

REPORTS



I

C [?][?][?][?][?][?][?][?] [?][?][?][?]

L D F H L K D

H
I

BA

C [?][?][?][?][?][?][?] [?][?][?]

specific shell proteins¹¹. Recently, high-throughput transcriptomic and proteomic analyses have been conducted to identify shell matrix proteins (SMPEs) in *Haliotis asinina*⁸, *Pinctada margaritifera*¹², *Pinctada maxima*¹², *Crassostrea gigas*¹³, *Lottia gigantea*¹⁴, *Cepaea nemoralis*⁶, *Mytilus coruscus*¹⁵, *Pinctada fucata*¹⁰, *Magellania venosa*², *Mytilus galloprovincialis*¹⁶, *Mya truncata*¹⁷, and *C. gigas*, *Mytilus edulis*, and *Pecten maximus*¹⁸. However, only a few global SMPEs comparisons have been performed in *C. nemoralis* and *M. venosa*, relating to similarities. Even in basic, different models have been proposed for calcification

| Gene name | Categories | Expression in mantle ^a | Best matched gene ID | BLAST best hit to NCBI accession (species) | Shell layer ^b | Domains ^c |
|-----------------|-------------------|-----------------------------------|----------------------|--|--------------------------|-----------------------------|
| ACCBP1 | CaCO ₃ | N | CGI_10024902 | EKC41060 (<i>C. gigas</i>) | . | NCLBD |
| ACCBP2 | CaCO ₃ | N | CGI_10024903 | EKC41058 (<i>C. gigas</i>) | . | NCLBD |
| BMSP | Chi in | N | CGI_10009194 | BAK86420 (<i>M. galloprovincialis</i>) | P,N | VWA;CBD |
| CaLP | CaCO ₃ | N | CGI_10011294 | P41041 (<i>Pneumocystis carinii</i>) | P,N | CaBEF |
| CaM | CaCO ₃ | N | CGI_10011293 | EKC20234 (<i>C. gigas</i>) | P,N | CaBEF |
| Cg_r1 | O. herc | S | CGI_10007793 | EKC29813 (<i>C. gigas</i>) | . | CBS |
| CgT_r2 | O. herc | S | CGI_10011913 | EKC18549 (<i>C. gigas</i>) | p | CBS |
| Chi in n_hare1 | Chi in | H | CGI_10009438 | AAY86556 (<i>Atrina rigida</i>) | P,N | MHD |
| Chi in n_hare2 | Chi in | N | CGI_10012656 | BAF73720 (<i>P. fucata</i>) | P,N | MHD |
| chi obiarc | Chi in | S | CGI_10007857 | H2A0L6 (<i>P. margaritifera</i>) | P,N | Gl_co_20 |
| Chi o_rioidare1 | Chi in | S | CGI_10024867 | AFO53261 (<i>Hyriopsis cumingii</i>) | P,N | Gl_co_20 |
| Chi o_rioidare2 | Chi in | S | CGI_10026605 | CAI96027 (<i>C. gigas</i>) | P,N | Gl_co_20 |
| Clp3 | Chi in | S | CGI_10026599 | H2A0L5 (<i>P. margaritifera</i>) | P | Gl_co_18 |
| CopAmO | ECM | H | CGI_10026457 | EKC31553 (<i>C. gigas</i>) | P | Copper amine oxidase domain |
| EGF-ZP1 | ECM | S | CGI_10017543 | P86785 (<i>C. gigas</i>) | . | EGF;ZP |
| EGF-ZP2 | ECM | S | CGI_10017544 | EKC41439 (<i>C. gigas</i>) | . | EGF;ZP |
| EGF-ZP3 | ECM | H | CGI_10017545 | P86954 (<i>C. gigas</i>) | . | EGF;ZP |
| Fibronec.in1 | ECM | H | CGI_10016964 | EKC41462 (<i>C. gigas</i>) | P | bronec.in. pe III |
| Fibronec.in2 | ECM | H | CGI_10016965 | EKC41461 (<i>C. gigas</i>) | P | bronec.in. pe III |
| Pero_idare1 | O. herc | S | CGI_10023200 | EKC34657 (<i>C. gigas</i>) | . | CaBS |
| Pero_idare2 | O. herc | S | CGI_10010240 | EKC26108 (<i>C. gigas</i>) | . | CaBS |
| PFMG9 | O. herc | H | CGI_10010153 | ADC52432 (<i>P. fucata</i>) | . | KAZAL_FS |
| Pif-like1 | Chi in | N | CGI_10014497 | AKV63183 (<i>P. fucata</i>) | P,N | VWA;CBD |
| Pif-like2 | Chi in | N | CGI_10017473 | BAK86420 (<i>M. galloprovincialis</i>) | P,N | VWA;CBD |
| SPARC | ECM | N | CGI_10005088 | AND99565 (<i>P. fucata</i>) | . | N-terminal acidic domain |

Table 1. Identification and characterization of SMPs from *C. gigas*. ^arepresented the gene expression level in mantle compared in other organs: S = special, H = high, N = no special or high. ^brepresented the shell layer: P, N means the SMP were found in both the layers; . - represented unknown. ^cAbbreviation: CBS = copper-binding site; MHD = motif in head domain; VWA = von Willebrand factor (WF) type A domain; CBD = chitin-binding domain; CaBS = calcium-binding site; EGF = Epidermal growth factor; ZP = zonapellucida; CaBEF = Ca²⁺-binding EF hand domain; NCLBD = neuronal ion-channel ligand-binding domain; KAZAL_FS = Kallikrein protease inhibitor and follistatin-like domains; Gl_co_18 = family 18 Glucosyl hydrolase; Gl_co_20 = family 20 Glucosyl hydrolase.

Since the SMP gene *Nacrein* has been cloned and characterized. Besides the shell's superior mechanical and remarkable biocompatibility properties, the high commercial value of pearl has made pearl one of the best studied biomineralization models. Here, we used these SMPs mainly from *P. fucata* to identify the homologues in *C. gigas* and analyzed their structural domains underlying biomineralization. A total of 58 SMP sequences were used to perform BLAST (Supplementar Table S3) and 25 protein sequences were annotated and characterized (Table 1). Furthermore, these SMPs were classified into four categories based on functions and domains: crystallization of CaCO₃, chitin related proteins, ECM related proteins and other proteins.

Crystallization of CaCO₃. *ACCBP*, *CaLP*, *Nacrein*. Amorphous calcium carbonate-binding protein (ACCBP) is a member of the acetylcholine-binding protein family that was first isolated from the cephalopod *P. fucata*³¹. It has been reported that ACCBP inhibits undesired crystal growth and plays a key role in forming the seeding order of microcrystals of nacre³¹. ACCBP contains an ID region near the N-terminus (ACCN). It has been implicated in regulation of CaCO₃ precipitation³². BLAST search identified nine genes similar to ACCBP (E-value < e-21, score > 100). Phylogenetic analysis of these nine genes formed two clades (Fig. 2), of which ACCBP, ACCBP-like, CGI_10024902, and CGI_10024903 were clustered in one clade. These genes have CGI_10024902 and CGI_10024903 diverged from a common ancestral gene of ACCBP. These two genes are located on the same scaffold146, suggesting an in-situ origin. Both proteins contain one neuronal ion-channel ligand-binding domain (NCLBD), as well as ACCBP and ACCBP-like.

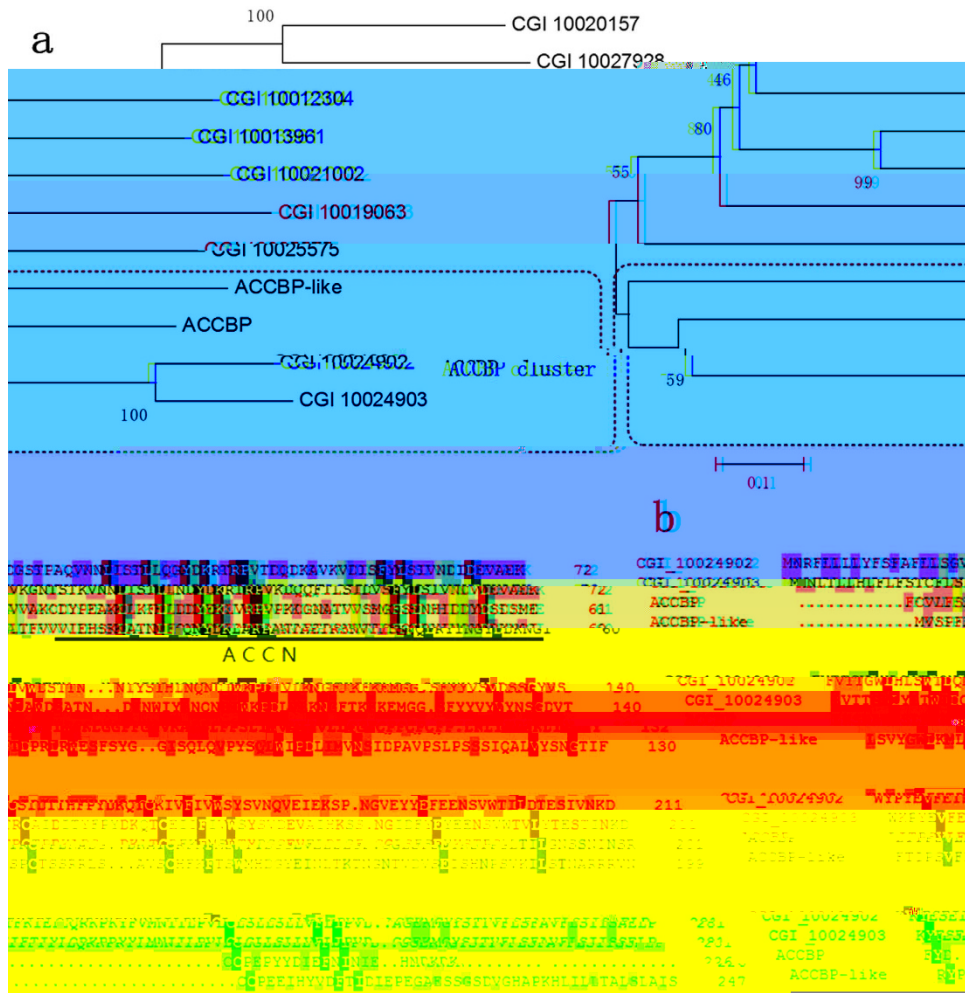


Figure 2. (a). Molecular phylogenetic tree of the nine ACCBP homologs in *C. gigas* and one ACCBP in *P. fucata*. A phylogenetic tree was inferred from the amino acid sequences using the neighbor-joining method. Bootstrap values from 1000 trials are indicated at each branch node. The scale bar indicates 0.1 amino acid replacements per site. b. Comparison of near-orthologous ion-channel ligand binding domains (NCLBD) of ACCBP and ACCN. Conserved residues are underlined.

Sequence alignments showed that the conserved residues are different from those residues between ACCBP and the nAChR family (NCLBD-containing proteins in *D. melanogaster*)³², especially in the ACCN region (Fig. 2). There are significant conserved residues among ACCBP, ACCBP-like and the other proteins in ACCN sequence, which are assumed to play a role in mineral ion acquisition.

Calmodulin-like protein (CaLP) is a multifunctional calcium sensor that belongs to a new member of the CaM superfamily^{33,34}. In biological systems, CaLP contains two Ca²⁺-binding EF-hand domains, each of which contains a pair of EF-hand motifs. Immunoblotting revealed that CaLP is localized in the organic layer and is shared between nacre (aragonite) and the prismatic layer (calcite) and the prismatic layer in *P. fucata*³⁵. These results suggest that CaLP might be involved in the growth of nacre layer and prismatic layer. We have identified three potential genes encoding CaLP using BLAST search (E-value < E-05, score > 50). BLAST search identified 26 gene models (E-value < E-05 and score > 100), which included the aforementioned three gene models in the best hit. Using the best hit protein sequence in BLAST, the phylogenetic tree showed that CGI_10011294 and CaLP consist of a single clade, as well as CGI_10011293 and CaM (Supplemental Fig. S2). Therefore, we annotated CGI_10011294 as CaLP, and CGI_10011293 as CaM.

Nacrein is the re-identified molluscan organic matrix component³⁶. Nacrein expressed throughout the entire mantle epithelium and its function in the production of both prismatic and nacreous layers^{37,38}. Nacrein contains two functional domains, a CA and a NG-repeat domain with a repeat sequence rich in Asn and Glu³⁶. Nacrein-related proteins have been found in some congeneric species (*P. maxima* and *P. margaritifera*) and also one gastropod, *Turbo marmoratus*³⁹. There has been a similar primary structure of nacrein in *P. fucata*¹¹. No nacrein-related genes could be identified using BLAST search with E-value < E-21 and score > 100. However, four nacrein-related genes were identified with E-value < E-05 and score > 80. Sequence alignments showed that the fourth hit has no NG-repeat domain and corresponding amino acids³⁹. Phylogenetic tree of the aforementioned protein sequences showed a different clade of nacrein group and hit group (Supplemental Fig. S3). These

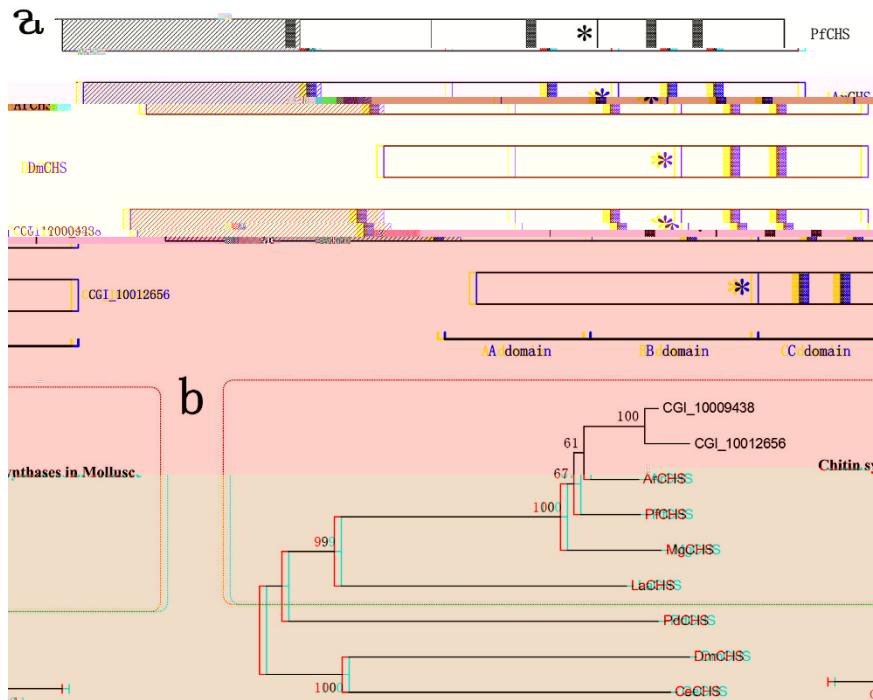


Figure 3. Diagrams of domain structures and phylogenetic tree analyses. (a). Diagrammatic representation of the domain structure of chitin synthase in different species. The main head domain is represented by a shaded box. Predicted coiled-coil regions are partially shaded. Asterisks represent the position of the chitin synthase acidic residue (QRRRW), which is used to align the diagram. The three domains in these proteins are denoted A, B, and C. PfCHS, ArCHS and DmCHS are chitin synthases from *Pinctada fucata*, *Atrina rigida* and *Drosophila melanogaster*⁶⁶. (b). A phylogenetic tree showing the evolutionary relationship of the predicted chitin synthases in *C. gigas* with known chitin synthases. The tree was constructed using maximum likelihood method. Sequences in the QRRRW catalytic domain were aligned to generate the tree. Numbers on the nodes indicate bootstrap values. ArCHS (AAY86556.1), PfCHS (BAF73720.1), MgCHS (ABQ08059.1), LaCHS (AHX26699.1), PdCHS (AHX26708.1), DmCHS (AAG09735.1) and CgCHS (NP_493682.2) are chitin synthases from *A. rigida*, *P. fucata*, *Mytilus galloprovincialis*, *Leptochiton asellus*, *Platynereis dumerilii*, *Drosophila melanogaster*, *Caenorhabditis elegans*.

for rhinoceros are not specifically or highly expressed in mantle. Taken together, the foregoing data for rhinoceros CA, but not nacrein, suggest the NG-repeat domain in CA proteins was acquired independently in the lineages of bivalves and gastropods.

These three SMPs have been reported to participate in concentration of Ca^{2+} and CO_3^{2-} , crystallization, and inhibition of crystallization.

Chitin Related Proteins: Chitin synthase, Clp. Chitin is a major component of the mollusk nacre and primary organic matrix, which plays an important role in biomineralization⁴⁰. Chitin synthase is expressed in mantle edge, contributing to the formation of the framework for shell calcification²⁹. Chitin synthase characterized so far comprises three domains, A, B, and C²⁹. It is noteworthy that the chitin synthase identified in Mollusca share a special feature—a main domain in the N-terminal^{29,41}. We have identified two predicted gene models (CGI_10009438, CGI_10012656) that are homologous to chitin synthase using BLASTn search (E-value < E^{-63} , score > 250). The two predicted gene models were located in different scaffolds. They encode two different predicted proteins that have also been identified by BLASTp with E-value = 0, score > 1800.

We aligned the amino acid sequences of chitin synthases and discovered that these two proteins included all domains (a main domain and A, B, C domains) in PfCHY (Fig. 3). COILS analysis of CGI_10009438 showed a strong potential for coiled coil formation for positions in PfCHS. Specifically, three positions showed a strong potential and one position showed a fair weak potential. COILS analysis of CGI_10012656 showed a strong potential for coiled coil formation at the same relative positions in DmCHS. The phylogenetic tree reconstructed has placed the chitin synthase from *C. gigas* belonging to the chitin synthase family (Fig. 3).

The mollusk chitin synthase group including PfCHS, ArCHS, MgCHS is separated from other chitin synthase groups forming an independent cluster. The position of CgCHS in the phylogenetic tree consists of the evolutionary position in Mollusca. These data suggest that the chitin synthase of mollusk species gained the coiled-coil sequence in the N-terminal region during evolution.

Chitinase-like proteins (Clp) were first described from *P. margaritifera* and *P. maxima*. Clp transcripts were localized in the mantle edge specifically implicated in the biomineralization of the prisms¹². Clp3 in *P. margaritifera* is located in a region of GH18_chitinase-like, which can hydrolyze chitin. Using the Clp3 sequences,

14 gene models are identified by BLASTp with E-value < E-39, Score > 150. Of these, three highly conserved CGI_10026599, CGI_10024867, CGI_10026605 are described as special examples in mantle of *C. gigas*¹³, which strongly suggest that these gene models are related to shell formation. These three highly conserved GH18 domain, showing high degree of conservation.

Both of chiitinase and chitinase play key roles in construction and reconstruction of chiitin framework, just like chitin-silk protein-acidic macromolecule model.

ECM related Proteins: Pif and BMSP, EGF-ZP, SPARC. Pif is an acidic matrix protein which regulates nacre formation. It was first identified in the pearl oyster *P. fucata*. The Pif gene encoded a precursor protein, which was processed into two cleaved products Pif 97 and Pif 80, respectively. Pif 97 has two conserved domains, a VWA domain for protein-protein interaction and chitin-binding domain. Pif 80 has aragonin-binding activity⁴². Sequence analysis revealed that the mollusk shell protein (BMSP) is a Pif homolog and a precursor protein in *M. galloprovincialis*. BMSP consists of a signal peptide and two proteins, BMSP 120 and BMSP 100, respectively⁴³. In addition, other Pif homologs from bivalves (*P. margaritifera*, *P. maxima*, and *Ptereria penguin*) and gastropods (*L. gigantea*) have been identified, suggesting a common ancestral gene duplication⁴³.

We have identified six Pif and BMSP genes by BLASTp search (E-value < E-05 and score > 100), but no Pif and BMSP gene hits using BLASTn search (E-value < E-05, score > 50). Of these, the gene CGI_10009194 is homologous to BMSP with E-value = 0, score = 1415. Another two genes, CGI_10014497 and CGI_10017473, consist of the VWA domain and a CBD. We have identified seven Pif genes by BLASTp (E-value < E-05 and score > 100) using Pif177, but none using BLASTn search (E-value < E-05, score > 50). In addition, the former three genes (CGI_10014497, CGI_10017473, CGI_10009194), one more gene (CGI_10006697) with both VWA domain and a CBD is considered to be homologous to Pif, otherwise have neither VWA domain nor CBD. In phylogenetic analysis, the CGI_10009194 and BMSP formed a separate clade, suggesting that CGI_10009194 should be annotated as BMSP. The CGI_10014497 and CGI_10017473 clustered with Pif genes from *L. gigantea* in a clade, while CGI_10006697 formed a single clade with others and showed an evolution, not being annotated as Pif (Fig. 4). Phylogenetic analysis of VWA domain showed that the former three VWA domains in CGI_10009194 formed a single clade consistent with the former three in BMSP, having evolved from BMSP-4 and CGI_10009194-4 (Fig. 4). These data indicated that BMSP and CGI_10009194 were likely evolved from an ancestral Pif for VWA domains. The phylogenetic relationship showed that VWA domains in CGI_10014497 and CGI_10017473 have evolved from CGI_10009194-4 and Pif177, being annotated as *Pif-like1* and *Pif-like2*. Schematic representation of Pif showed the common VWA domains, CBDs, chitin-binding like domains and aragonin C-terminal aragonin-binding sequences (Fig. 4).

Epidermal growth factor (EGF) domain-containing SMPs were first identified from *C. gigas*⁴⁴ and named as Cgiga-IMSP-2. They were first described in *P. maxima*, *P. margaritifera*, and *L. gigantea*^{12,14}. Generally, the SMP consists of both EGF-like domain and one non-pellucida (ZP) domain. The presence of both domains in one protein is uncommon¹². EGF-like domains, which are characterized by cysteine-linked in a characteristic pattern of disulfide bonds and small loops, are among the smallest biomolecules distributed in extracellular proteins. EGF-like domains occur in a variety of proteins associated with diverse biological functions such as cell adhesion, signaling, and Ca²⁺-binding⁴⁵. The ZP domains are present in a range of extracellular laminar matrix proteins. The ZP domain is characterized by eight conserved cysteine residues, which are involved in protein polymerization. The specific function of the EGF- and ZP-containing SMPs in calcified shell biomineralization is still unknown. The EGF-ZP genes were identified by BLASTp using the known IMSP-2 genes of *C. gigas* with (E-value < E-48, score > 100). CGI_10017543, CGI_10017544, and CGI_10017545 are located on the same scaffold 120. In combination with their relationship, high degree of sequence identity (37.53%), strongly suggest that they originated from a gene duplication event. A sequence alignment of these EGF-ZP illustrates a strong conservation of each domain (signal peptide, EGF, ZP), suggesting a fundamental role in biomineralization (Fig. S4).

SPARC (secreted protein, acidic, rich in cysteine), also known as BM-40 or osteonectin, is a major noncollagenous matrix protein of bone and a common mineralization-related protein of vertebrates and molluscs⁴⁵.

The primary structure of SPARC is characterized by the presence of three functional domains: the N-terminal acidic domain I; the follistatin-like domain II with 10 conserved cysteine residues; and the C-terminal domain III, which is involved in interactions with collagen molecules⁴⁶. By searching the *C. gigas* genome, we could only identify a single SPARC gene (CGI_10005088). Thus, the *C. gigas* genome contains one SPARC gene as observed in most triploblastic organisms. The SPARC protein consists of a region Kα-like proteinase inhibitor, follistatin-like domain (KAZAL_FS), and a region extracellular Ca²⁺ binding domain (SPARC_EC). KAZAL_FS can inhibit serine proteases and play an important role in its specific regulation. SPARC_EC function is related to cell-matrix interactions and binding to proteins such as, acidic, rDNA rkn1.60000038(P)84-6.199999861(εa)0000038 (

and oxidation of *o*-diphenols. It is also classified as a phenolase. Tyrosinase is well known for its biological role in melanin biosynthesis via transformation of tyrosine to L-DOPA. Tyrosinase functions in pigmentation and innate immunity⁴⁷. In addition, other products of the melanin pathway participate in cyclodextrin formation in insects⁴⁹.

In Mollusca, tyrosinase has been suggested in pigmentation and biomineralization of sea shells. Cephalopod tyrosinases are expressed in the ink sac, suggesting an important role in melanin production⁵⁰. In *P. fucata*, three tyrosinase genes have been characterized, P1 and P2 are suggested of function in primary formation and OT47 is proposed to influence the periodic form formation^{51,52}. In *C. gigas*, *Cgtyr1* was cloned and proposed to specifically function in the initial phase of the larval shell biogenesis²³. *CgTyr2* was also cloned and showed high levels of expression in mantle edge. It has been suggested to play a role in the formation of periodic form/pigmentation²⁴. These reports strongly suggest that tyrosinase plays diverse roles in stages when sea shells are constructed, and correlated with the periodic form.

Ten tyrosinase genes were identified from *C. gigas* genome. Two tyrosinase genes CGI_10007793 and CGI_10011913 were found to be identical to the reported *Cgtyr1* and *CgTyr2*, respectively (Fig. 5). The tyrosinase gene family can be further classified into three types: secreted form with signal peptides (Type A), cytosolic form (Type B) and member-bound form (Type C). According to SignalP 4.0, and TMHMM Server 2.0, there are six Type A tyrosinase genes, 15 Type B tyrosinase genes and five Type C tyrosinase genes. The phylogenetic tree of 26 tyrosinase genes showed that there are eight pairs of duplication genes. Among them, only one pair of CGI_10021076 Type C and CGI_10021075 Type A, CGI_10009319 Type A and CGI_10009318 Type B are located in the same scaffold/paralog, belonging to in-paralog duplication. The phylogenetic tree showed that the clusters of tyrosinase

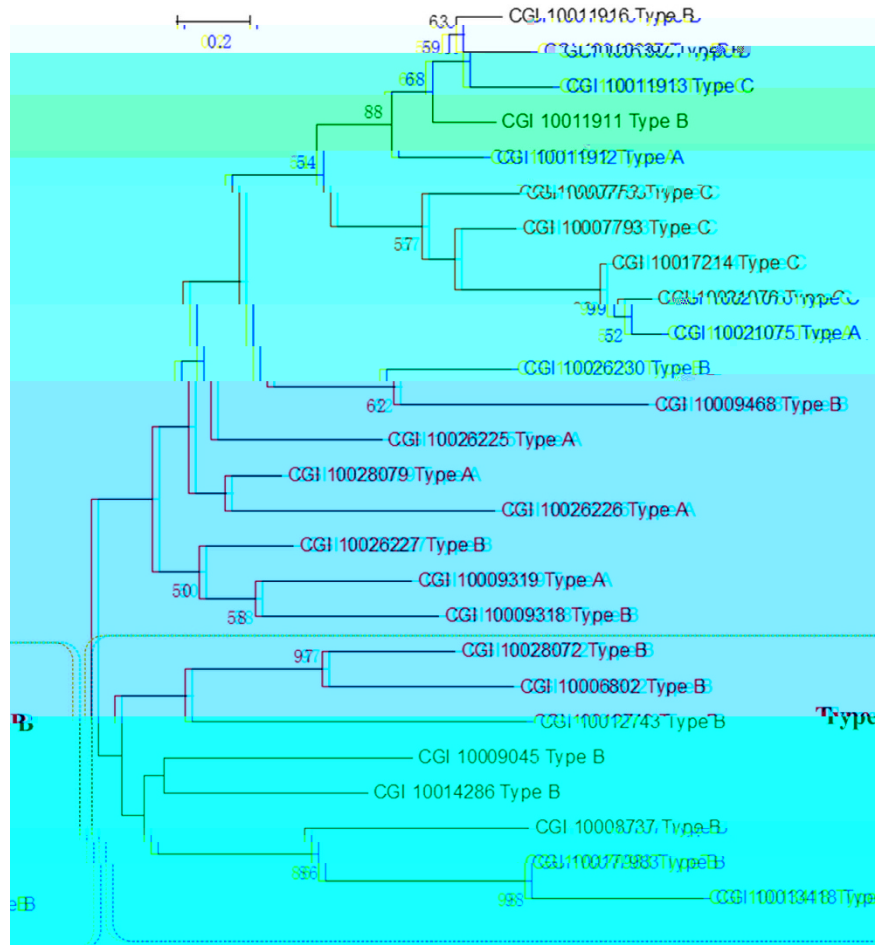


Figure 5. Tyrosinase genes diversity in *C. gigas*. The phylogenetic tree of tyrosinase genes constructed by the maximum likelihood method. Numbers on the nodes indicate bootstrap values. The tyrosinase genes are named by accession number and type.

Type B, CGI_10011913 Type C, CGI_10011911 Type B and CGI_10011912 Type A are clustered, all of which are located in Scaffold 43702. Above all, it suggests that the tyrosinase has evolved through both inorganic duplication and inorganic duplication.

Peroxidases are iron proteins that catalyze the oxidation of many aromatic amines and phenols by hydrogen peroxide. In addition, they are involved in the DOPA reaction in the mantle of *Lymnaea stagnalis* catalyzed by peroxidase, suggesting that peroxidase is involved in the quinone-tanning of periostracum proteins⁵³. Peroxidase is also expressed in the ink gland of *Sepia officinalis*, likely involved in melanin biosynthesis⁵⁴. Typical peroxidase is characterized by a unique conserved domain that contains histidine (proximal and distal histidines) and one calcium-binding site, which are suggested to function in maintaining the protein structure in the heme-iron center^{54,55}. Peroxidase has been retrieved from the shell matrix of *P. margaritifera* and *L. gigantea*. It is proposed to be involved in biomaterial hydrogel formation via protein matrix framework assembly¹². We identified 26 peroxidase genes by BLASTp using the known peroxidase genes of *P. margaritifera* as a query (E-value < E-23, score > 100). The peroxidase genes have gone through an expansion as shown for tyrosinase in *C. gigas*. Nine of the 26 peroxidase genes are special or highly expressed in the mantle of *C. gigas*. The ones are CGI_10023200 and CGI_10010240 are special call expressed in the mantle¹³. All nine peroxidase consist of two histidines and one calcium-binding site. Phylogenetic analysis showed that these peroxidases formed a cluster that could be identified as melanin biosynthesis group and the shell formation group (Fig. 6). Peroxidase from *Drosophila melanogaster*, *Bombyx mori* and *Sepia officinalis* have been implicated in melanin synthesis form melanin polymer⁵⁴.

We noted in our bioinformatics analysis that in many cases, BLASTn gave no hits, however, BLASTp gave hits (Table S2). It suggests that the genes encoding SMPs have undergone more variation than SMP proteins sequences. Until now, literature on biomaterials is scarce and it has been identified SMPs. Obviously, SMPs do more than providing a framework for crystallization. They could be involved in interactions with other macromolecular components of the matrix and cells, giving a feedback between the shell and the calcifying mantle epithelium.

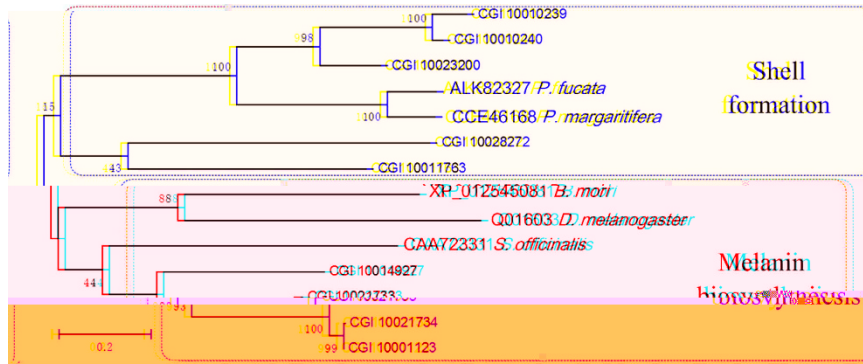


Figure 6. Peroxidase genes diversity in *C. gigas*. The phylogenetic tree of peroxidase genes constructed by the maximum likelihood method. Nucleotide sequence

C **E** **C** In *C. gigas*, a SMPE consisting of 53 proteins was recently published¹⁸. The SMPE was constructed from peptide fragments in shell, which can be mapped to the genome. In addition, a SMPE in *C. gigas* consisting of 259 proteins was also analyzed, which was constructed by Zhang and colleagues using the same method. Many shell-keeping proteins, such as elongation factors 1 and ribosomal proteins, were found in these 259 SMPs, which are significantly more than SMPs identified from other molluscan shell proteins. Generally, shell-keeping proteins should not be found in shell matrix proteins and they do not play a special role in biomineralization. Given that shell-keeping genes are generally expressed at a relatively constant level in most non-paleontological situations, the identified genes have specific or highly expressed mantle locomotion. The inference from shell-keeping proteins. An integrated SMPE consisting of 76 SMPs was constructed. The resulting SMPE was constructed from the intersection of (1) 259 proteins isolated from *C. gigas* shells; (2) 492 genes that were specifically or highly expressed in mantle, of which, the highly expressed genes were defined as having RPFM. All of a least 5 and a least 10 times of other organa era¹³. In principle, the conserved/disordered (ID) domains were predicted by IUPRED⁵⁶ and XSTREAM⁵⁷. IUPRED was used to recognize disordered regions from the amino acid sequence of SMPs based on the estimated pairwise energy content. XSTREAM was used to detect proteins in tandem-arranged repeats in the default settings.

C Based on BLASTp, the global similarity comparison of *C. gigas* SMPE was performed again. 443 SMPs derived from nine other biocalcifying mollusks were included: 53 *M. coruscus* proteins¹⁵, 75 *P. fucata* proteins¹⁰, 45 *P. margaritifera* proteins¹², 26 *P. maxima* proteins¹², 39 *L. gigantea* proteins¹⁴, 14 *H. asinina* proteins⁸, 59 *C. nemoralis* proteins⁶, 67 *M. truncata¹⁷, and 65 *M. venosa² proteins. All protein sequences are aligned by mapping to the reference sequences, ESTs in biocalcifying organs or genome assemblies. The alignment threshold was set to 0.1e-06. The sequence comparisons were made using blast+⁵⁸; blastp -p 0.01 -e 1e-10 -db XX.change.fasta -o XX.blp -of 6 -e 1e-06 -n m_hread: 10. The *.blp files generated by blast+ were modified using cscript and then passed to Circos in order to generate an ideogram⁵⁹. The *.blp files in the default settings of the BLASTp results are provided (Supplemental Table S4).**

I **C** The SMP searches in the *C. gigas* gene models (osteopetrolin_1) and indexed protein sequences (osteopetrolin_1) were performed using Oysterbase¹³. Shell formation complementary DNAs (cDNAs) in *C. gigas* were BLASTn and BLASTp searched. Identifications of the hit gene models associated with the original cDNAs were confirmed by sequence alignment. Similarly, SMPs identified from other molluscan species were BLAST searched, and the obtained gene models were reciprocally BLASTp searched against the NCBI nonredundant (nr) database to confirm the best hit sequence.

C **C** The conserved structural domains were examined using the SMART⁶⁰ and InterProScan⁶¹. The amino acid sequences were aligned using MEGA5⁶² or DNAMAN (Lnon Bio). For phylogenetic analysis, poorly aligned positions were checked and removed manually. The phylogenetic tree was rapidly constructed using MEGA5. Protein sequences for phylogenetic analysis were retrieved from GenBank, Swiss-Prot, or the Oysterbase. In the case of secretory proteins or peptides, the presence of a signal peptide was predicted by SignalP 4.0⁶³. COILS was used to detect coiled coil formation⁶⁴. To identify membrane proteins, the TMHMM Server 2.0 prediction algorithm was used for transmembrane helices⁶⁵.

C We compared the different SMPEs of *C. gigas* consisting of 76 SMPs and 53 SMPs and chose the latter one to perform a broad level comparison by bioinformatic analysis. The SMPE were characterized by having a high proportion of ID proteins, especially RLCD proteins. We used a marine SMPE in *C. gigas* to perform a broad comparison again. 443 SMPs from nine other species using BLASTp. Overall, the earlier findings have the SMPs similarity depends not only on the evolutionary distance but also in the encoded biomineralog of shell,

parallel evolution, adaptation of the ironmen gene. The highly conserved proteins *rocinase* and *chironin* are identified in bivalve, and the related conserved proteins in the domains of CA, VWA, CBD, IG-like and LaG are identified from all ten species. 25 genes encoding SMPs were annotated and characterized. They are chitin related or ECM related proteins involved in crystallization of CaCO₃. These conserved SMPs and nacreal domain enrich the molecular knowledge of shell formation mechanism in *C. gigas*, paving for a refined shell formation model including both chitin and ECM-related proteins.

1. Simion, P. & Wilbrink, M. *Biomineralization* (Elsevier, 2012).
2. Jackson, D. J. *et al.* The *Magellania venosa* biomineralizing prokaryote: a window into brachiopod shell evolution. *Genome biology and evolution* **7**, 1349–1362 (2015).
3. Söll, A. H. Biomineralization and evolution. *Reviews in mineralogy and geochemistry* **54**, 329–356 (2003).
4. Moll, A., Aguilera, F., McDougall, C., Jackson, D. & M. Degnan, B. Shell diversity and rapid diversification of biomineralization. *Frontiers in Zoology* **13**, 1 (2016).
5. Marie, B., Leclercq, N., Zanella-Cleon, L., Becchi, M. & Marin, F. Molecular evolution of mollusc shell proteins: insight from prokaryotic analogs of the edible mussel *Mytilus*. *J Mol Evol* **72**, 531–46 (2011).
6. Mann, J. & Jackson, D. J. Characterization of the pigmented shell-forming prokaryote of the common green snail *Cepaea nemoralis*. *BMC genomics* **15**, 1 (2014).
7. Aguilera, F., McDougall, C. & Degnan, B. M. Co-option and de novo gene evolution underlie molluscan shell diversity. *Molecular Biology and Evolution* doi: <https://doi.org/10.1093/molbev/mwz294> (2016).
8. Marie, B. *et al.* Prokaryotic analogs of the organic matrix of the abalone *Haliotis asinina* calcified shell. *Proteome Sci* **8**, 54 (2010).
9. Marin, F. & Leclercq, N., G. Molluscan shell protein. *Comptes Rendus Palevol* **3**, 469–492 (2004).
10. Li, C. *et al.* In-depth prokaryotic analogs of shell matrix proteins of *Pinctada fucata*. *Scientific reports* **5** (2015).
11. Miamoto, H. *et al.* Diversity of shell matrix proteins: Genome-wide investigation of the pearl oyster, *Pinctada fucata*. *Zoological science* **30**, 801–816 (2013).
12. Marie, B. *et al.* Diverse secretory proteins control the biomineralization processes of prism and nacre deposition of the pearl oyster shell. *Proceedings of the National Academy of Sciences* **109**, 20986–20991 (2012).
13. Zhang, G. *et al.* The oyster genome reveals diversity and complexity of shell formation. *Nature* **490**, 49–54 (2012).
14. Marie, B. *et al.* The shell forming prokaryote of *Lottia gigantea* reveals both deep conservation and lineage specificity. *FEBS Journal* **280**, 214–232 (2013).
15. Liao, Z. *et al.* In-depth prokaryotic analogs of nacre, prism, and mollusc calcification of *Mytilus* shell. *Journal of Proteomics* **122**, 26–40 (2015).
16. Gao, P. *et al.* Late Prokaryotic Analogs of *Mytilus galloprovincialis* Shell. *PLOS ONE* **10** (2015).
17. Arialan, J. *et al.* Shell matrix proteins of the clam, *Mya truncata*: roles of shell formation through prokaryotic diversity. *Marine Genomics* **27**, 69–74 (2016).
18. Arialan, J. *et al.* Insight from the shell prokaryote: biomineralization adaptation. *Molecular Biology and Evolution* **34**, 66–77 (2016).
19. Furuhachi, T., Schäringer, C., Mitsui, I., Smirnov, M. & Beran, A. Molluscan shell evolution in the evolution of shell calcification. *Comparative biochemistry and physiology Part B: Biochemistry and molecular biology* **154**, 351–371 (2009).
20. Johnson, M. *et al.* Cellular orchestrated biomineralization of crystalline calcium carbonate on implanted surfaces. *The eukaryotic cell*, *Crassostrea virginica* (Gmelin, 1791). *Journal of Experimental Marine Biology and Ecology* **463**, 8–16 (2015).
21. Monaghan, A. S., Wheeler, A. P., Paradis, P. & Snider, D. Hemocyanin-mediated shell mineralization in the eukaryotic cell. *Science* **304**, 297 (2004).
22. Yang, H., Zhao, X. & Li, Q. Genome-wide identification and characterization of long intergenic noncoding RNAs and their potential association in larval development in the Pacific oyster. *Scientific reports* **6** (2016).
23. Han, P., Li, G., Wang, H. & Li, B. Identification of a *rocinase* gene potentially involved in early larval shell biogenesis of the Pacific oyster *Crassostrea gigas*. *Development genes and evolution* **223**, 389–394 (2013).
24. Yang, X. *et al.* Molecular cloning and differential expression analysis of a *rocinase* gene in the Pacific oyster *Crassostrea gigas*. *Molecular biology reports* **41**, 5403–5411 (2014).
25. Evans, J. S. Aragonite-associated biomineralization proteins are disordered and contain intracellular motifs. *Bioinformatics* **28**, 3182–3185 (2012).
26. Jackson, D. *et al.* Parallel evolution of nacre building genes in molluscs. *Molecular biology and evolution* **27**, 591–608 (2010).
27. Feng, D., Li, Q., Yang, H., Zhao, X. & Gong, L. Comparative Transcriptome Analysis of the Pacific Oyster *Crassostrea gigas* Characterized by Shell Color: Identification of Genetic Basis of Polymorphic Pigmentation. *PLoS one* **10** (2015).
28. Aguilera, F., McDougall, C. & Degnan, B. M. Evolution of the *rocinase* gene family in bivalve molluscs: independent expansion of the mantle gene repertoire. *Acta biomaterialia* **10**, 3855–3865 (2014).
29. Saito, M., Saito, S. & Nagasawa, H. Identification of chitin in the primary structure of the shell and a chitinase gene from the Japanese pearl oyster, *Pinctada fucata*. *Bioscience, biotechnology, and biochemistry* **71**, 1735–1744 (2007).
30. Weis, I. M., Schniener, V., Eichner, N. & Sömpfer, M. The chitinase in molluscs evolved in marine bivalve molluscs shell formation contains a motif domain. *FEBS letters* **580**, 1846–1852 (2006).
31. Ma, Z. *et al.* A novel extracellular protein involved in the morphology of nacre lamellae in the pearl oyster, *Pinctada fucata*. *Journal of Biological Chemistry* **282**, 23253–23263 (2007).
32. Amos, F. E., Ndao, M. & Evans, J. S. Evidence of mineralization activity and a prokaryotic assembly of the N-terminal sequence of ACCBP, a biomineralization protein having a homologous choline binding protein family. *Biomacromolecules* **10**, 3298–3305 (2009).
33. Li, S. *et al.* Cloning and expression of a potential calcium mobilizing regulator: calmodulin involved in shell formation from pearl oyster (*Pinctada fucata*). *Comparative Biochemistry and Physiology Part B: Biochemistry and Molecular Biology* **138**, 235–243 (2004).
34. Li, S., Xie, L., Ma, Z. & Zhang, J. cDNA cloning and characterization of a novel calmodulin-like protein from pearl oyster *Pinctada fucata*. *Febs Journal* **272**, 4899–4910 (2005).
35. Yan, Z. *et al.* Biomineralization: function of calmodulin-like protein in the shell formation of pearl oyster. *Biochimica et Biophysica Acta (BBA)-General Subjects* **1770**, 1338–1344 (2007).
36. Miamoto, H. *et al.* A carbonic anhydrase from the nacreous layer in the pearl oyster. *Proceedings of the National Academy of Sciences* **93**, 9657–9660 (1996).
37. Miamoto, H., Mitsui, F. & Ohno, J. The carbonic anhydrase domain protein nacrein is expressed in the epithelial cells of the mantle and acts as a negative regulator in calcium ion homeostasis in the mollusc *Pinctada fucata*. *Zoological science* **22**, 311–315 (2005).
38. Tanaka, T. & Endo, T. Biphasic and dual coordinated expression of the genes encoding major shell matrix proteins in the pearl oyster *Pinctada fucata*. *Marine biotechnology* **8**, 52–61 (2006).
39. Nori, I., M. & Saito, T. Diversity and function of the nacrein-related proteins inferred from a comprehensive analysis. *Marine Biotechnology* **10**, 234–241 (2008).

40. Frerhashi, T. *et al.* Role of GC/MS and I²C¹⁸ in chitin analysis of mollusc shells. *Bioscience, biotechnology, and biochemistry* **73**, 93–103 (2009).
41. Weis, I. M. & Schuster, V. The distribution of chitin in the larval shells of the bivalve mollusc *Mytilus galloprovincialis*. *Journal of structural biology* **153**, 264–277 (2006).
42. Saito, M. *et al.* An acidic matrix protein, Pif, is a macromolecular for nacre formation. *Science* **325**, 1388–1390 (2009).
43. Saito, M. *et al.* Identification and Characterization of Calcium Carbonate Binding Protein, Bivalve Matrix Shell Protein (BMSP), from the Nacre Layer. *Chembiochem* **12**, 2478–2487 (2011).
44. Marie, B., Zanella-Clon, I., Gichard, N., Becchi, M. & Marin, F. Novel proteins from the calcifying shell matrix of the Pacific oyster *Crassostrea gigas*. *Marine biotechnology* **13**, 1159–1168 (2011).
45. Marer, P. & Hohenecker, E. Structural and functional aspects of calcium binding in extracellular matrix proteins. *Matrix biology* **15**, 569–580 (1997).
46. Sasaki, T., Hohenecker, E., Ghring, W. & Timpl, R. Crystal structure and mapping binding site directed mutagenesis of the collagen binding epitope of an acidic form of BM-40/SPARC/osteonon. *The EMBO Journal* **17**, 1625–1634 (1998).
47. Garcia Borrn, J. C. & Solano, F. Molecular Analysis of Troponin and its Related Protein: Beyond the Histidine Bond Metal Cationic Center. *Pigment Cell Research* **15**, 162–173 (2002).
48. Aguilera, F., McDougall, C. & Degnan, B. M. Origin, evolution and classification of pe-3 copper proteins: lineage-specific gene expansion and loss across the Metazoa. *BMC evolutionary biology* **13**, 96 (2013).
49. Andersen, S. O. Insect cuticular sclerotization: a review. *Insect biochemistry and molecular biology* **40**, 166–178 (2010).
50. Narao, A. T. *et al.* Purification, characterization and molecular cloning of troponin from the cephalopod mollusc, *Illex argentinus*. *European Journal of Biochemistry* **270**, 4026–4038 (2003).
51. Nagai, Y., Yano, M., Morimoto, Y. & Miamoto, H. Troponin localisation in mollusc shells. *Comparative Biochemistry and Physiology Part B: Biochemistry and Molecular Biology* **146**, 207–214 (2007).
52. Zhang, C., Xie, L., Huang, J., Chen, L. & Zhang, X. An oyster troponin isolated in periodontal form from the pearl oyster (*Pinctada fucata*). *Biochemical and biophysical research communications* **342**, 632–639 (2006).
53. Timmermans, L. P. S. *et al.* The shell formation in molluscs. *Netherlands Journal of Zoology* **19**, 413–523 (1968).
54. Gerardo, I., Aniello, F., Branno, M. & Palumbo, A. Molecular cloning of a perodonta NA-specific expressed in the gland of *Sepia officinalis*. *Biochimica et Biophysica Acta (BBA)-Gene Structure and Expression* **1353**, 111–117 (1997).
55. Shiro, Y., Iwano, M. & Morishima, I. Presence of endogenous calcium ion and its functional and structural regulation in horseradish peroxidase. *Journal of Biological Chemistry* **261**, 9382–9390 (1986).
56. Doñi, Z., Číž, V., Tompa, P. & Simon, I. IUPred: web server for the prediction of intrinsically disordered regions of proteins based on estimated energy content. *Bioinformatics* **21**, 3433–3434 (2005).
57. Newman, A. M. & Cooper, J. B. XST-EAM: A practical algorithm for identification and architecture modeling of tandem repeats in protein sequences. *BMC Bioinformatics* **8**, 382 (2007).
58. Camacho, C. *et al.* BLAST+: architecture and applications. *BMC bioinformatics* **10**, 1 (2009).
59. Riesen, M. *et al.* Circo: an information aesthetic for comparative genomics. *Genome research* **19**, 1639–1645 (2009).
60. Lenic, I., Doerflinger, T. & Borner, P. SMA-T: recent developments and applications in 2015. *Nucleic Acids Research* **43**, 257–260 (2015).
61. Quesillon, E. *et al.* InErProScan: protein domain identification. *Nucleic Acids Research* **33**, 116–120 (2005).
62. Tamura, K. *et al.* MEGA5: molecular evolutionary genetics analysis using maximum likelihood, evolutionary distance, and maximum parsimony methods. *Molecular biology and evolution* **28**, 2731–2739 (2011).
63. Petersen, T. N., Brannas, S., von Heijne, G. & Nielsen, H. SignalP 4.0: discriminating signal peptides from transmembrane regions. *Nature methods* **8**, 785–786 (2011).
64. Lupas, A., Van Donge, M. & Söck, J. Predicting coiled coils from protein sequences. *Science* **252**, 1162–1164 (1991).
65. Miller, S., Croning, M. D. & Apweiler, R. Evaluation of methods for the prediction of membrane spanning regions. *Bioinformatics (Oxford, England)* **17**, 646–653 (2001).
66. Tellam, L., Vocolo, T., Johnson, S. E., Jarman, J. & Pearson, D. Insect chitinase. *European Journal of Biochemistry* **267**, 6025–6043 (2000).

A

This study was supported by the grant from National Natural Science Foundation of China (31372524), Shandong Seed Project, Shandong Province (2016ZDJS06A06), and Qingdao National Laboratory for Marine Science and Technology (2015ASKJ02).

A C

D.D.F. and S.H.D. analyzed the data and wrote the paper. Q.L., H.Y. and L.F.K. conceived and designed the study.

A I

Supplementary information accompanies this paper at <http://www.nature.com/srep>

Competing Interests: The authors declare no competing financial interests.

How to cite this article: Feng, D. *et al.* Identification of conserved proteins from dipteran shell matrix proteinome in *Crassostrea gigas*: characterization of genetic background regulating shell formation. *Sci. Rep.* **7**, 45754; doi: 10.1038/srep45754 (2017).

Publisher's note: Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

