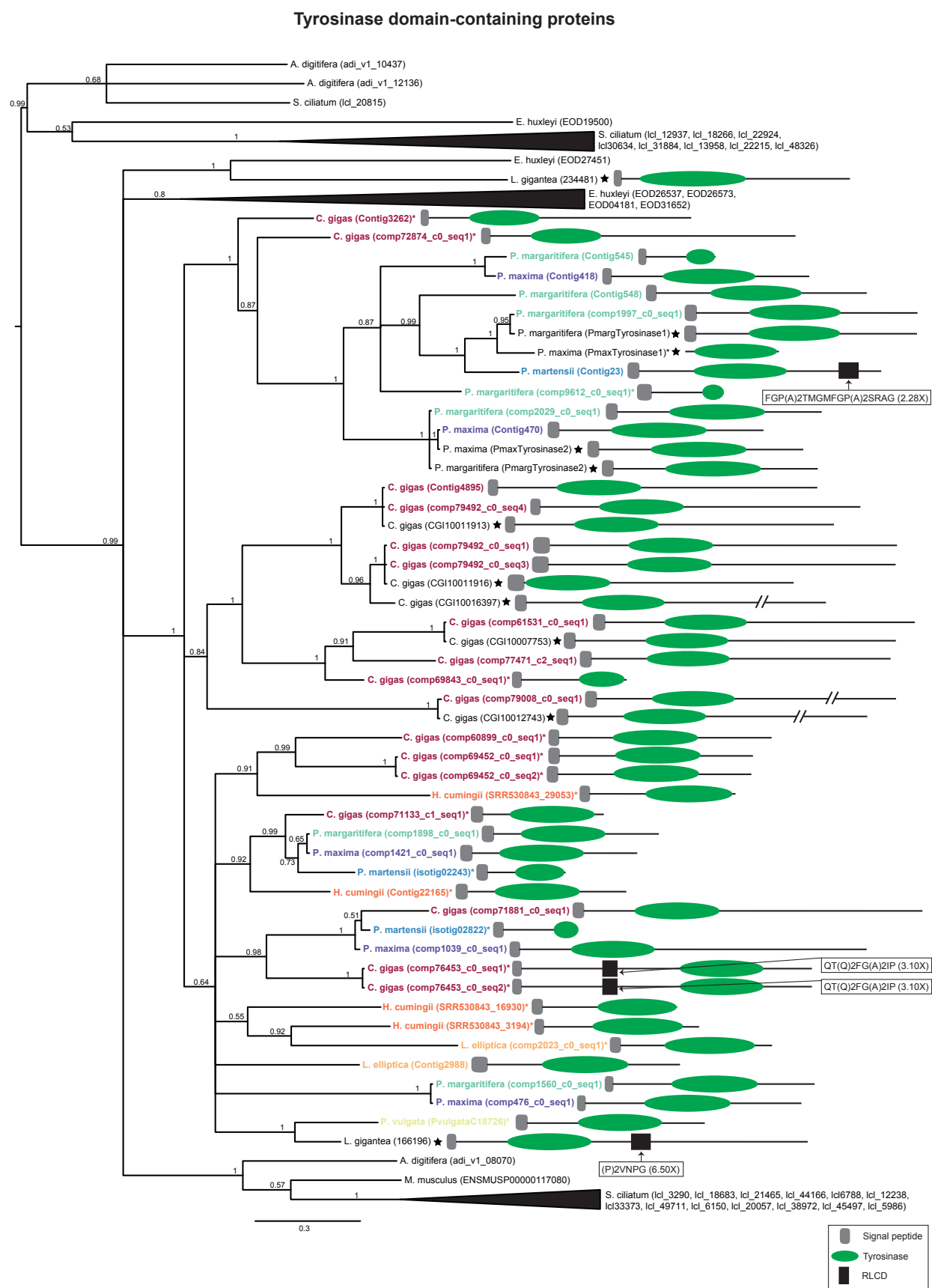
# SPARC domain-containing proteins



# Zona pellucida-like domain-containing proteins

**Figure S8. Phylogenetic analyses of SPARC domain-containing proteins and zona pellucida-like domain-containing proteins.** Bayesian phylogenetic trees were obtained under the WAG+G and LG+I+G substitution models, respectively. The secreted protein acidic and rich in cysteine Ca binding domain-containing protein tree is rooted using the calcareous sponge sequence as outgroup, while the Zona pellucida domain-containing protein tree is rooted using the midpoint-rooted option. Statistical support for each node is indicated as posterior probabilities after 1,000,000 and 6,000,000 generations, respectively. Domain architectures of molluscan secreted proteins are shown on the right of each sequence, and domain nomenclature is shown in the rectangles. Proteins previously extracted from shells and/or known to be involved in shell formation are indicated by a black star. For each RLCD, the consensus repeat sequence is shown, followed by the copy number (in brackets).
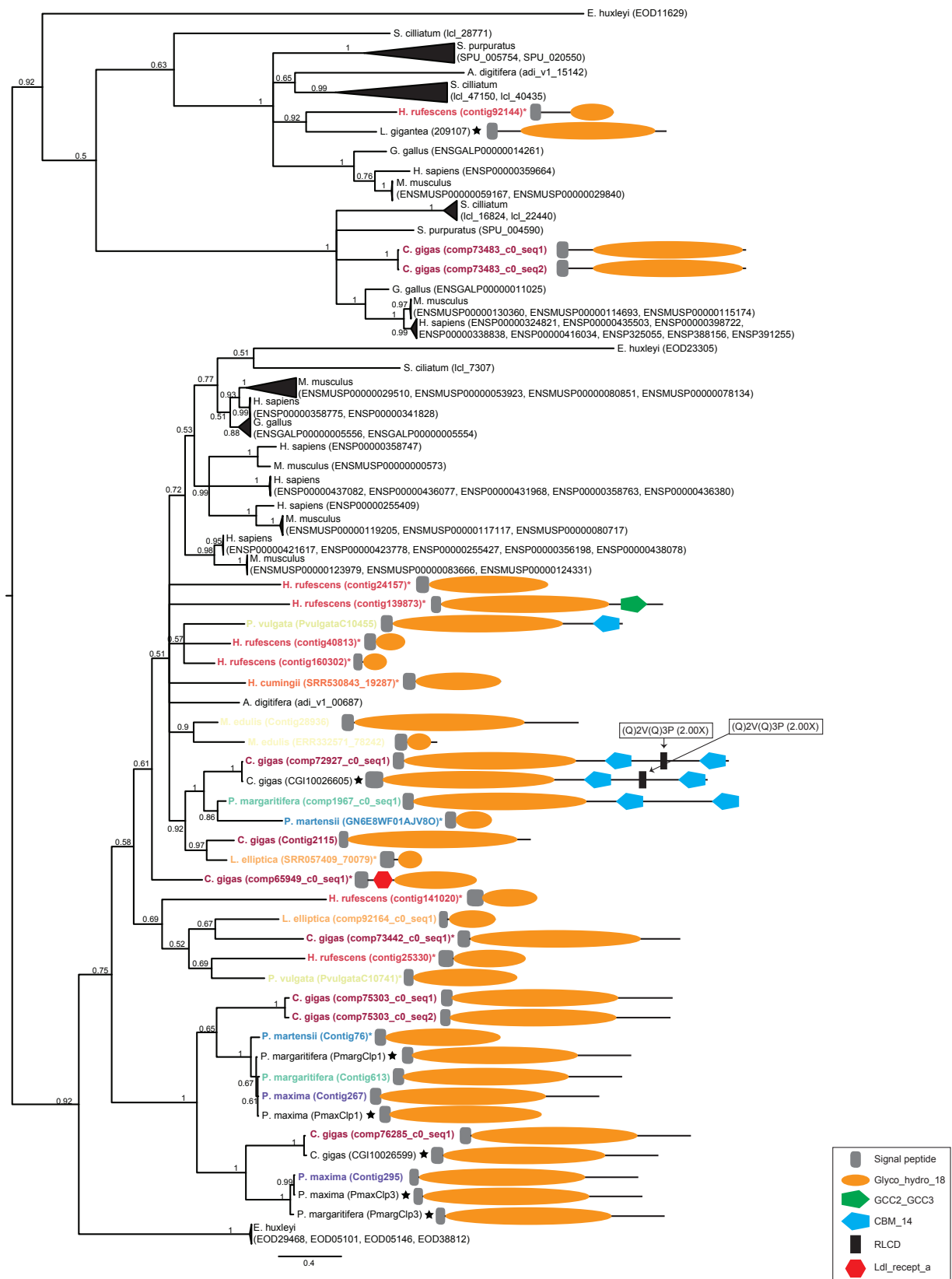
# Tyrosinase domain-containing proteins



**Figure S9. Phylogenetic analysis of tyrosinase domain-containing proteins.** Bayesian phylogenetic tree was obtained under the LG+I+G substitution model. The tree is rooted

using the midpoint-rooted option. Statistical support for each node is indicated as posterior probabilities after 5,000,000 generations. Domain architectures of molluscan secreted proteins are shown on the right of each sequence, and domain nomenclature is shown in the rectangle. Proteins previously extracted from shells and/or known to be involved in shell formation are indicated by a black star. For each RLCD, the consensus repeat sequence is shown, followed by the copy number (in brackets).
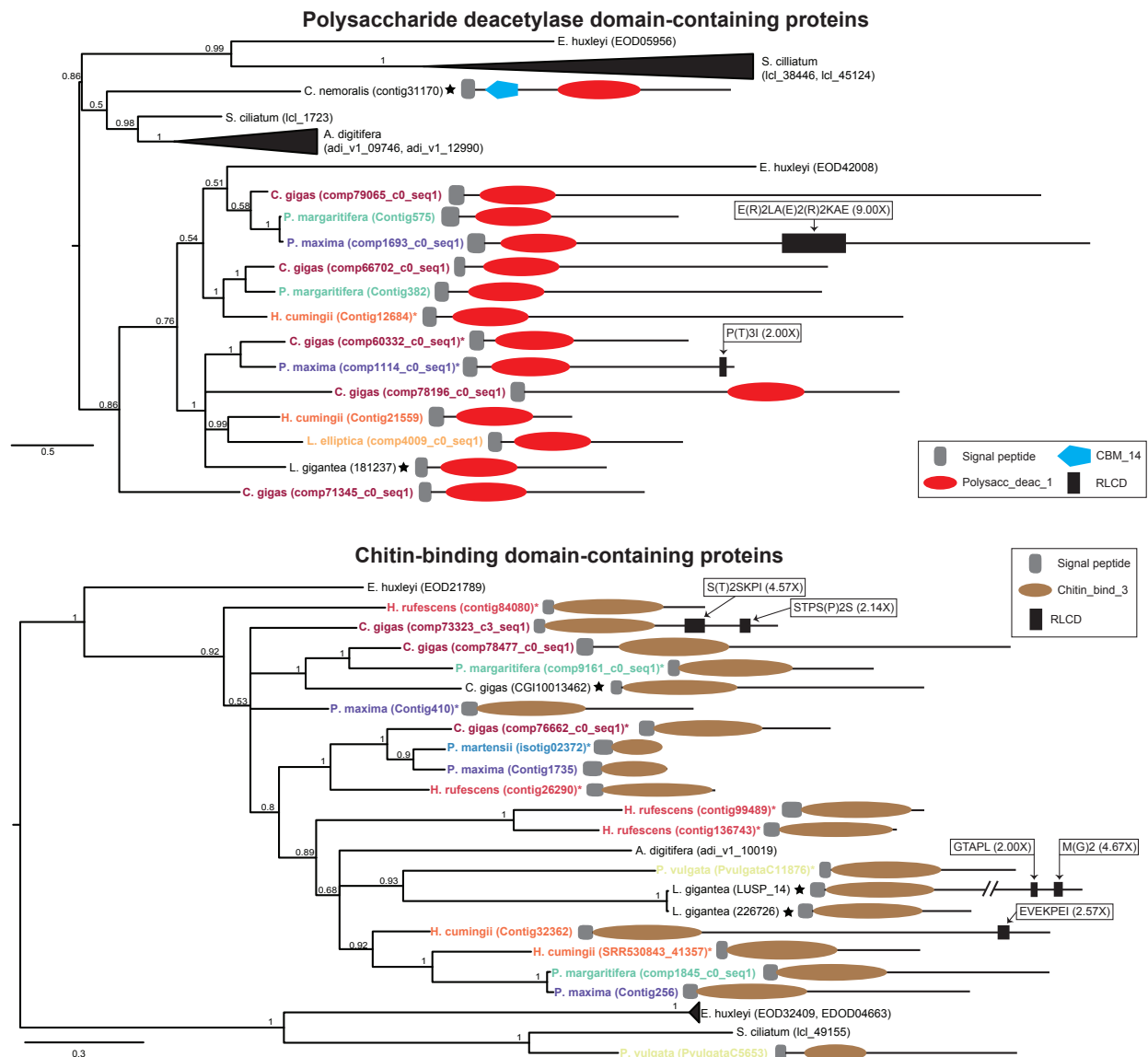
**Glycoside hydrolase family 18 domain-containing proteins**

**Figure S10. Phylogenetic analysis of glycoside hydrolase family 18 domain-containing proteins.** Bayesian phylogenetic tree was obtained under the LG+G

substitution model. The tree is rooted using the midpoint-rooted option. Statistical support for each node is indicated as posterior probabilities after 4,000,000 generations. Domain architectures of molluscan secreted proteins are shown on the right of each sequence, and domain nomenclature is shown in the rectangle. Proteins previously extracted from shells and/or known to be involved in shell formation are indicated by a black star. For each RLCD, the consensus repeat sequence is shown, followed by the copy number (in brackets).



**Figure S11. Phylogenetic analyses of polysaccharide deacetylase domain-containing proteins and chitin-binding domain-containing proteins.** Bayesian phylogenetic trees were obtained under the LG+I+G and WAG+G substitution models, respectively. Both trees were rooted using the midpoint-rooted option. Statistical support for each node is indicated as posterior probabilities after 2,000,000 and 1,000,000 generations,

respectively. Domain architectures of molluscan secreted proteins are shown on the right of each sequence, and domain nomenclature is shown in the rectangles. Proteins previously extracted from shells and/or known to be involved in shell formation are indicated by a black star. For each RLCD, the consensus repeat sequence is shown, followed by the copy number (in brackets).
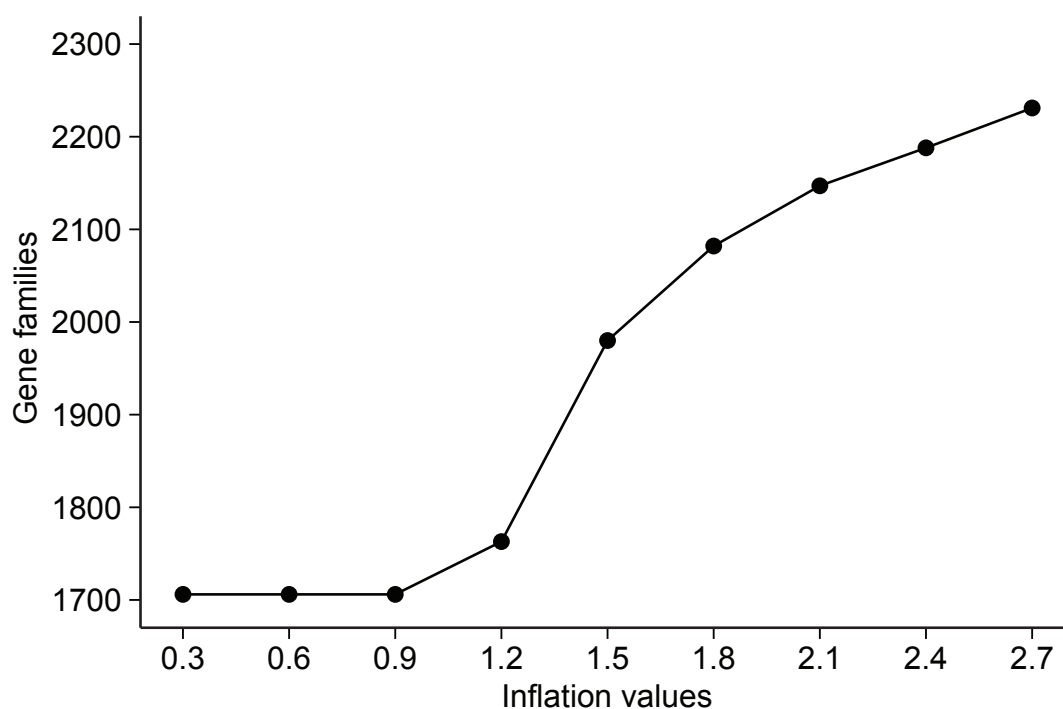


**Figure S12. Changes in the prediction of gene families using different MCL inflation values.** The effect of varying the MCL inflation value on the number of gene families are shown.