

Università degli studi di Modena e Reggio Emilia

Dipartimento di Scienze della Vita

Corso di laurea Magistrale
Biologia sperimentale ed applicata

**Analisi di trascrittomi tessuto-specifici di *Pomacea
canaliculata* tramite sequenziamento Nanopore:
Ipervariabilità ed approccio AI**

Relatore
Prof. Nicola Franchi

Tesi di laurea di
Filippo Bertolasi

Anno Accademico 2023/2024

Indice

INDICE	1
RIASSUNTO.....	2
1. INTRODUZIONE	3
1.1 IMMUNITÀ INNATA	3
1.1.1 Immunità cellulo-mediata	4
1.1.2 Immunità umorale.....	4
1.2 SISTEMA DEL COMPLEMENTO	4
1.2.1 FUNZIONI NON IMMUNITARIE DEL COMPLEMENTO	6
1.2.2 SISTEMA DEL COMPLEMENTO NEGLI INVERTEBRATI	7
1.2.3 REGOLATORI DEL COMPLEMENTO	8
1.3 <i>POMACEA CANALICULATA</i>	9
2. SCOPO DELLA TESI	10
3. RISULTATI.....	11
3.1 IPERVARIABILITÀ DA SPLICING ALTERNATIVI.....	11
.....	11
3.2 ANALISI 3D DA ALGORITMO AI.....	12
4 DISCUSSIONE.....	14
5 CONCLUSIONI	17
6 MATERIALI E METODI.....	18
6.1 ANIMALI	18
6.2 SOFTWARE E STRUMENTI BIOINFORMATICI	18
ALLEGATI.....	20
.....	20
BIBLIOGRAFIA.....	25

Riassunto

Il sistema del complemento è un componente fondamentale della difesa immunitaria sia per vertebrati che per invertebrati. Indagare circa le sue origini ed evoluzione è un passaggio chiave per comprendere come sono cambiate le sue proteine in modo da adattarsi a nuove condizioni ambientali, cicli vitali o, nel caso dei vertebrati, all'interazione con l'immunità adattativa. Le proteine responsabili della sua regolazione vengono definite Regolatori del Complemento, sono sempre più studiate perché è dimostrato il loro coinvolgimento in numerosi processi fisiologici e patologici. In *Pomacea canaliculata* queste indagini possono essere utili anche per comprendere un eventuale coinvolgimento del sistema del complemento nei processi definiti “non canonici” come la rigenerazione.

Questo studio ha dimostrato che l'espressione dei regolatori del complemento in *P. canaliculata* è tessuto-specifica, almeno per quanto riguarda emociti e rene posteriore. È stato inoltre dimostrato quanto l'ipervariabilità di questi sia determinata da fenomeni di splicing alternativo. È stato poi con successo applicato il modello AI ProstT5 permettendo la creazione di due dataset contenenti sequenze caratterizzate da un simile significato biologico a delle già note proteine regolatrici del complemento di *P. canaliculata*. L'analisi dei dati ottenuti ha individuato che la maggior parte delle sequenze non è identificabile e non presenta domini noti ai database SMART, gettando le fondamenta per future analisi circa potenziali proteine regolatrici del complemento con struttura diversa dalle convenzionali presentanti CCP.

1. Introduzione

Tutti i metazoi hanno evoluto dei meccanismi cellulari ed umorali per difendersi dal “non-self”. Storicamente, il sistema immunitario è stato suddiviso in due: immunità innata ed immunità adattativa. L’immunità innata fornisce una difesa aspecifica e rappresenta la prima forma di difesa che andrà ad agire contro l’organismo patogeno e/o qualsiasi altra molecola riconosciuta come non appartenente all’ospite, al “self”. L’immunità adattativa invece rappresenta una seconda linea di difesa, specifica per i singoli organismi/molecole antigeniche (Alberts et al., 2009).

1.1 Immunità innata

L’immunità innata viene detta anche naturale, il nome deriva dal fatto che fornisce protezione da agenti patogeni senza necessità di una prima esposizione ad essi. In altre parole, quando il sistema immunitario innato incontra un antigene non-self, reagirà istantaneamente per isolarlo o rimuoverlo dall'ospite (Kimbrell, 2001). Questa funzione è molto importante dato che la risposta immunitaria adattativa e specifica richiede anche una settimana per diventare efficace. In generale ci sono tre linee di difesa immunitaria innata che possono agire: la prima è rappresentata dalle barriere fisiche e chimiche che impediscono il semplice ingresso di microrganismi all’interno del corpo. Queste barriere comprendono lo strato di tegumento esterno, le giunzioni strette tra cellule epiteliali, il pH acido nello stomaco e i componenti degli strati di muco che inibiscono la colonizzazione ed uccidono i batteri estranei. La flora simbiote ha un ruolo anche nel proteggere le superfici corporee dagli invasori, competendo per la stessa nicchia ecologica e limitando così la colonizzazione. La seconda linea di difese innate comprende le risposte intrinseche della cellula, tramite cui una singola cellula può riconoscere una cellula bersaglio e rispondere di conseguenza con misure atte alla fagocitosi o all’attività citotossica. La maggior parte delle cellule che hanno internalizzato un batterio tramite fagocitosi indotta da patogeni, per esempio, dirigerà immediatamente il fagosoma a fondersi con un lisosoma, esponendo il microrganismo invasore agli enzimi digestivi. Un altro meccanismo di difesa intrinseco è la capacità di cellule ospiti di degradare l’RNA a filamento doppio, che è un intermediario comune nella replicazione virale; le cellule infette degraderanno anche qualsiasi RNA a singolo filamento che ha innescato il meccanismo. La terza linea di difese immunitarie dipende da una serie specializzata di proteine e cellule ad attività immunitaria (fagociti e cellule citotossiche) che riconoscono le caratteristiche conservate dei patogeni e si attivano velocemente per aiutare

a distruggere gli invasori. Tra queste troviamo le cellule specializzate degli invertebrati e neutrofili, macrofagi e cellule natural-killer dei Vertebrati oltre a sistemi molecolari complessi come il sistema del complemento (Alberts et al., 2009).

1.1.1 Immunità cellulo-mediata

Il riconoscimento di pattern molecolari conservati dei patogeni causa l'attivazione, nei Vertebrati, dei leucociti: granulociti neutrofili, macrofagi, monociti, cellule dendritiche. Il loro ruolo principale è quello di immunità cellulo-mediata che include: fagocitosi, lisi della cellula bersaglio, ed immunomodulazione. Per l'attivazione di queste cellule vengono riconosciuti i PAMP (Pathogen Associated Molecular Patterns), si tratta di pattern molecolari, costituiti ad esempio da zuccheri sulla superficie del microrganismo, che durante l'evoluzione si sono conservati e quindi permettono all'organismo di individuare la presenza di un patogeno e generare una risposta aspecifica, ma selettiva (Alberts et al., 2009).

1.1.2 Immunità umorale

La componente umorale della risposta immunitaria innata comprende: il sistema del complemento, il CAS (Contact Activation System) e le lectine. Il CAS è una cascata di proteasi plasmatiche che legano la superficie batterica e promuovono coagulazione e infiammazione, molti componenti di questo sistema inoltre promuovono il taglio proteolitico di C3 e C5, questo suggerisce un certo livello di collaborazione tra i due sistemi (Shishido, 2012). Molte lectine fanno parte di famiglie di PRR (Pattern Recognition Receptors) che riconoscono specifici zuccheri sulla superficie delle cellule e quindi diversi patogeni e ligandi intrinseci come cellule apoptotiche e componenti della matrice extracellulare (Bassi, 2009. Hirschfield, 2003). Queste proteine vengono riconosciute dai macrofagi ed altre cellule dell'immunità innata facilitando la clearance dei patogeni e dei detriti cellulari tramite fagocitosi (Shishido, 2012).

1.2 Sistema del complemento

Il sistema del complemento, pur essendo parte dell'immunità umorale innata, ha un ruolo chiave sia nella risposta immunitaria innata che in quella adattativa (Kohl et al., 2006). È formato da più di 50 glicoproteine solubili o di membrana che sono coinvolte in numerose interazioni proteina-proteina, causando l'assemblaggio ed attivazione di complessi

enzimatici e la generazione di frammenti bioattivi che iniziano svariate risposte cellulari mediante il legame di recettori del complemento e regolatori. (Ricklin et al., 2010). Il sistema del complemento contribuisce a numerose funzioni come fagocitosi, lisi cellulare, infiammazione, solubilizzazione degli immunocomplessi, rimozione delle cellule apoptotiche e stimolazione della risposta immune umorale (Shmidt, 2000).

Nei vertebrati sono conosciute 3 vie di attivazione: la classica, la via lectinica e la via alternativa. Tutte comprendono al loro interno una molecola capace di riconoscere e legare la superficie microbica ed una serina proteasi, che quando attivata, esegue il taglio della proteina C3 in C3a e C3b. C3 rappresenta l'elemento fondamentale del sistema del complemento dato che ha la capacità di esercitare un ruolo nell'opsonizzazione mediante legame di C3b sulla superficie microbica, ma può anche reclutare gli immunociti sul sito d'infezione tramite rilascio di C3a ed eventualmente portare alla lisi cellulare tramite formazione del Complesso di Attacco alla Membrana (MAC) (Song et al., 2000).

Nella via alternativa C3b richiede l'interazione con il fattore B (Bf), una proteasi, per formare la C3 convertasi che promuove la formazione di C3a e C3b. Il legame di C3b alla convertasi attivata causa un cambiamento nell'affinità per il substrato del complesso enzimatico che esegue, nei soli Vertebrati, il taglio di C5 in C5a e C5b che avrà un ruolo iniziale nell'assemblaggio del MAC. La via lectinica inizia con il riconoscimento di zuccheri su superfici cellulari da parte di lectine come le MBLs (Mannose-Binding Lectins), che portano all'attivazione di una serina proteasi associata alle MBLs (MASPs) la quale eseguirà il taglio di C3 in C3a e C3b (Wallis, 2007). La via classica, a differenza delle precedenti, inizia con lo specifico riconoscimento microbico da parte degli anticorpi che interagiranno con C1q che è associato a due serin-proteasi (C1r e C1s), questa via è esclusiva dei Vertebrati (Walport, 2001).

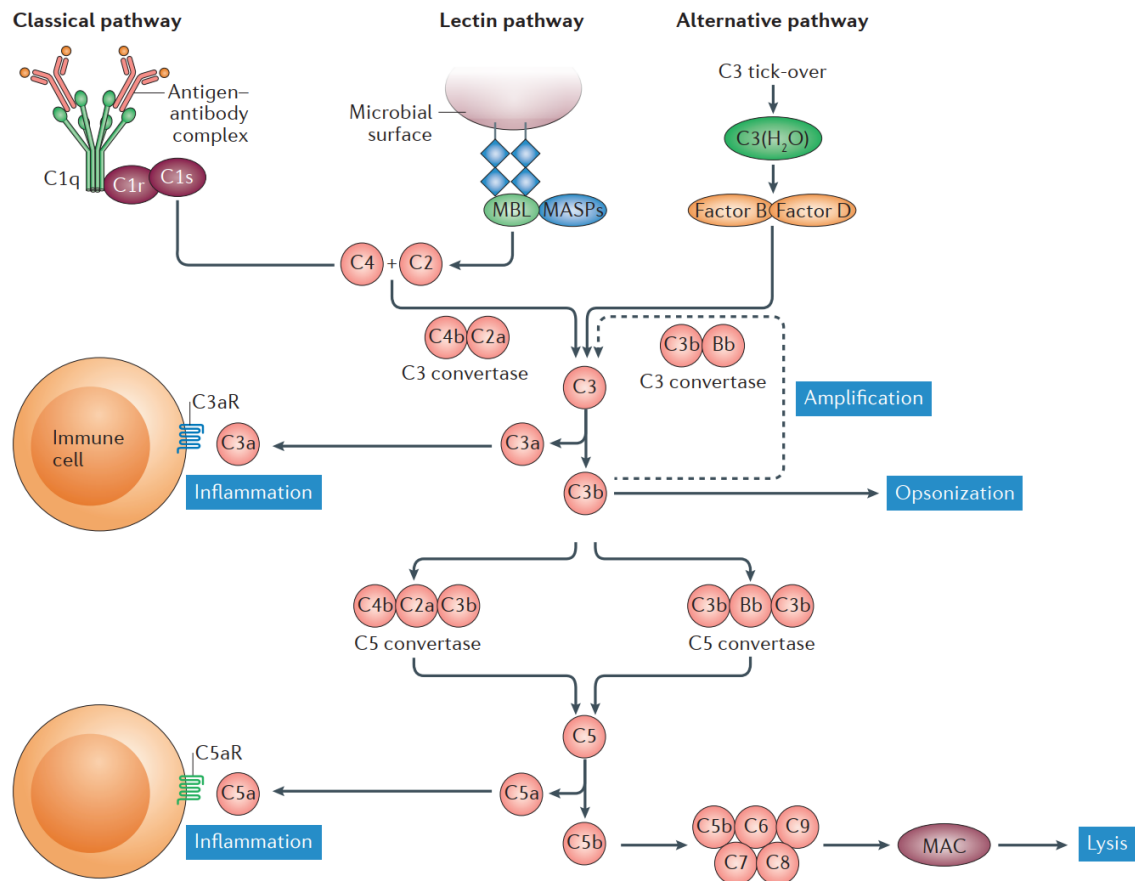


Figura 1 Sistema del complemento. Modificata da: Trouw, 2017.

L'attivazione del Sistema del Complemento può avvenire anche in maniera non canonica indipendente dalla convertasi. Proteasi del sistema cinina-callicreina, o delle cascate coagulativa e fibrinolitica, come plasmina e trombina, possono portare al clivaggio di C3, C5 o del fattore B, generando frammenti che possono assemblare una C3 convertasi funzionante o partecipare a risposte effettrici downstream (Irmscher et al., 2018. Amara et al., 2010) . Queste vie convertasi-indipendenti sono coinvolte in meccanismi di amplificazione dell'attivazione del complemento o agiscono separatamente dalle normali vie. È necessario specificare che la loro rilevanza fisiologica è dipendente dal contesto e dovrebbe essere considerata in relazione alle specifiche patologie (Nilsson et al., 2021).

1.2.1 Funzioni non immunitarie del complemento

Con l'avanzare della ricerca sono state evidenziate molte funzioni definite anche “non canoniche” per il sistema del Complemento tra le quali: controllo nella morfogenesi dei tessuti, riparazione delle ferite e pruning sinaptico. Come costante, la presenza di molecole

associate al complemento garantisce la prevenzione della risposta autoimmune, ad esempio, in spermatozoi e durante lo stadio di blastocisti. In particolare, negli spermatozoi, alcune molecole del Sistema del Complemento aiutano l'interazione oocita-spermatozoo e facilitano la fecondazione (Harris et al., 2006. Rooney et al., 1996), allo stadio di blastocisti contribuiscono allo sviluppo della massa cellulare interna e del trofoblasto (Taylor et al., 1996), durante l'organogenesi (più specificatamente durante la neurulazione) facilitano la chiusura del tubo neurale e contribuiscono alla migrazione delle cellule delle creste neurali (Danny et al., 2016), inoltre partecipano al pruning sinaptico ossia il rimodellamento di alcune connessioni sinaptiche per la maturazione dei circuiti neurali (Stevens et al., 2007) mentre nella formazione della placenta sembrano avere un compito di prevenzione contro eventuali patogeni (Stevens et al., 2007. Faulk et al., 1980.).

È Stato inoltre visto che C1q ed il Fattore H hanno ruoli nella proliferazione di cellule tumorali, immunosoppressione in ambito di microambiente tumorale, angiogenesi e progressione neoplastica. L'attività del complemento a livello intracellulare è inoltre stata collegata alla regolazione dell'autofagia e della secrezione dell'insulina a livello delle isole pancreatiche (King et al., 2023).

A livello di organismo adulto le proteine del complemento sono coinvolte in processi vitali delle cellule andando ad influenzarne proliferazione, differenziazione e collocazione (Hawksworth et al., 2018).

Sono sempre più numerose le prove che sottolineano come le proteine del complemento siano profondamente coinvolte nella regolazione di processi biologici che sono molto distanti dall'immunosorveglianza e dalla difesa dell'organismo (Mastellos et al., 2013. Stephan et al., 2012).

1.2.2 Sistema del Complemento negli invertebrati

La prima indagine circa il sistema del complemento a livello di invertebrati aveva portato all'ipotesi che potesse esservi un complemento primitivo o archeo-complemento auto-attivato dal "tick-over" di C3 o dal taglio proteolitico da parte di proteasi di patogeni (Lachmann, 1979). Questo è stato poi confermato da studi seguenti che hanno portato ad individuare una sequenza parziale C3-simile (Smith et al., 1996), seguita dal rispettivo cDNA completo (Al-Sharif et al., 1998) la cui espressione è stata localizzata in cellule correlate all'immunità: i celomociti di *Strongylocentrotus purpuratus* (Clow et al., 2000).

Dalle prime individuazioni di analoghi di C3 in invertebrati numerosi studi hanno poi portato ad identificarne la presenza nei principali phyla (Peng et al., 2016).

A seguito venne individuato anche un omologo del Fattore B (Smith et al., 1996) ricostruendone la relativa sequenza di cDNA completa, espressa anch'essa in celomociti (Terwilliger et al., 2004). Col progredire degli studi l'idea di archeo-complemento si è rivelata corretta ma il concetto è stato modificato ed espanso, aggiungendo una proposta di ulteriore via la cui attivazione è dipendente da lectine. Questo inoltre presenta un feedback d'amplificazione del segnale tipico della via alternativa corredato dalle proteine chiave della via come le proteasi che portano all'attivazione del Fattore B e aumentano l'attivazione di C3 (Smith et al., 1999, 2001). Risultano assenti invece le proteine C4 e C5 perché frutto di una duplicazione genica avvenuta nel clade dei Vertebrati (Franchi, non pubblicato).

1.2.3 Regolatori del complemento

C3 nella sua forma attiva (C3b) espone il sito tioestere che è in grado di legare covalentemente gruppi idrossilici ed amminici, potenzialmente presenti sulle cellule self, ed attivare la cascata del complemento. Perciò, a livello di vertebrati, le proteine regolatrici del complemento sono presenti sulla superficie cellulare e nel liquido extracellulare per proteggere l'organismo da un attacco autoimmune da parte della cascata del Complemento (Merle et al., 2015a). Tutte le proteine codificate dal cluster genico dei regolatori del complemento (RCA) in genomi di mammifero presentano una notevole somiglianza in struttura e composizione. Sono principalmente costituite da domini denominati Short Consensus Repeats (SCRs) o domini CCP (Complement Control Proteins) i quali presentano una struttura tridimensionale fortemente conservata tra loro (Krych-Goldberg and Atkinson, 2001). Sono varie le funzioni attribuite a differenti domini CCP di ogni proteina regolatrice, spaziano da attività inibitorie a legame con superfici cellulari o ad altre componenti del complemento. Questa organizzazione modulare dei geni permette di codificare per proteine ibride con diverse serie di domini CCP assolvendo a differenti funzioni (Pouw et al., 2015).

Informazioni circa i regolatori del complemento a livello degli invertebrati sono scarse. La ricerca di queste proteine nei non mammiferi ed invertebrati è particolarmente complessa a causa della modularità dei geni che porta ad una semplice riorganizzazione dei domini CCP e a causa della similarità nella struttura proteica che portano a difficoltà nel predire parentele con sequenze di mammifero. Questo potrebbe essere dovuto alla variabilità della sequenza insorta nel corso dell'evoluzione originatasi sia tramite: mutazioni puntiformi, duplicazioni

geniche, delezioni e/o riorganizzazione di domini ed esoni, in aggiunta a varianti di splicing portando a geni che codificano per proteine contenenti domini CCP multipli. Un'ulteriore problematica, che vale anche per tutte le altre proteine, è rappresentata dal fatto che non è possibile individuare similarità funzionali tra le proteine regolatrici del complemento tramite discendenza evolutiva o similarità di sequenza. I domini CCP potrebbero variare di numero all'interno e tra le diverse specie, e possono includere differenze dovute a splicing alternativo (Smith et al., 2023).

Ad oggi, nei Vertebrati, sono state individuate una serie di proteine con comprovata azione di regolazione del complemento: Decay Accelerating Factor (DAF), Membrane Cofactor Protein (MCP), Fattore H (FH) e C4b-Binding Protein (C4BP). Queste hanno una forte influenza sulla regolazione della localizzazione e deposizione del complemento (Medzhitov and Janeway, 2002). Data però la probabile ipervariabilità di questi trascritti, anche nei Vertebrati, vengono costantemente scoperte sempre più proteine con azione regolatrice del complemento caratterizzate dal pattern di domini CCP, tra cui CUB and SUSHI Multiple Domains 1 (CSMD1) (Escudero-Esparza et al., 2013; Håvik et al., 2011; Kraus et al., 2006). Per questa proteina, in particolare, è stato recentemente dimostrato un ruolo diretto nella gestione dei processi di eliminazione sinaptica durante lo sviluppo, un dato significativo verso la conclamata funzione del complemento in ambiti non immunitari (Baum et al., 2024).

1.3 *Pomacea canaliculata*

Pomacea canaliculata (*Pc*) è un mollusco gasteropode d'acqua dolce interessante da più punti di vista, soprattutto per la sua grande resistenza a condizioni di stress e spiccate capacità rigenerative nell'organismo adulto (Bever et al., 1988. Accorsi et al., 2017). L'elevata fecondità e capacità di adattamento di *Pc* ad un'ampia varietà di ambienti la rende una delle 100 specie più invasive secondo l'ISSG (Invasive Species Specialist Group) (View 100 of the world's worst invasive alien species, 2000). *Pc* è anche un ospite intermedio del nematode parassita umano *Angiostrongylus cantonensis* perciò è molto studiata per la sua eradicazione (Yang et al., 2013). Il sistema immunitario di *Pc* è composto da una componente cellulare, gli emociti, ed una componente umorale (Smith et al., 2016). Attualmente esistono informazioni circa la morfologia (Accorsi et al., 2013. Cueto, 2015) e la proliferazione (Accorsi, 2014) degli emociti ma sono pochi i dati molecolari.

2. Scopo della tesi

Numerosi studi hanno provato l'importanza del sistema del complemento ed i suoi regolatori in processi fisiologici immunitari e non, puntualizzando anche quanto alterazioni dei regolatori possano portare a condizioni patologiche. Con il seguente lavoro l'obiettivo è quello di valutare l'ipotesi che l'ipervariabilità dei regolatori del complemento sia dovuta a splicing alternativo. Disponendo poi delle strutture 3D delle proteine note aventi domini CCP e con sussidio di un'intelligenza artificiale (AI) cercare di raggruppare le proteine dedotte dai trascrittomi degli emociti e del rene posteriore di *P. canaliculata* secondo la loro struttura tridimensionale costituendo così un dataset da poter indagare per individuare nuovi regolatori potenzialmente composti anche da domini non convenzionali.

3. Risultati

3.1 Ipervariabilità da splicing alternativi

Sono state individuate all'interno dei trascrittomi di *P. canaliculata* ottenuti con sequenziamento Nanopore 11 sequenze putative per Regolatori del Complemento contenenti CCP, 3 di queste appartenenti al trascrittoma di emociti e le restanti 8 del rene posteriore (Figura 2) dimostrando che i Regolatori in *Pomacea*, almeno per i tessuti presi in esame, risultano essere tessuto-specifici.

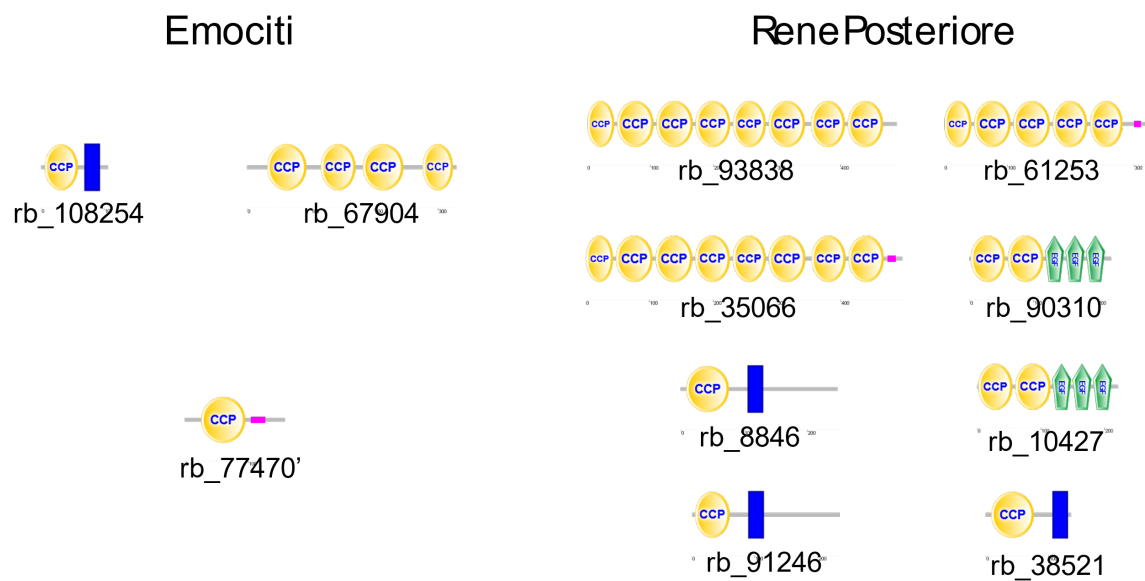


Figura 2 Sequenze putative per Regolatori del Complemento individuate dall'analisi BLAST sui rispettivi trascrittomi.

Allineando queste sequenze con il genoma di *P. canaliculata* depositato in GeneBank sono riuscito ad individuare due geni specifici all'interno delle regioni *SZHN201 Linkage group LG10 ASM307304v1* e *SZHN201 Linkage group LG2 ASM307304v1*. I trascritti di interesse dal rene posteriore risultano generati da splicing alternativo da due geni: PcPKRCA1 e PcPKRCA2. In particolare, sono state individuate 5 varianti di splicing dal gene PcPKRCA1 e 3 varianti di splicing dal gene PcPKRCA2 (Figura 3). L'analisi dei singoli domini CCP delle diverse varianti di splicing conferma questa constatazione (Allegato 1).

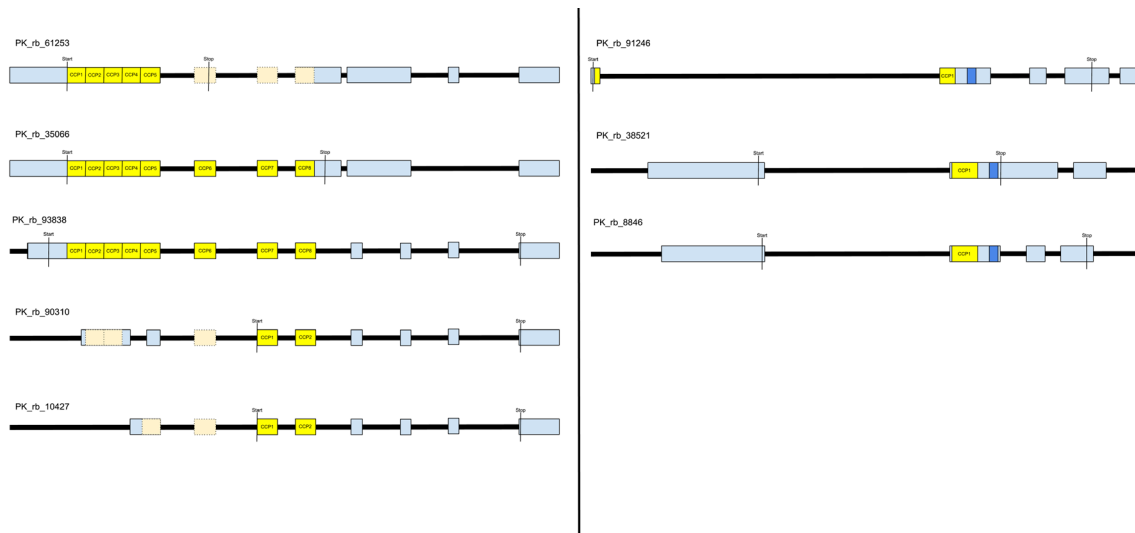


Figura 3 Rappresentazione grafica dei geni e splicing.

Nel caso del trascrittoma di emociti non è stato possibile indagare in modo specifico la presenza di varianti di splicing poiché è stata individuata solo una regione genomica compatibile con uno dei tre trascritti: *SZHN2017 linkage group LG2, ASM307304v1* (Figura 4). Gli altri due trascritti non risultano come splicing alternativo di questa regione genomica.

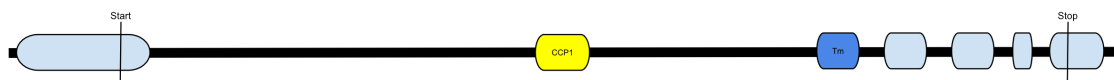


Figura 4 Unica sequenza individuata dal confronto con i database nel trascrittoma di Emociti

3.2 Analisi 3D da algoritmo AI

L'output del modello ProstT5 è stato rappresentato, dopo accurata riduzione dimensionale con metodo Umap, in un grafico bidimensionale a dispersione dove ogni punto è assegnato ad un cluster e la maggior vicinanza tra punti indica una maggiore similarità in termini di significato biologico (Figura 5). Ho quindi identificato il cluster di appartenenza delle sequenze con CCP note individuandole nei cluster 407 e 362, questi risultano vicini tra loro garantendoci un maggior grado di sicurezza nel proseguire con l'analisi.

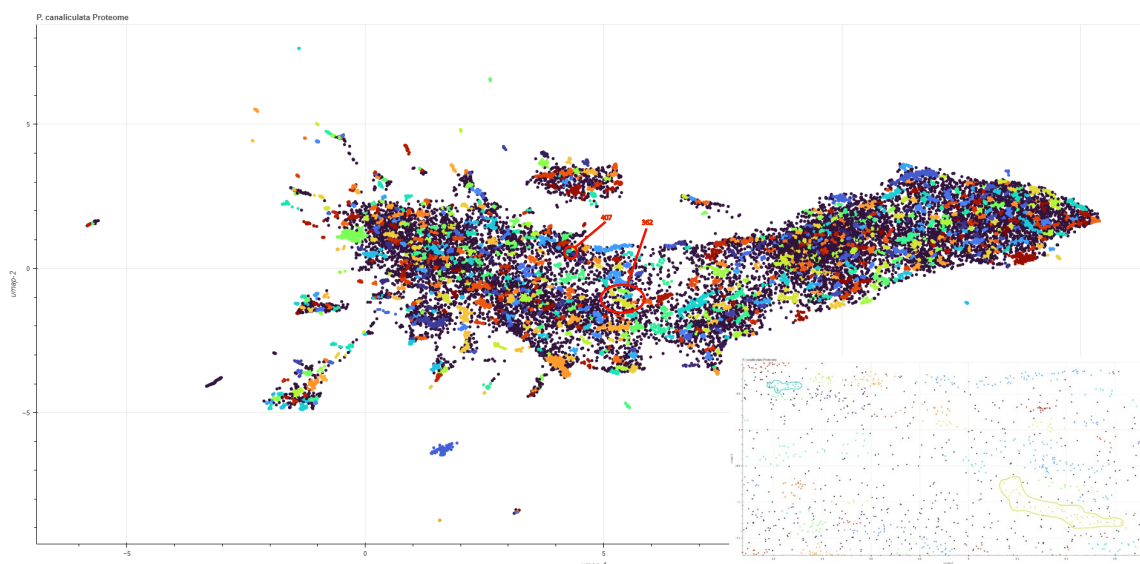


Figura 5 rappresentazione 2D dei dati ricavati tramite approccio AI. Evidenziati i cluster 407 e 362 contenenti le sequenze CCP note.

E' stata quindi eseguita un analisi Blast per verificare l'identità delle sequenze contenute nei due cluster, seguita poi da un'analisi SMART per determinarne l'organizzazione in termini di domini proteici. Questo ci ha portato a costituire un dataset con 14 sequenze appartenenti al cluster 407 e 36 appartenenti al cluster 362 (Allegato 2). Osservando i dati è stato individuato che, la maggior parte delle sequenze non restituisce alcun dominio noto dall'analisi SMART e per una buona parte di queste non è stata possibile l'identificazione con Blast perché non sono presenti sequenze compatibili all'interno dei database (Tabella 1).

Tabella 1 Analisi del dataset

	Cluster 407	Cluster 362
Sequenze identificate	14	36
"Uncharacterized"	4 (28.57%)	16 (44.44%)
No domini noti con SMART	7 (50%)	27 (75%)

Avendo ottenuto un riscontro positivo dal primo approccio con il modello AI, sono state inserite nel trascrittoma delle sequenze umane notoriamente appartenenti ai Regolatori del Complemento (CFAH, CD55, TNF) che potessero funzionare da controllo interno. Questo ci ha permesso di costituire un secondo dataset estremamente più abbondante (Allegato 3). Analizzando i dati abbiamo individuato che le sequenze note con CCP sono localizzate nei cluster 231 e 232, estremamente vicine tra di loro, mentre quelle umane in 207, 308 e 374.

4 Discussione

I regolatori del complemento presentano la caratteristica conservata di essere costituiti da ripetizioni di domini CCP e la diversa disposizione di questi ultimi determina alterazioni funzionali (Pouw et al., 2015). È nota la difficoltà nell'analisi di sequenze che codificano per regolatori del complemento soprattutto per il fatto che la tecnologia di sequenziamento più utilizzata ad oggi è basata sulla piattaforma Illumina che genera delle reads molto brevi che necessitano di essere poi assemblate sul genoma per formare cDNA completi. Questa metodica risulta inefficace per valutare sequenze ipervariabili come appunto le varianti di splicing (Smith et al., 2023). Il nostro approccio è stato quindi quello di impiegare la tecnologia di sequenziamento Nanopore (*MinION mk1*) che permette di generare reads molto più lunghe che possano essere assemblate facilmente senza l'ausilio di sequenze genomiche di riferimento. In questo modo emergono chiaramente i reali trascritti e quindi tutte le possibili varianti dei diversi geni. Con questo approccio sono stati costituiti due trascrittomi: il primo da rene posteriore ed il secondo da emociti derivati dallo stesso individuo. Abbiamo effettuato un'analisi Blast utilizzando come chiave di ricerca delle sequenze da *P. canaliculata* contenenti domini CCP e precedentemente individuate nei database pubblici, questo ci ha portato a identificare una serie di sequenze con un ottimo grado di similarità. Abbiamo proseguito analizzando le sequenze restituite dall'analisi per individuare quali potessero codificare per proteine contenenti domini CCP. Per questo sono stati valutati due approcci differenti: il primo metodo si serviva del software InterPro Scan, uno strumento bioinformatico capace di identificare domini proteici, famiglie e caratteristiche funzionali partendo da sequenze nucleotidiche o amminoacidiche tramite confronto della sequenza indagata con vari database di famiglie proteiche e modelli di domini. Questo metodo si è rivelato estremamente esoso in termini di tempo e risorse computazionali e non ha restituito dei buoni risultati. Siamo quindi passati al secondo approccio tramite SMART, un database e strumento bioinformatico in grado di identificare domini proteici e motivi modulari maggiormente incentrato verso la comprensione funzionale delle proteine prese in analisi, questo ha restituito 11 sequenze con caratteristiche tipiche dei Regolatori del Complemento (Figura 2) e si è rivelato il metodo migliore.

Numerosi studi hanno evidenziato la presenza di fenomeni di splicing alternativo in alcune delle principali proteine regolatrici del complemento a livello di vertebrati (Post et al., 1991, Russel et al., 1992, Mannes et al., 2020), ma sono poche ad oggi le informazioni circa questo meccanismo a livello degli invertebrati. Ottenute quindi le sequenze abbiamo mappato

queste sul genoma pubblicato in NCBI di *Pomacea canaliculata* descrivendo la composizione in introni ed esoni dei singoli geni. Quello che è emerso immediatamente è che, almeno per i tessuti presi in esame in questo studio, i Regolatori del Complemento di *P. canaliculata* risultano essere tessuto-specifici: 3 trascritti unici negli emociti ed 8 nel rene posteriore. Per quanto riguarda il rene posteriore le sequenze individuate sono risultate appartenere a due geni distinti: PcPKRCA1 e PcPKRCA2. Analizzando le sequenze e creando delle rappresentazioni grafiche di questi trascritti abbiamo evidenziato come questi siano verosimilmente originati da splicing alternativo (Figura 2). A conferma di questo abbiamo generato un albero che raggruppa le sequenze dei diversi CCP sulla base della loro similarità evidenziando come gli stessi CCP si possano trovare in alcune delle varianti trascrizionali sottolineandone la provenienza dal medesimo esone (Allegato 1).

Fino ad oggi, per condurre ricerche di sequenze specifiche su invertebrati, si è sempre utilizzato l'approccio dell'omologia o dissimilarità di sequenza con trascritti dai Vertebrati. Per quanto questo metodo si sia rivelato esaustivo in molti settori rimane sempre vincolato alla ricerca di qualcosa che deve essersi conservato nel corso dell'evoluzione appunto fino ai Vertebrati ed escludendo la possibilità di trovare qualcosa di nuovo. Oggi si apre una nuova frontiera di calcolo computazionale con algoritmi che possono andare molto oltre il semplice confronto di sequenze lineari. Questi nuovi approcci si basano sulla capacità di *machine learning* ed Intelligenza Artificiale (AI).

L'intelligenza artificiale è diventata uno strumento essenziale per l'elaborazione di grandi quantità di dati. Grazie al machine learning, le AI sono capaci di analizzare in modo molto rapido dataset complessi portando all'individuazione di pattern ed informazioni che difficilmente potrebbero essere colte in altro modo (Russell, 2021). L'adattamento di queste tecniche in ambito proteomico tramite sostituzione delle parole con amminoacidi e frasi con sequenze proteiche ha aperto la strada a nuovi e potenti strumenti bioinformatici per la modellizzazione di sequenze proteiche detti *protein Language Models* (pLMs) (Heinzinger et al., 2019; Chen et al., 2023). Il principale, AlphaFold2, ha reso disponibili le previsioni delle strutture 3D di più di 200 milioni di proteine al pubblico. Questo grande database ha consentito di allenare un metodo denominato ProstT5 che, attribuendo un vettore numerico di dimensione 512/1024/2048 ad ogni amminoacido, riesce ad attribuirvi un significato biologico e ricostruire la struttura 3D della sequenza indagata (Heinzinger et al., 2023). Sono stati quindi uniti ed analizzati i due trascrittomi con il pLM ProstT5. Sono stati ottenuti una serie di vettori ciascuno ascrivibile ad una diversa proteina. Sono poi stati rappresentati

graficamente sfruttando la tecnica Umap (*Uniform Manifold Approximation and Projection for Dimension Reduction*) che ha permesso di ridurre l'informazione in un piano bidimensionale dato che il vettore costituisce un dato di tante dimensioni quanti sono i valori di cui è formato. È stato ottenuto quindi un grafico a dispersione dove i punti rappresentano le singole proteine e la loro vicinanza indica un simile significato biologico (Figura 5). Questo approccio porta ad una perdita sostanziosa di informazioni data la riduzione di dimensionalità, ma rappresenta comunque il metodo più idoneo in quanto non è possibile impiegare la PCA (Principal Component Analysis) visto che disponiamo di vettori con molti caratteri. Il metodo ci ha consentito di individuare due cluster dove ricadono le sequenze note contenenti CCP. Il contenuto è stato analizzato prima con Blast per stabilire l'identità delle sequenze e poi con SMART ottenendo un risultato estremamente interessante poiché vengono raggruppate insieme principalmente proteine con risaputa struttura da Regolatore del Complemento, ma anche proteine con predizione di similarità ignota o proteine in cui non si riescono a riconoscere domini proteici. Non va dimenticato che il riconoscimento di un dominio specifico è possibile solo nel caso in cui quel dominio sia stato registrato e gli sia stata data una sequenza consenso nei diversi database. L'assenza di dominio, quindi, può significare una porzione di proteina di collegamento, senza funzione specifica o la presenza di una proteina la cui funzione è effettivamente ignota. È quest'ultima ipotesi che rende interessante il nostro dato allargando la possibile indagine non solo a proteine con domini CCP, ma anche ad altre sequenze. Chiaramente solo studi futuri potranno far luce se queste proteine ignote e se queste abbiano effettivamente un ruolo nella regolazione del Complemento, ma risulta quanto mai interessante riuscire a poter ipotizzare una funzione per gruppi di proteine che non hanno alcuna collocazione funzionale nei database. In seguito, a scopo di controllo, è stata eseguita nuovamente l'analisi utilizzando lo stesso trascrittoma ma con l'aggiunta di sequenze umane che fanno riferimento a regolatori del complemento noti (CFAH, CD55, TNF) in modo da valutare come venisse influenzata la clusterizzazione. Analizzando l'output è stato visto che le sequenze di *pomacea* hanno mantenuto la clusterizzazione indipendente da quelle umane, questo è con molta probabilità dovuto alla marcata differenza delle proteine tra invertebrati e vertebrati. Inoltre, andando ad analizzare il contenuto dei singoli cluster si riscontrano parecchie somiglianze strutturali tra le sequenze comprese nei singoli raggruppamenti a rafforzare l'efficacia del metodo d'indagine. Questa seconda analisi ha in ogni caso permesso di ottenere un altro abbondante dataset, contenente molte sequenze “*Uncharacterized*” e con previsione dei domini ignota (Allegato 3), che si presta ad ulteriori e più approfondite future analisi.

5 Conclusioni

Lo studio ha individuato 11 sequenze che potrebbero essere dei regolatori del complemento in *P. canaliculata*. Questi, almeno per i tessuti presi in analisi, hanno un'espressione tessuto-specifica e che la loro ipervariabilità è definita da fenomeni di splicing alternativo a partire da, almeno, 3 geni distinti.

Lo studio con AI, tramite applicazione del pLM ProstT5, ha portato alla generazione di due dataset molto popolati. Questi sono un importante risultato nell'ambito dello studio di regolatori del complemento in invertebrati data la scarsità di informazioni e dati a riguardo. Si tratta quindi di uno studio preliminare che ha gettato le basi per l'indagine dei regolatori del complemento in *P. canaliculata*.

6 Materiali e metodi

6.1 Animali

Le specie di *P.canaliculata* sono state acquistate da T.A.F. Trans Aquarium Fish S.r.l. (Torino, Italia) e mantenute in acquari riempiti d'acqua areata a 23.1°C. Organismi adulti ed attivi nella riproduzione (diametro della conchiglia 35-55 mm) sono stati mantenuti in un acquario separato. Gli animali sono stati nutriti *ad libitum* con verdure e frutta. Il 10-20% del contenuto dell'acquario è stato regolarmente sostituito. Il benessere degli animali è stato accertato dalla continua valutazione della riproduzione in acquario. Per i prelievi di campioni gli animali sono stati lasciati in una vasca separata per 24 ore prima di proseguire al prelievo.

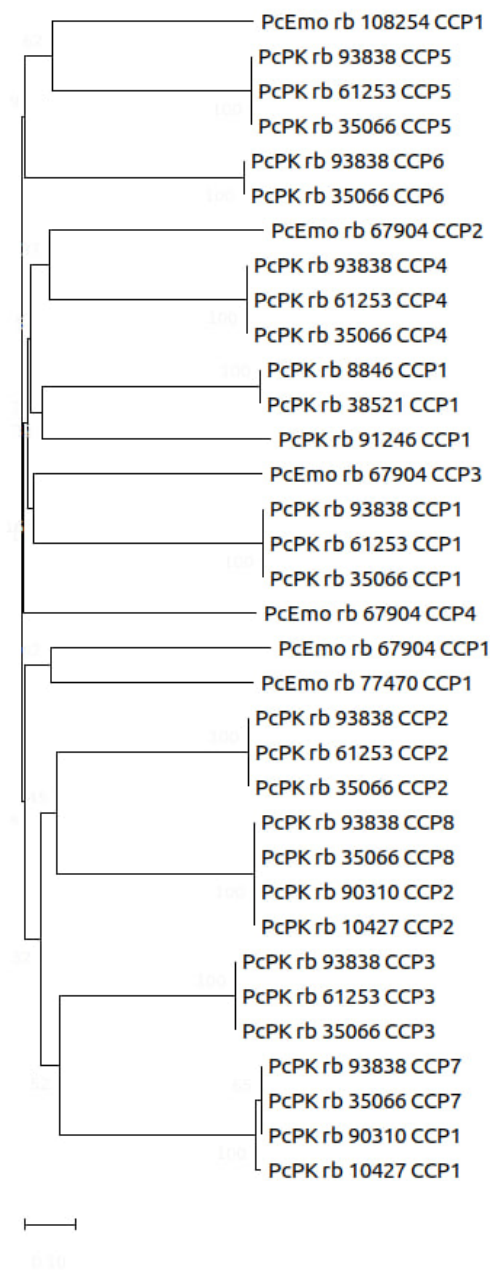
6.2 Software e Strumenti bioinformatici

Blast	Allineamento sequenze.	https://blast.ncbi.nlm.nih.gov/Blast.cgi
InterPro Scan	Caratterizzazione funzionale di sequenze proteiche e nucleotidiche.	https://www.ebi.ac.uk/interpro/about/interproscan/
SMART	Identificazione ed annotazione di domini proteici	https://smart.embl.de/
MEGA	Suite di strumenti per l'analisi di sequenze proteiche e nucleotidiche	https://www.megasoftware.net/

ProstT5	pLM capace di tradurre tra sequenza amminoacidica e struttura proteica	https://github.com/mheinzinger/ProstT5
Umap	Tecnica di riduzione dimensionale dei dati	https://umap-learn.readthedocs.io/en/latest/
Dorado	Basecaller ufficiale di Nanopore	https://github.com/nanoporetech/dorado
Pychopper	Strumento capace di identificare, orientare e tagliare le reads cDNA.	https://github.com/epi2me-labs/pychopper
RNAbloom	Assemblatore trascrittomico	https://github.com/bcgsc/RNA-Bloom

Allegati

Allegato 1 Rappresentazione grafica della similarità di sequenza dei domini CCP individuati



Allegato 2 Dataset analisi AI senza sequenze umane

Cluster 402	Query Sequence	Query SMART Analysis	Best Hit Description	Best Hit Sequence	Best Hit SMART Analysis
PomcanEvm0027031	MAVWILWLGVALF		nectin-1 precursor [Rattus norvegicus]	MARMGLAGAAGRW	
PomcanEvm0015671	MARGLVTAVGRLTC		uncharacterized protein LOC112568583[Pomacea canaliculata]	MARGLVTAVGRLTCH	
PomcanEvm0013521	MSKIGRLMLFFRAT		limbic system-associated membrane protein isoform 1 preproprotein [Homo sapiens]	MVRVQDPDRKQLPLV	
PomcanEvm0010621	GSQIFSDNRGYRLD		sushi, von Willebrand factor type A, EGF and pentraxin domain-containing protein 1	MWVPLAFQWGLALV	
PomcanEvm0198791	MPYDDGVVSDVDR		low-density lipoprotein receptor-related protein 1B precursor [Mus musculus]	MSQLLIALLTSLQLFN	
PomcanEvm0137641	MRWLLCLAYFFACL		DUF1553 domain-containing protein [Aporhodoirellula aestuarii]	MHRSLJALVWLAFA	
PomcanEvm0035501	MCKASNSYTHRGCP		hemocentin-2 precursor [Mus musculus]	MTPGAQLRLPLVAISTA	
PomcanEvm0014801	MFLIACHRTLAIRMIS		uncharacterized protein LOC112568581 isoform X1 [Pomacea canaliculata]	MATSPDSLSTNLIVK	
PomcanEvm0046581	MPRMFAVAVITSESL		Ig-like domain-containing protein [Priestia flexa]	MPYNFRFLGNANGAITI	
PomcanEvm0059701	MHQISFDIVNTNTEVT		nuclear pore membrane glycoprotein 210-like [Pomacea canaliculata]	MLSSLKDVAYKGVTPS	
PomcanEvm0424971	MVDITVIDNHKQFTI		uncharacterized protein LOC112571911[Pomacea canaliculata]	MSTOGQHLYKQTKTM	
PomcanEvm0015311	METWKSMLLCTVYV		uncharacterized protein LOC112568136 isoform X2 [Pomacea canaliculata]	MBCKSMVLCTVYVHLL	
PomcanEvm0008011	MRRIITVTGVMGCV		adhesion G protein-coupled receptor E2 precursor [Rattus norvegicus]	MTPACRLLSMLSLRLA	
PomcanEvm0142581	MVQNAIGTCENNTC		teneurin-3 isoform 1 [Homo sapiens]	MDVKERRFYCSLTKSR	

Cluster 362	Query Sequence	Query SMART Analysis	Best Hit Description	Best Hit Sequence	Best Hit SMART Analysis
PomcanEvm0071521	MTPSRFRSTLSS		heat shock-related 70 kDa protein 2 [Rattus norvegicus]	MSARGPAIGDLGTYS	
PomcanEvm0727371	MRSPFYLRLGMLKIF		GTP-binding protein 1 [Mus musculus]	MAABRFRSPVDSPPA	
PomcanEvm0198871	MFESRNIFGLYAVQA		actin-related protein 3C isoform a precursor [Homo sapiens]	MFESRNIFGLYAVQA	
PomcanEvm0477901	MNRFGLKAVAVQRI		E3 ubiquitin-protein ligase HEW2-like [Pomacea canaliculata]	MVKAKMSRNSDYSS	
PomcanEvm0074921	MSRRLQCAVDVTH		spermatogenesis-associated 6-like protein [Rattus norvegicus]	MPLEVVLEQIRASCP	
PomcanEvm0079091	MTAKLYKPLFAVKIL		benzoate transporter [Chloroflexus aurantiacus]	MNRRPFRALSTLAVIL	
PomcanEvm0040141	MLRWARHVPYKRS		hypothetical protein CQ70_17910[Pomacea canaliculata]	MDWESSSCOEYETIKV	
PomcanEvm0052581	MRNPATDVGVD		uncharacterized protein PB18B9.04c-like isoform X2 [Pomacea canaliculata]	MBWKGGLRPVAVNI	
PomcanEvm0109541	MRFSVJLTLFAV		hypothetical protein I4U23_029159[Adineta vaga]	MRLSLVLLCFTLVLVN	
PomcanEvm0180211	MSASRPYPTELTFS		hypothetical protein SAR_C06627[Sphaerofarma arctica JP610]	MKVAAALVLFVGSQS	
PomcanEvm0081821	MARASECALVFWV		fibronectin type III domain-containing protein [Aquisphaera giovannonii]	MEDLEPRLVMSGAVLT	
PomcanEvm0167681	MPBHQEKIYTABKIL		E3 SUMO-protein ligase CBX4 [Mus musculus]	MELPAVGBHFAVESEI	
PomcanEvm0049491	MPGLSSLINGLWM		uncharacterized protein LOC112559808 isoform X1 [Pomacea canaliculata]	MPGLSSLINGLWMRF	
PomcanEvm0049501	MPGLSSLINGLWM		uncharacterized protein LOC112559808 isoform X1 [Pomacea canaliculata]	MPGLSSLINGLWMRF	
PomcanEvm0041521	MMKISVRLLVFTFM		uncharacterized protein LOC112559049 isoform X1 [Pomacea canaliculata]	MMKISVRLLVFTFLC	
PomcanEvm0242131	MKSPFEEBPREPE		hypothetical protein QJ45_000725[Haematococcus lacustris]	MSPEVMHVDALWT	
PomcanEvm0044681	MVLYLGLPDKKRTG		hypothetical protein CQ70_06994[Pomacea canaliculata]	MOGLASVSESVAIVF	
PomcanEvm0053931	DCRNWQNCNHLI		uncharacterized protein LOC112568838 isoform X1 [Pomacea canaliculata]	MRGLIITVGVLLTVC	
PomcanEvm0073691	KTIRRLNLTYPPIKRF		uncharacterized protein LOC112568756 isoform X2 [Pomacea canaliculata]	MLTVRKVTRADNGTIF	
PomcanEvm0114521	MAKRKGAGGGGAL		uncharacterized protein LOC112574796[Pomacea canaliculata]	MVKRKAGAGGGGSDS	
PomcanEvm0083581	MPKECAAPDCQKGG		THAP domain-containing protein 6 [Rattus norvegicus]	MVKCSAIGCASRCLP	
PomcanEvm0050071	MSPVRIPGLTSALVC		uncharacterized protein LOC112569258[Pomacea canaliculata]	MSPVQLPGLTSALVQL	
PomcanEvm0046181	MEANYFMLTLGLFL		uncharacterized protein LOC112571572[Pomacea canaliculata]	MEANYFMLTLGLFL	
PomcanEvm0112921	MQFFKSDPPQDLI		leucine-rich repeats and immunoglobulin-like domains protein 2 isoform X1 [Pomacea canaliculata]	MKITYVTLTALIGFMW	
PomcanEvm0132031	MIAMKABKQSEIRLP		uncharacterized protein LOC112556833 isoform X1 [Pomacea canaliculata]	MAAKATVQCPLDASA	
PomcanEvm0050701	TYISQVGLPAPWR		uncharacterized protein LOC112568447[Pomacea canaliculata]	MARDPKASKQKRFRI	
PomcanEvm0173291	MBSFYRFLRLTAV		pituitary tumor-transforming gene 1 protein-interacting protein precursor [Mus musculus]	MAPANLGLTTHWMM	
PomcanEvm0093781	MGCNFQCDNSLTJM		uncharacterized protein LOC112568712[Pomacea canaliculata]	MGCNFQCDSSKVMNF	
PomcanEvm0086141	MPFGKPKCVSKVTS		GATA zinc finger domain-containing protein 1 [Homo sapiens]	MPLGLKPTCSVCKTTS	
PomcanEvm0058571	MAWQLLLVLLLVST		uncharacterized protein LOC112571063 isoform X3 [Pomacea canaliculata]	MAWQLLLVLLLVSTA	
PomcanEvm0083881	METLRNMLLRIRYKI		uncharacterized protein LOC112575090 isoform X1 [Pomacea canaliculata]	MGGNKVGVSSAWETL	
PomcanEvm0054141	MGNKRTSNACIGFIF		uncharacterized protein LOC112568641[Pomacea canaliculata]	MKLRPQTAKPDDQK	
PomcanEvm0054581	MSHKSGSETCADLE		hypothetical protein CQ70_12045[Pomacea canaliculata]	MTHVEFAGYKDGDTVI	
PomcanEvm0057501	MRVREBLLLVVTV		uncharacterized protein LOC112573282 isoform X1 [Pomacea canaliculata]	MRVREBLLLVVTV	
PomcanEvm0043831	MAQRITWQQTTLTL		mucin-5B precursor [Homo sapiens]	MGAFSACRTLVLAALAA	
PomcanEvm0091211	MGRVGVFGGFLN		hypothetical protein CQ70_20392[Pomacea canaliculata]	MPFVWSKSKFEMPS	

Allegato 3 Dataset analisi AI con sequenze umane di controllo

231	Query Sequence	Query SMART Analysis	Best Hit Description	Best Hit Sequence	Best Hit SMART Analysis
PomcanEvM0036741	MNENTSTSGITGEC		uncharacterized protein LOC112574294 [Pomacea canaliculata]	SRQHNSYLKDNLDI	
PomcanEvM0041261	MENDTNAVWNIPG		uncharacterized protein LOC112558992 [Pomacea canaliculata] >ref XP_0251	MENDTNAVWNIPG	
PomcanEvM0118501	MVKNSPWICMMLL		uncharacterized protein LOC112559493 [Pomacea canaliculata] >ref XP_0251	MVKNSPWICMMLL	
PomcanEvM0146161	MVLSLRSPDVLCLG		hypothetical protein BaRGS_00003034 [Batillaria attramentaria]	LQTDFFVSSAQVSA	
PomcanEvM0071521	MTPRSRRRTSTLSS		hsp70-like protein [Diploporaceae sp. PMI_573]	PKIQKLLSDFDGGK	
PomcanEvM0114691	MSSVIFVTIDWLKIV		uncharacterized protein LOC112559256 isoform X1 [Pomacea canaliculata]	MSSVIFVTIDWLKLL	
PomcanEvM0063471	MEQGTTFKIVMSN		-	-	-
PomcanEvM0198871	MFESFNIPGLYIAVQ		prolyl oligopeptidase (S09 family) [Schistosoma mansoni]	MFETFNIPGLYIAVQ	
PomcanEvM0074001	METLANLTPTDIENC		unnamed protein product [Rotaria sordida] >emb CAF3688405.1 unnamed p	MEKLANLTQNDINN	
PomcanEvM0477901	MNRFLGKAVVAQVQ		hypothetical protein C0Q70_10123 [Pomacea canaliculata]	MNRFLGKAVVAQVR	
PomcanEvM0194801	MLSFDLHKLTCILMI		uncharacterized protein LOC112568809 [Pomacea canaliculata]	MLSFDLHKLTCILMI	
PomcanEvM0074921	MSRRALRCVADVITI		spermatogenesis-associated protein 6-like isoform X3 [Pomacea canaliculata]	MSRRALRCVADVITI	
PomcanEvM0049191	MGDAATCCTVETPC		uncharacterized protein LOC112576411 [Pomacea canaliculata] >ref XP_0251	MGDAATCCTVETPC	
PomcanEvM0396451	MGQVGEWSCSTCS		thrombospondin type-1 domain-containing protein [Chitinophages bacterium VGPWSPWACDTS	MGQVGEWSCSTCS	
PomcanEvM0066361	MNSFEHSTRELCLL		uncharacterized protein LOC112569225 isoform X3 [Pomacea canaliculata]	MNSLGHSAARELL	
PomcanEvM0144781	MMTKLLFLLTVVCI		unnamed protein product [Rotaria sp. Silwood1] >emb CAF3540664.1 unnar	KVFFIFFATVFLPCT-	
PomcanEvM0400491	SLSILRHFIPRSRT		unnamed protein product [Adineta ricciae]	TISILFFGRIPSYSR	
PomcanEvM0040141	MLRWARHVPYRK		hypothetical protein C0Q70_17910 [Pomacea canaliculata]	LNLYTLVWRNDSL	
PomcanEvM0218051	MKEKSFFCVFFTIK		unnamed protein product [Rotaria sp. Silwood2]	KNKPLCIPYFTIKHY	
PomcanEvM0099091	MADIQNLFLKELKSG		hypothetical protein I4U23_026529 [Adineta vaga]	SRSPIRSSPVYPNEN	
PomcanEvM0175241	MELDEEDERKEDLL		MYCBP-associated protein-like isoform X1 [Pomacea canaliculata]	MELDEEDERKEDLLA	
PomcanEvM0143941	MVNEEVGTAEIDINT		unnamed protein product [Moneuplotes crassus]	FSNDEVGSQTLDINT	
PomcanEvM0236981	MKRKFLTVVIFSLIC		unnamed protein product [Adineta steineri]	KNYLIVVITIALFSTV	
PomcanEvM0311991	GAYYNLPRDDKTLI		unnamed protein product [Rotaria sp. Silwood1] >emb CAF1653488.1 unnar	LKSILPLISQRSNCT	
PomcanEvM0154911	MVLFINVYVGQRVE		hypothetical protein C0Q70_20847 [Pomacea canaliculata]	MVLFINVYVGQRVEV	
PomcanEvM0172141	MEGVAQSDSDPEI		uncharacterized protein LOC112573858 [Pomacea canaliculata]	MEGVAQSDSDPEH	
PomcanEvM0052581	MRYNPTAPDVGVGF		hypothetical protein C0Q70_05358 [Pomacea canaliculata]	MRYNPTAPDVGVGF	
PomcanEvM0076571	MVVVCSVSEREKIRR		uncharacterized protein LOC112569182 [Pomacea canaliculata] >ref XP_0251	LCLGVFLWLMPLA	
PomcanEvM0109541	MRLSLVLTLLFAV		hypothetical protein I4U23_029159 [Adineta vaga]	MRLSLVLTLLFAV	
PomcanEvM0170241	MPPLRQIPDSIKALS		LEM domain-containing protein 1-like [Pomacea canaliculata]	MPPLRQIPDSIKALS	
PomcanEvM0180211	MSASRPYPPTLFT		hypothetical protein SARC_06627 [Sphaerofoma arctica JP610] >gb KNC810	PPYTHSVIFQAPCFIS	
PomcanEvM0081821	MARASECALVFVVV		uncharacterized protein LOC112567540 [Pomacea canaliculata]	MARASECALVFVVV	
PomcanEvM0167681	MPHEQAEAKYTAEKI		unnamed protein product [Rotaria socialis] >emb CAF3247914.1 unnamed p	MPDDKEAPVYTAEKI	
PomcanEvM0086071	MLEDEEDTACEVC		PHD and RING finger domain-containing protein 1-like [Pomacea canaliculata]	MLEDEEDTACEVC	
PomcanEvM0049491	MPGILSSLLKGLWM		uncharacterized protein LOC112559808 isoform X1 [Pomacea canaliculata]	MPGILSSLLKGLWM	
PomcanEvM0045461	MTEGVDTDDILNVS		uncharacterized protein LOC112574827 [Pomacea canaliculata]	MTEGVDTDDILNNG	
PomcanEvM0116181	MVATRVRLDALVLA		uncharacterized protein LOC112559070 isoform X3 [Pomacea canaliculata]	MMATRVRLDALVLA	
PomcanEvM0234681	MGTFIQFNKITGNP		-	-	-
PomcanEvM0049501	MPGILSSLLKGLWM		uncharacterized protein LOC112559808 isoform X1 [Pomacea canaliculata]	MPGILSSLLKGLWM	
PomcanEvM0068421	MIRTEGQSLFLFW		uncharacterized protein LOC112567779 [Pomacea canaliculata] >ref XP_0251	MIRSGSFCDSTAGY	
PomcanEvM0229331	MAKTFLPYFVSLTI		transmembrane protein 9-like [Pomacea canaliculata]	QASGFDRCKCVCP	
PomcanEvM0054621	MTEYCTVFTVEIDS		uncharacterized protein LOC112564904 isoform X1 [Pomacea canaliculata] >	MTEYCTVFTVEIDS	
PomcanEvM0060201	MTWREGHRIVDMC		zinc finger protein 62 homolog [Pomacea canaliculata]	MTWREGHRIVDMC	
PomcanEvM0041521	MMKISVFLLVFTFE		uncharacterized protein LOC112559049 isoform X1 [Pomacea canaliculata] >	MMKISVFLLVFTFE	
PomcanEvM0242131	MKSPREEVPREPPE		hypothetical protein QJ45_000725 [Haematooccus lacustris]	LSKEDDVRKYVNTYR	
PomcanEvM0327661	MMFVGMCMCNKRRL		uncharacterized protein LOC112553726 [Pomacea canaliculata] >ref XP_0251	MMFVGMCMCNKRRL	
PomcanEvM0201661	IEYEEVIRDIHYHN		LOW QUALITY PROTEIN: calyntenin-1-like [Pomacea canaliculata]	IEYEEVIRDIHYHN	
PomcanEvM0109991	MFIMKITVLGRYFIV		hypothetical protein C0Q70_08092 [Pomacea canaliculata]	MCHEASSNMSRCC	
PomcanEvM0044681	MVLYLGLPDKKRTG		hypothetical protein C0Q70_06994 [Pomacea canaliculata]	ALGADFQIEKLHASK	
PomcanEvM0124551	MKQIGRSQMLLVLV		delta-like protein A [Pomacea canaliculata]	MNQIGRSQMLLVLV	
PomcanEvM0053931	DCRNWQNCNHLI		uncharacterized protein LOC112568838 isoform X1 [Pomacea canaliculata]	FFVRRRRAAQTKTSE	
PomcanEvM0063511	MNCSLQATSSGHEI		uncharacterized protein LOC112567783 [Pomacea canaliculata]	FTCANPGKPLDIYE	
PomcanEvM0080131	MGDKVRAALVTLMI		uncharacterized protein LOC112575946 isoform X2 [Pomacea canaliculata]	MGDKVRAALVTLMF	
PomcanEvM0073691	KTTRLNLTYPPPKP		uncharacterized protein LOC112568756 isoform X2 [Pomacea canaliculata]	KTTRLNLTYPPPKPP	
PomcanEvM0114521	MAKRKGAGGGGGA		uncharacterized protein LOC112574796 [Pomacea canaliculata]	RVKVVYLHVTNDEA	
PomcanEvM0083581	MPKECAAPDCQEK		uncharacterized protein LOC112554666 isoform X1 [Pomacea canaliculata]	MPKECAAPDCQEK	
PomcanEvM0057951	MAHWSQWFNLDDP		uncharacterized protein LOC112572815 [Pomacea canaliculata]	MAHWSQWFNLDDP	
PomcanEvM0050071	MSPVRHPLGTSALV		uncharacterized protein LOC112569258 [Pomacea canaliculata] >ref XP_0251	MSPVCLPLGTSALV	
PomcanEvM0046161	MEANYFIVLFTGLF		LOW QUALITY PROTEIN: uncharacterized protein LOC112571572 [Pomacea ca	MEANYFIVLFTGLFL	
PomcanEvM0088281	MKPGLLMAGVVRD		hypothetical protein C0Q70_12470 [Pomacea canaliculata]	SLTFHIEGNKTSFTFP	
PomcanEvM0112921	MQFFKSDPPCLDL		leucine-rich repeats and immunoglobulin-like domains protein 2 isoform X1 [P	FRSDPPCLDLTYTT	
PomcanEvM0176851	LKPTRWSARCTIQS		hypothetical protein C0Q70_10661 [Pomacea canaliculata]	REGKAVHGVVWPKA	
PomcanEvM0191101	MEDDVVTAARRIM		uncharacterized protein LOC112558065 [Pomacea canaliculata]	MEDDVVTAARRIM	
PomcanEvM0132031	MIAMKAKEKISEILI		uncharacterized protein LOC112556833 isoform X1 [Pomacea canaliculata]	AKATVQCPLASAS	
PomcanEvM0071151	MWKGDDGQSNGTI		uncharacterized protein LOC112568994 isoform X3 [Pomacea canaliculata]	MWKGDDGQSNGTI	
PomcanEvM0067581	MRHVLVIAIVCCV		uncharacterized protein LOC112569162 isoform X2 [Pomacea canaliculata]	MRHVLVIAIVCCV	
PomcanEvM0242711	MVPSVFLMLMLAV		uncharacterized protein LOC112570348 [Pomacea canaliculata]	VVPTHGRSCPPG	
PomcanEvM0138851	MTLSTSSLLTLVMLF		uncharacterized protein LOC112560378 [Pomacea canaliculata] >ref XP_0251	VMLAIVTSPSLAAEQ	
PomcanEvM0045011	MTDNNEGNLAGRR		uncharacterized protein LOC112574294 [Pomacea canaliculata]	SRQHNSYLKDNLDI	
PomcanEvM0150671	MPSTSLCASATAV		tumor necrosis factor receptor superfamily member 13B-like [Pomacea canali	MPSTSLCASATAVV	

PomcanEVm01756711 MVTKYFIDSNNNK
PomcanEVm01280111 MSGETSLKFLIPDSS
PomcanEVm01947111 MVSFVFLVLMALVV
PomcanEVm00427611 MKISVFLLVFTVMV
PomcanEVm00858511 MKNLLVHLVSLPLFF
PomcanEVm02252911 MSSSCACCSKQTS
PomcanEVm02337211 MLVQAVASQLLEAV
PomcanEVm00570711 TYISQSVCLPAPWF
PomcanEVm00746211 MKMSCSTAHQGAS
PomcanEVm01732511 MEISFYRIFLFLITA
PomcanEVm00937811 MGCNFQCDNSLTM
PomcanEVm00585711 MAWCLLLVLVLLVS
PomcanEVm01396511 MTCYIESSRLHCTST
PomcanEVm00838811 METLRNMLLRPPYK
PomcanEVm00541411 MGKNRTSNACIGFI
PomcanEVm00545611 MSIHKSSGELCADLE
PomcanEVm02312411 MSLGLTSVAMVALL
PomcanEVm00575011 MRVREETLLLVVTV
PomcanEVm00438311 MAQRTWQVQFQTLT
PomcanEVm00912111 MGRVCVVGCGFL
PomcanEVm00531711 MTTTEEMKKFIAASE

uncharacterized protein LOC112568770 [Pomacea canaliculata] IVTKYFIDSNNNKE
uncharacterized protein LOC112567443 [Pomacea canaliculata] MSGETSLKFLIPDSS
uncharacterized protein LOC112570348 [Pomacea canaliculata] VVPTHTGRSPCPG
uncharacterized protein LOC112559049 isoform X3 [Pomacea canaliculata] >r MKI-
uncharacterized protein LOC112568269 [Pomacea canaliculata] >ref|XP_025: MKNLLVHLVSLPLFF
uncharacterized protein LOC112561109 isoform X2 [Pomacea canaliculata] MFSSCACCSKQTS
zinc finger C2HC domain-containing protein 1A-like isoform X1 [Pomacea cana RFGVKAGMSPSLLL
uncharacterized protein LOC112568447 [Pomacea canaliculata] EISATRRRATCSH-
uncharacterized protein LOC112568957 [Pomacea canaliculata] MKMSCSTAHQGASL
LOW QUALITY PROTEIN: pituitary tumor-transforming gene 1 protein-interactin MSNVTDVASTTEVLT
uncharacterized protein LOC112568712 [Pomacea canaliculata] >ref|XP_025: MGCNFQCDNSLKM
uncharacterized protein LOC112571063 isoform X3 [Pomacea canaliculata] MAWCLLLVLVLLST
uncharacterized protein LOC112567786 [Pomacea canaliculata] >ref|XP_025: MTCYIESSRLHCTST
uncharacterized protein LOC112575090 isoform X1 [Pomacea canaliculata] METLRNMLLRPPYK
uncharacterized protein LOC112568641 [Pomacea canaliculata] >gb|PVD273 DDQNTVSTLSMFTFR
hypothetical protein C0Q70_12045 [Pomacea canaliculata] ADLEFHGFRDGETV
hypothetical protein C0Q70_14184 [Pomacea canaliculata] MSLGLTSVVMVALLT
uncharacterized protein LOC112573282 isoform X1 [Pomacea canaliculata] >g MRVREETLLLVVTV
uncharacterized protein LOC112555430 [Pomacea canaliculata] >ref|XP_025f MAQRTRAATWQRFR
hypothetical protein C0Q70_20392 [Pomacea canaliculata] MPFVWSPKSIKIFEM
uncharacterized protein LOC112558245 [Pomacea canaliculata] MTTTEEMKKFIAASD

232

Query Sequence	Best Hit Description	Best Hit Sequence
PomcanEVm01287611 MTTVEVFLTAFAVM	uncharacterized protein LOC112567706 isoform X2 [Pomacea canaliculata]	MTTVEVFLTAFAVMT
PomcanEVm01850411 MSCSFNHLAHLMC	uncharacterized protein LOC112569223 isoform X1 [Pomacea canaliculata] >r	MSCSTEVEEKSQFLC
PomcanEVm00624711 MWSPRRLASPTTG	predicted protein, partial [Nematostella vectensis]	GGLVCSMNRLSLTR
PomcanEVm01793711 MSGETFVKFLIPDSS	uncharacterized protein LOC112567443 [Pomacea canaliculata]	MSGETSLKFLIPDSS
PomcanEVm01242911 MPLNLPIVELVLDRC	synaptotagmin-2-binding protein-like isoform X1 [Pomacea canaliculata]	MPLNLPIVELVLDRC
PomcanEVm05356211 MVLMPVFMHKGIE	40S ribosomal protein S3 [Aplysia californica]	LMIHSGEPLNDYVD
PomcanEVm0285211 MKADPPVSWEGEV	uncharacterized protein LOC112568675 isoform X2 [Pomacea canaliculata]	EGNPPPVSWEGGG
PomcanEVm01288411 MKTDLVILSLSLIS	uncharacterized protein LOC112567701 [Pomacea canaliculata]	MKTDLVILSLSLIS
PomcanEVm00952711 YKPLRSRYSGACSF	hypothetical protein C0Q70_03286 [Pomacea canaliculata]	VAYPPVGLYGVWQ
PomcanEVm00790911 MTAKEYIPLFVAKIL	mucin-5AC-like isoform X2 [Pomacea canaliculata]	MTAKYIPLFVAKIL
PomcanEVm04149211 MIKPLRVFAGHTHQ	hypothetical protein I4U23_009426 [Adineta vaga]	MIKPLRVFAGHTHQ
PomcanEVm01476411 MSTNDHNDPDISST	unnamed protein product [Rotaria sp. Silwood2] >emb CAF2528129.1 unnar	MSSNKQNNQEISTS
PomcanEVm02845111 MTLIRSNLFIETLHV	emp24/gp25L/p24 family protein [Dictyostelium discoideum]	LTDPKNTIFERLHF
PomcanEVm01581411 MSGEMILKQIPDSS	uncharacterized protein LOC112567439 isoform X2 [Pomacea canaliculata]	MSGETSLKFLIPDSS
PomcanEVm00538911 MAVRLPFLVTVTVL	hypothetical protein C0Q70_04880 [Pomacea canaliculata]	MALRQPFVTVTVL
PomcanEVm02883311 MAVGECSTLPSVLD	uncharacterized protein LOC112557320 [Pomacea canaliculata]	MAVGECSTLPSVLD
PomcanEVm01695311 MTCEIKSHHPVQIK	uncharacterized protein LOC112568155 [Pomacea canaliculata]	LDRRLRRQTEGFI
PomcanEVm00486611 MAFNFRHOSAKAM	hypothetical protein C0Q70_04882 [Pomacea canaliculata]	MYGRNQGARATLTS
PomcanEVm00846311 MIRKHTFDCVLLCV	hypothetical protein C0Q70_07607 [Pomacea canaliculata]	TTRVQRKCEDPDVV
PomcanEVm02607311 MTRSSKRFRFRYEYK	uncharacterized protein LOC112564392 [Pomacea canaliculata] >ref XP_025f ERREFRIEYNKPGQK	ERREFRIEYNKPGQK
PomcanEVm00996011 MLVVRPVLITLIFVY	unnamed protein product [Rotaria sp. Silwood1] >emb CAF4684318.1 unnar	NPPESTSLERAHWC
PomcanEVm00446911 MAAPAYVHSENAQ	protein draper-like [Pomacea canaliculata]	ECVCVEGWGTGAYC
PomcanEVm01298111 MREQMVLDEKIGLE	unnamed protein product [Adineta ricciae]	MREQMVLDEKIGV
PomcanEVm01300611 MSCSPGYSAEGACF	platelet endothelial aggregation receptor 1-like [Pomacea canaliculata]	MCSQPGYSAEGACF
PomcanEVm02126411 MSGETSVKFLIPDSS	uncharacterized protein LOC112567443 [Pomacea canaliculata]	MSGETSLKFLIPDSS
PomcanEVm01606211 MNFVGKILNEANPL	hypothetical protein M9434_003998 [Picochlorum sp. BPE23]	VNFFGVVLKPDAPTA
PomcanEVm01371111 MSVRTQSREPLLRNI	uncharacterized protein LOC112568424 [Pomacea canaliculata]	VISFLQLALVLSFASG
PomcanEVm00922111 MHIMCASSTLTSLS	uncharacterized protein LOC112569050 [Pomacea canaliculata]	MSGFFCLVLVTMTY
PomcanEVm03551011 MQRQLTLTRISCVS	PLAC8 family-domain-containing protein [Jimmerdemannia flammicorona]	LGDGSCFSQGAICY
PomcanEVm02483311 MDTDQPVQRHTAV	unnamed protein product [Rotaria sp. Silwood2] >emb CAF2538398.1 unnar	FSEDFLNRTFHHNP
PomcanEVm01560711 MSGILFYVNNRNESI	uncharacterized protein LOC112567463 [Pomacea canaliculata]	GQKHDANVCKLIVT
PomcanEVm01214011 MIIYFITIHCLKAHC	uncharacterized protein LOC112567789 [Pomacea canaliculata]	VQAGQGPDLYPSPF
PomcanEVm01672211 MRCITTTSTPFLHDC	uncharacterized protein LOC112569064 [Pomacea canaliculata] >ref XP_025: MRCASTRITAPFHD	MRCASTRITAPFHD
PomcanEVm01278511 MNCISLEFDGLRSS	hypothetical protein C0Q70_12062 [Pomacea canaliculata]	VEFSGLGSGENWTEI
PomcanEVm01120611 MSLYILLVWTLGAV	hypothetical protein C0Q70_13127 [Pomacea canaliculata]	MSLYILLVWTLGAV
PomcanEVm01159811 VENCAELVFGNDTD	uncharacterized protein LOC112568328 isoform X1 [Pomacea canaliculata]	VENCAELVFGNDTD
PomcanEVm00712811 MVTCVACGVCLMLI	uncharacterized protein LOC112568313 [Pomacea canaliculata] >ref XP_025: MVTCVACGVCLMLI	MVTCVACGVCLMLI
PomcanEVm01118711 MEQTGLLHGHYFVI	uncharacterized protein LOC112569216 [Pomacea canaliculata]	EMTKSSSQESSQMT
PomcanEVm06009411 MQVSDLVKVRITPRF	histone H2A.V-like [Lagopus leucura]	DLKVKRITPRHLQAI
PomcanEVm01139711 MLSSLAVYRALKKEW	uncharacterized protein LOC112568183 isoform X1 [Pomacea canaliculata]	APRLPGNTDHPSLIK
PomcanEVm01398111 MMFSIWLRLHTAISR	uncharacterized protein LOC112569144 [Pomacea canaliculata] >ref XP_025: EEVSLTCTGFENINE	EEVSLTCTGFENINE
PomcanEVm01769211 MQRMELFYRGDISR	uncharacterized protein LOC112569144 [Pomacea canaliculata] >ref XP_025: MHFSSDAHLK	MHFSSDAHLK
PomcanEVm01242211 MFTSFCLFTFVLFTL	hypothetical protein C0Q70_12137 [Pomacea canaliculata]	TEDDSPRPVGVVAV
PomcanEVm01362211 MLRNHDDSVMQQC	uncharacterized protein LOC112569996 [Pomacea canaliculata]	SRASPSTSYRKK----
PomcanEVm00760611 VVFERVKCAQKEKC	uncharacterized protein LOC112569625 isoform X2 [Pomacea canaliculata]	LVRTIVAYTHLENCA
PomcanEVm01147611 MMSTQGIQSELYLTI	uncharacterized protein LOC112567729 isoform X1 [Pomacea canaliculata]	MMSAQGIQSELYLTI
PomcanEVm01665811 KHPLKQNVSTPSSAA	furin-like protease kpc-1 isoform X5 [Pomacea canaliculata]	KHPLKQNVSTPSSAA

PomcanEVm007776t1 MFLSAAKRWVQV
PomcanEVm006007t1 MTSRLFWCTFVAV
PomcanEVm014346t1 MTCEIKKCPFEIQ
PomcanEVm014235t1 MRCSTTRTPFHLDX
PomcanEVm014348t1 MSFSTATTQFLERR
PomcanEVm008955t1 MSTTRPSPEVSLTN
PomcanEVm005626t1 SGKEKYLDRCIGIC
PomcanEVm018370t1 MSGETSLLELPDSS
PomcanEVm011503t1 MMKIALLVLTATVT
PomcanEVm012464t1 MFIFMYLALLAFGC
PomcanEVm014238t1 MATYGICATDKTCA
PomcanEVm007921t1 MDTPSTPQGTSSG
PomcanEVm004292t1 MCLHRLSHNKLKTL
PomcanEVm005714t1 MAAILPFSIVAMTL
PomcanEVm014158t1 MSCSTRHQATHQLC
PomcanEVm016944t1 MLFCTTHLSTGVID
PomcanEVm001416t1 MELEKLYFVDLFLF
PomcanEVm010717t1 MTGHFYSCDDSRIT
PomcanEVm010990t1 MSCQLKNTGCPFEI
PomcanEVm008564t1 MNPEELFKRFFGFT
PomcanEVm010835t1 MEHVHRAVLRFVMI
PomcanEVm009842t1 MAATTAYLLVFTVM
PomcanEVm006395t1 MHERVLDRLCTNVI
PomcanEVm012206t1 MHAAGVCRHSLSLR
PomcanEVm018311t1 MSCSFNLHLALMC
PomcanEVm011721t1 MRNFSNYQLLLIT
PomcanEVm009652t1 MTLVLSRLAFVVIIL
PomcanEVm008722t1 MAYTSRASFLNRF

hypothetical protein COQ70_06059 [Pomacea canaliculata]
uncharacterized protein LOC112574344 isoform X1 [Pomacea canaliculata]
uncharacterized protein LOC112569680 [Pomacea canaliculata]
uncharacterized protein LOC112569064 [Pomacea canaliculata] >ref|XP_025_MRCASRTAPFHD
uncharacterized protein LOC112569075 [Pomacea canaliculata]
LOW QUALITY PROTEIN: uncharacterized protein LOC112567035 [Pomacea canaliculata]
uncharacterized protein LOC112568766 isoform X1 [Pomacea canaliculata]
uncharacterized protein LOC112567439 isoform X2 [Pomacea canaliculata]
hypothetical protein COQ70_06060 [Pomacea canaliculata]
uncharacterized protein LOC112568201 [Pomacea canaliculata]
uncharacterized protein LOC112569050 [Pomacea canaliculata]
uncharacterized protein LOC112569182 [Pomacea canaliculata] >ref|XP_025_LGVLFLWLPLANC
uncharacterized protein LOC112567763 [Pomacea canaliculata]
hypothetical protein COQ70_04879 [Pomacea canaliculata]
uncharacterized protein LOC112569223 isoform X1 [Pomacea canaliculata] >ref|XP_025_MSCITEVEEKSQIFC
uncharacterized protein LOC112569520 [Pomacea canaliculata]
uncharacterized protein LOC112568578 [Pomacea canaliculata]
unnamed protein product [Adineta steineri]
hypothetical protein COQ70_12375 [Pomacea canaliculata]
uncharacterized protein LOC112553847 isoform X1 [Pomacea canaliculata]
uncharacterized protein LOC112572452 isoform X2 [Pomacea canaliculata] >ref|XP_025_MEHVHRAVLRFVMI
uncharacterized protein LOC112568212 [Pomacea canaliculata]
uncharacterized protein LOC112558823 [Pomacea canaliculata]
uncharacterized protein LOC112569555 isoform X1 [Pomacea canaliculata]
uncharacterized protein LOC112569223 isoform X1 [Pomacea canaliculata] >ref|XP_025_MSCITEVEEKSQIFC
uncharacterized protein LOC112568273 isoform X3 [Pomacea canaliculata]
uncharacterized protein LOC112569165 [Pomacea canaliculata]
uncharacterized protein LOC112572874 [Pomacea canaliculata]
uncharacterized protein LOC112572874 [Pomacea canaliculata]

207	Query Sequence	Best Hit Description	Best Hit Sequence
PomcanEVm007646t1	MPVKEFNDNINVG	uncharacterized protein LOC112555791 [Pomacea canaliculata] >gb PVD190	MPVKEFNDNINVG
PomcanEVm045415t1	ISAKIALSTCFICSTT	hypothetical protein COQ70_06414 [Pomacea canaliculata]	FILENDLSFYVGVGD
PomcanEVm034938t1	MLLLISLNLVAKSS	hypothetical protein [Philodina roseola]	ISHMKKQPIRAVLN
PomcanEVm036462t1	RSFLRCNVASSLVIA	hypothetical protein I4U23_020465 [Adineta vaga]	MGSDDSKSKKRATQ
PomcanEVm053774t1	SFGTVVKPKLKAIV	non-canonical purine NTP pyrophosphatase, partial [Escherichia coli]	ILRFARGAHGFGYDP
PomcanEVm009016t1	MSAQPIVSNVSDSH	uncharacterized protein LOC112560396 [Pomacea canaliculata] >gb PVD350	MSAQPIVSNVSDSH
PomcanEVm006685t1	MEHPMSQQVLFEDA	uncharacterized protein LOC112555076 [Pomacea canaliculata] >gb PVD187	MEHPMSQQVFEDA
PomcanEVm015891t1	MPLSRPVFYRLILFL	unnamed protein product [Rotaria magnacalcarata]	KRKLHFACFI
PomcanEVm012091t1	MTTRNEKAPRYFV	hypothetical protein SeMB42_g04412 [Synchytrium endobioticum] >gb TPX46	THAINCSDKSLVR
PomcanEVm008466t1	MTGTVHFDPVPTAI	tumor necrosis factor ligand superfamily member 11-like [Pomacea canaliculata] >ref XP_025_MPPSRIPP	GADLLTC
PomcanEVm077461t1	MEKFNIEKDIAIYKI	hypothetical protein INT44_005811 [Umbelopsis vinacea]	LEKYNIEKDIAIHKR
PomcanEVm010351t1	MTIYKTPGGKTEK	integral membrane protein 2C-like [Pomacea canaliculata]	LCDRLTGVDKEVC
PomcanEVm049432t1	MQEKDDTLQRVRL	myosin heavy chain, embryonic smooth muscle isoform-like [Pomacea canaliculata] >ref XP_025_MVGSVLS	INTPNKADD
PomcanEVm011325t1	MVGSVLSNNAKSD	unnamed protein product [Rotaria sp. Silwood1] >emb CAF1616688.1 unna	MVGSVLSNNKADD
PomcanEVm061939t1	DEIRFGKQFPAVI	hypothetical protein SteCoe_35774 [Stentor coeruleus]	DEIRFGPNRIKNEP
PomcanEVm02269t1	MAMLSNHELIRRV	cyclic nucleotide-binding domain-containing protein [Burkholderiales bacteriu	MAMLSNHELIRRVPI
PomcanEVm008678t1	MQRTSGHVSWPT	uncharacterized protein LOC112555273 [Pomacea canaliculata] >ref XP_025_MQTSTSGHVS	WPTN
PomcanEVm014712t1	MLCIIQIFTEPNILS	alpha-N-acetylgalactosamine-specific lectin-like isoform X1 [Pomacea canaliculata] >ref XP_025_MATRSRPR	FTRTRVA
PomcanEVm010699t1	MECVKRSIHVALVFI	tumor necrosis factor ligand superfamily member 10-like [Pomacea canaliculata] >ref XP_025_QERGF	TKTYDNHITL
PomcanEVm016587t1	MLGVVDQAKGLYRI	microtubule-actin cross-linking factor 1-like isoform X6 [Pomacea canaliculata] >ref XP_025_MLGVVDQ	AKGLYRIN
PomcanEVm015528t1	WDDRPNSVLLASGI	RING finger protein 150-like [Pomacea canaliculata] >gb PVD35007.1 hypoth	MNPCTSIFFSVIVLL
PomcanEVm009195t1	MGFSTLFGSNTITTA	tumor necrosis factor ligand superfamily member 11-like [Pomacea canaliculata] >ref XP_025_MGFSTL	FGSKTITTD
PomcanEVm005429t1	MAATERQWSIREDC	LOW QUALITY PROTEIN: baculoviral IAP repeat-containing protein 6-like [Pom	MAATERQWSIREDC
PomcanEVm009897t1	MSESRSVSLASGSG	hypothetical protein COQ70_11300 [Pomacea canaliculata]	MSDSRVSSLASGSG
PomcanEVm020932t1	MLETTEYTLQVPSFI	hypothetical protein Clob_013751 [Chrysochromulina tobinii]	NWMGVKLEEDSFG
PomcanEVm006118t1	MSVSLQDLGVREK	uncharacterized protein LOC112555073 isoform X7 [Pomacea canaliculata]	MSVSLQDLGVREK
PomcanEVm009258t1	MKTRLSSESMDST	tumor necrosis factor-like [Pomacea canaliculata]	MKTRLSSESMDST
PomcanEVm012947t1	MPSRLHTGILRGRGI	hypothetical protein COQ70_14854 [Pomacea canaliculata]	VGTFVRFSADPTNV

NP_000585.2

308	Query Sequence	Best Hit Description	Best Hit Sequence
PomcanEVm001120t1	MMDFNSKVFFFLII	unnamed protein product [Rotaria sp. Silwood2] >emb CAF4070867.1 unna	MNFIVKTSFEILLVIL
PomcanEVm009674t1	MAMLSMDSTVLLV	unnamed protein product [Rotaria sp. Silwood1] >emb CAF1605507.1 unna	EDGFEAVGCCCLCCA
PomcanEVm007060t1	MACARNWILLVV	uncharacterized protein LOC112568895 isoform X3 [Pomacea canaliculata]	DVVTSCPEYFIAGQV
PomcanEVm018791t1	MQDRRLNLVTDRIE	uncharacterized protein LOC112571157 [Pomacea canaliculata]	MQDRRLNLATDRIE

CD55

374	Query Sequence	Best Hit Description	Best Hit Sequence
PomcanEVm002703t1	MAVWILWLCGVAL	uncharacterized protein LOC112568828 isoform X1 [Pomacea canaliculata]	LWRLWLCGVALLAI
PomcanEVm000688t1	MYRLIASCLLYLLIQ	uncharacterized protein LOC112569288 isoform X3 [Pomacea canaliculata]	MYRLIASCLLYLLIP
PomcanEVm010228t1	MCVLVLLTAAQEAC	uncharacterized protein LOC112568119 [Pomacea canaliculata]	MCVLVLLTAAQEAC
PomcanEVm037242t1	MTLRKGGGGRPM	cubilin-like isoform X2 [Pomacea canaliculata]	MTLRKGGGGRPM
PomcanEVm001480t1	MFLIACHRTLAIIFI	uncharacterized protein LOC112568581 isoform X1 [Pomacea canaliculata]	MFLIACHRTLAIISLQ
PomcanEVm008328t1	MKFCSEYFRICQLQ	uncharacterized protein LOC112560063 [Pomacea canaliculata]	MGLEYKPDQKQVYC
PomcanEVm001993t1	MRLLSCVIIVAAAL	tollid-like protein 2 [Pomacea canaliculata]	MESCINCSDFLAL

CPAH_HUMAN

Bibliografia

A.Accorsi, E. Ross, E. Ottaviani, A. S. Alvarado, Eur. J. Histochem. 2017, 61, 1.

Accorsi, A. Soluble Factors in the Immune-Neuroendocrine System of Invertebrate Models. Ph.D. Thesis, University of Modena and Reggio Emilia, Modena, Italy, 2015.

Accorsi, Alice, Enzo Ottaviani, and Davide Malagoli. "Effects of repeated hemolymph withdrawals on the hemocyte populations and hematopoiesis in *Pomacea canaliculata*." *Fish & Shellfish Immunology* 38.1 (2014): 56-64.

Accorsi, Alice, et al. "Comparative analysis of circulating hemocytes of the freshwater snail *Pomacea canaliculata*." *Fish & shellfish immunology* 34.5 (2013): 1260-1268.

Accorsi, Alice, et al. "Complete Regeneration of a Camera-type Eye in the Research Organism *Pomacea canaliculata*." *The FASEB Journal* 32 (2018): 232-4.

Akira, Shizuo, Satoshi Uematsu, and Osamu Takeuchi. "Pathogen recognition and innate immunity." *Cell* 124.4 (2006): 783-801.

Alberts, Bruce, et al. *Biologia molecolare della cellula*. Bologna: Zanichelli, 2009.

Al-Sharif, W.Z., Sunyer, J.O., Lambris, J.D., Smith, L.C., 1998. Sea urchin coelomocytes specifically express a homologue of the complement component C3. *J. Immunol.* 160, 2983–299.

Amara, U. et al. Molecular intercommunication between the complement and coagulation systems. *J. Immunol.* 185, 5628–5636 (2010).

B. Chen *et al.*, "xTrimoPGLM: Unified 100B-Scale Pre-trained Transformer for Deciphering the Language of Protein." *bioRxiv*, p. 2023.07.05.547496, Jul. 06, 2023.

Bassi, N., et al. "Pentraxins, anti-pentraxin antibodies, and atherosclerosis." *Clinical Reviews in Allergy & Immunology* 37.1 (2009): 36-43.

Baum, Matthew L., et al. "CSMD1 regulates brain complement activity and circuit development." *Brain, Behavior, and Immunity* 119 (2024): 317-332.

Bergamini, Giulia, et al. "A New Protocol of Computer-Assisted Image Analysis Highlights the Presence of Hemocytes in the Regenerating Cephalic Tentacles of Adult *Pomacea canaliculata*." *International journal of molecular sciences* 22.9 (2021): 5023.

Beutler, Bruce. "Innate immunity: an overview." *Molecular immunology* 40.12 (2004): 845-859.

Bever, Michele Miller, and Richard B. Borgens. "Eye regeneration in the mystery snail." *Journal of Experimental Zoology* 245.1 (1988): 33-42.

Brehélin, Michel, and Patricia Roch. "Specificity, learning and memory in the innate immune response." *Invertebrate Survival Journal* 5.2 (2008): 103-109.

- Clow, L.A., Gross, P.S., Shih, C.-S., Smith, L.C., 2000. Expression of SpC3, the sea urchin complement component, in response to lipopolysaccharide. *Immunogenetics* 51, 1021–1033.
- Cueto, J. A., Rodriguez, C., Vega, I. A., & Castro-Vazquez, A. (2015). Immune defenses of the invasive apple snail *Pomacea canaliculata* (Caenogastropoda, Ampullariidae): phagocytic hemocytes in the circulation and the kidney. *PLoS One*, 10(4), e0123964.
- Denny, Kerina J., et al. "C5a receptor signaling prevents folate deficiency–induced neural tube defects in mice." *The Journal of Immunology* 190.7 (2013): 3493-3499.
- Escudero-Esparza, A., Kalchishkova, N., Kurbasic, E., Jiang, W.G., Blom, A.M., 2013. The novel complement inhibitor human CUB and Sushi multiple domains 1 (CSMD1) protein promotes factor I-mediated degradation of C4b and C3b and inhibits the membrane attack complex assembly. *FASEB J.* 27, 5083–5093.
- Faulk, W. P., et al. "Immunological studies of human placentae: complement components in immature and mature chorionic villi." *Clinical and experimental immunology* 40.2 (1980): 299. Akira, Shizuo, Satoshi Uematsu, and Osamu Takeuchi. "Pathogen recognition and innate immunity." *Cell* 124.4 (2006): 783-801.
- Ghosh, Julie, et al. "Invertebrate immune diversity." *Developmental & Comparative Immunology* 35.9 (2011): 959-974.
- Harris, Claire L., Masashi Mizuno, and B. Paul Morgan. "Complement and complement regulators in the male reproductive system." *Molecular immunology* 43.1-2 (2006): 57-67. Accorsi, A. Soluble Factors in the Immune-Neuroendocrine System of Invertebrate Models. Ph.D. Thesis, University of Modena and Reggio Emilia, Modena, Italy, 2015.
- Håvik, B., Le Hellard, S., Rietschel, M., Lybæk, H., Djurovic, S., Mattheisen, M., Mühleisen, T.W., Degenhardt, F., Priebe, L., Maier, W., Breuer, R., Schulze, T.G., Agartz, I., Melle, I., Hansen, T., Bramham, C.R., Nothen, M.M., Stevens, B., Werge, T., Andreassen, O.A., Cichon, S., Steen, V.M., 2011. The complement control related genes CSMD1 and CSMD2 associate to schizophrenia. *Biol. Psychiatry* 70, 35–42.
- Hawthornthwaite, Owen A., et al. "Complement in stem cells and development." *Seminars in Immunology*. Vol. 37. Academic Press, 2018.
- Heinzinger, Michael, et al. "Prost5: Bilingual language model for protein sequence and structure. bioRxiv." (2023): 2023-07.
- Hirschfield, Gideon M., et al. "Human C-reactive protein does not protect against acute lipopolysaccharide challenge in mice." *The Journal of Immunology* 171.11 (2003): 6046-6051.
- Irmscher, S. et al. Kallikrein cleaves C3 and activates complement. *J. Innate Immun.* 10, 94–105 (2018)

Kimbrell, Deborah A., and Bruce Beutler. "The evolution and genetics of innate immunity." *Nature Reviews Genetics* 2.4 (2001): 256-267.

King, B. C. & Blom, A. M. Intracellular complement: evidence, definitions, controversies, and solutions. *Immunol. Rev.* 313, 104–119 (2023). Comprehensive overview of the mechanisms and functions of intracellular complement, which also highlights experimental and conceptual challenges.

Klein, Jan, and Vaclav Horejsi. *Immunology*. Wiley-Blackwell, 1999.

Köhl, Jörg. "The role of complement in danger sensing and transmission." *Immunologic research* 34.2 (2006): 157-176.

Kolev, Martin, Gaelle Le Friec, and Claudia Kemper. "Complement—tapping into new sites and effector systems." *Nature Reviews Immunology* 14.12 (2014): 811-820.

Kraus, D.M., Elliott, G.S., Chute, H., Horan, T., Pfenninger, K.H., Sanford, S.D., Foster, S., Scully, S., Welcher, A.A., Holers, V.M., 2006. CSMD1 is a novel multiple domain complement-regulatory protein highly expressed in the central nervous system and epithelial tissues. *J. Immunol.* 176, 4419–4430.

Krych-Goldberg, M., Atkinson, J., 2001. Structure-function relationships of complement receptor type 1. *Immunol. Rev.* 180, 112–122.

Lachmann, P.J., 1979. An evolutionary view of the complement system. *Behring Institute Mitteilungen* 63, 25–37.

M. Heinzinger *et al.*, "Modeling aspects of the language of life through transfer-learning protein sequences," *BMC Bioinformatics*, vol. 20, no. 1, p. 723, Dec. 2019.

Mannes, Marco, et al. "Tuning the functionality by splicing: factor H and its alternative splice variant FHL-1 share a gene but not all functions." *Frontiers in Immunology* 11 (2020): 596415.

Mastellos, D. C., Deangelis, R. A. & Lambris, J. D. Complement-triggered pathways orchestrate regenerative responses throughout phylogenesis. *Semin. Immunol.* 25, 29–38 (2013).

Medzhitov, R., Janeway Jr, C.A., 2002. Decoding the patterns of self and nonself by the innate immune system. *Science* 296, 298–300.

Merle, N.S., Church, S.E., Fremeaux-Bacchi, V., Roumenina, L.T., 2015a. Complement system part I - molecular mechanisms of activation and regulation. *Front. Immunol.* 6, 262.

Merle, Nicolas S., et al. "Complement system part I—molecular mechanisms of activation and regulation." *Frontiers in immunology* 6 (2015): 262.

Nilsson, P. H. et al. A conformational change of complement C5 is required for thrombin-mediated cleavage, revealed by a novel ex vivo human whole blood model preserving full thrombin activity. *J. Immunol.* 207, 1641–1651 (2021)

- Nonaka, Masaru, et al. "Opsonic complement component C3 in the solitary ascidian, *Halocynthia roretzi*." *The Journal of Immunology* 162.1 (1999): 387-391.
- Ochsenbein, Adrian F., et al. "Control of early viral and bacterial distribution and disease by natural antibodies." *Science* 286.5447 (1999): 2156-2159.
- Peng, M., Niu, D., Wang, F., Chen, Z., Li, J., 2016. Complement C3 gene: expression characterization and innate immune response in razor clam *Sinonovacula constricta*. *Fish Shellfish Immunol.* 55, 223–232.
- Peng, Maoxiao, et al. "Complement C3 gene: expression characterization and innate immune response in razor clam *Sinonovacula constricta*." *Fish & shellfish immunology* 55 (2016): 223-232.
- Post, Theodore W., et al. "Membrane cofactor protein of the complement system: alternative splicing of serine/threonine/proline-rich exons and cytoplasmic tails produces multiple isoforms that correlate with protein phenotype." *The Journal of experimental medicine* 174.1 (1991): 93-102.
- Pouw, R.B., Vredevoogd, D.W., Kuijpers, T.W., Wouters, D., 2015. Of mice and men: the factor H protein family and complement regulation. *Mol. Immunol.* 67, 12–20.
- Ricklin, D., Hajishengallis, G., Yang, K. & Lambris, J. D. Complement: a key system for immune surveillance and homeostasis. *Nat. Immunol.* 11, 785–797 (2010).
- Ricklin, Daniel, et al. "Complement: a key system for immune surveillance and homeostasis." *Nature immunology* 11.9 (2010): 785-797.
- Rooney, Isabelle A., John E. Heuser, and John P. Atkinson. "GPI-anchored complement regulatory proteins in seminal plasma. An analysis of their physical condition and the mechanisms of their binding to exogenous cells." *The Journal of clinical investigation* 97.7 (1996): 1675-1686.
- Russell, S. (2021). *Artificial Intelligence: A Modern Approach*. Pearson.
- Russell, Sarah M., et al. "Tissue-specific and allelic expression of the complement regulator CD46 is controlled by alternative splicing." *European journal of immunology* 22.6 (1992): 1513-1518.
- Schmidt, Bela Z., and Harvey R. Colten. "Complement: a critical test of its biological importance." *Immunological reviews* 178 (2000): 166-176.
- Shishido, Stephanie N., et al. "Humoral innate immune response and disease." *Clinical immunology* 144.2 (2012): 142-158.
- Smith, L. Courtney, et al. "The echinoid complement system inferred from genome sequence searches." *Developmental & Comparative Immunology* 140 (2023): 104584.
- Smith, L.C., Azumi, K., Nonaka, M., 1999. Complement systems in invertebrates. The ancient alternative and lectin pathways. *Immunopharmacology* 42, 107–120.

Smith, L.C., Chang, L., Britten, R.J., Davidson, E.H., 1996. Sea urchin genes expressed in activated coelomocytes are identified by expressed sequence tags. Complement homologues and other putative immune response genes suggest immune system homology within the deuterostomes. *J. Immunol.* 156, 593–602.

Smith, L.C., Clow, L.A., Terwilliger, D.P., 2001. The ancestral complement system in sea urchins. *Immunol. Rev.* 180, 16–34.

Smith, Valerie J., Alice Accorsi, and Davide Malagoli. "Hematopoiesis and hemocytes in pancrustacean and molluscan models." *The evolution of the immune system*. Academic Press, 2016. 1-28.

Song, W. C., Sarrias, M. R. & Lambris, J. D. Complement and innate immunity. *Immunopharmacology* 49, 187–198 (2000).

Stephan, A. H., Barres, B. A. & Stevens, B. The complement system: an unexpected role in synaptic pruning during development and disease. *Annu. Rev. Neurosci.* 35, 369–389 (2012).

Stevens, Beth, et al. "The classical complement cascade mediates CNS synapse elimination." *Cell* 131.6 (2007): 1164-1178. Accorsi, Alice, et al. "Complete Regeneration of a Camera-type Eye in the Research Organism *Pomacea canaliculata*." *The FASEB Journal* 32 (2018): 232-4.

Stevens, Beth, et al. "The classical complement cascade mediates CNS synapse elimination." *Cell* 131.6 (2007): 1164-1178.

Taylor, Clare T., and Peter M. Johnson. "Complement-binding proteins are strongly expressed by human preimplantation blastocysts and cumulus cells as well as gametes." *Mol Hum Reprod* 2.1 (1996): 52-59. Accorsi, Alice, Enzo Ottaviani, and Davide Malagoli. "Effects of repeated hemolymph withdrawals on the hemocyte populations and hematopoiesis in *Pomacea canaliculata*." *Fish & Shellfish Immunology* 38.1 (2014): 56-64.

Terwilliger, D.P., Clow, L.A., Gross, P.S., Smith, L.C., 2004. Constitutive expression and alternative splicing of the exons encoding SCRs in Sp152, the sea urchin homologue of complement factor B. Implications on the evolution of the Bf/C2 gene family. *Immunogenetics* 56, 531–543.

Trouw, Leendert A., Matthew C. Pickering, and Anna M. Blom. "The complement system as a potential therapeutic target in rheumatic disease." *Nature Reviews Rheumatology* 13.9 (2017): 538-547.

V. J. Smith, A. Accorsi, D. Malagoli, in *The Evolution of the Immune System: Conservation and Diversification*, 1st ed. (Ed: D. Malagoli), Academic Press/Elsevier, Cambridge, Massachusetts, USA 2016, Ch. 1.

View 100 of the world's worst invasive alien species, December 2000. Updated and reprinted version: November 2004. https://www.issg.org/worst100_species.html.

Wallis, Russell. "Interactions between mannose-binding lectin and MASPs during complement activation by the lectin pathway." *Immunobiology* 212.4-5 (2007): 289-299.

Walport, Mark J. "Complement first of two parts." *N Engl J Med.* 344 (2001): 1058-1066.

Yang, Ting-Bao, Zhong-Dao Wu, and Zhao-Rong Lun. "The apple snail *Pomacea canaliculata*, a novel vector of the rat lungworm, *Angiostrongylus cantonensis*: its introduction, spread, and control in China." *Hawai'i Journal of Medicine & Public Health* 72.6 Suppl 2 (2013): 23.