

Key to Gene Annotations

- **General organization:** The medial acid-base dyads (ABDs) are “parsed” into “strings that initiate with a dyad. N-term domain is denoted the head, C-term the tail.
- **Highlights:** Yellow, tyrosine; Green, cysteine; Pink, VPV; Gray, ABD in predicted head and tail domains.
- **Colored font:** Red, glycine residues in ABD domains; Bold-faced red, last G residue in head and first G residue in tail; Blue, alanine residues in ABD domains; Green, repeated string domains.
- **Underscores:** Acid-base triads or tetrads.
- **Commentary:** Notes on distinctive gene features (e.g. orthologues, localization patterns) in green font at top of some pages. Predicted homology domains (e.g. PDZ, coiled-coil) are given in Supplement Table 2.
- **Secondary structure:** PSIPRED predictions for a subset of proteins. Yellow, amino acids predicted in β -strand; pink, amino acids predicted in α -helix; no highlight, amino acids predicted in disordered (random coil) domain.

Articulin

Predicted secondary structure next slide

MSYLAGAPVTAAASYPVVPGAAAYPVGASYPAVAQASVLRANAQFAGQQRVIR

KEVVG^RNVEY^YIPVQHNHEVELIEREFIV^{PV}EKIVQRHV^{PV}PVERIVQRRV^{PV}PRQ^VPV^PQ^RVEIP^VPVERIQHRQVP^YPVEQIVEKRI^PV^PTQ^TVEQ^AVE^VPV^PPV^HRRVIQ^QV^PP^HAV^VREVIRHEP^YPVTKEVTRQ^VP^VEVPREVVRQVTVD^VP^VQVPQH^VQVP^YPVEKVVHRQVP^YPVEKVVQRQVP^YPVQKIVERQVQVP^YEVLVP

ERVEIPVPHEVITH

RD^VP^VPQE^VIRT^VQ^VP^VP^VEQIV^HRDVP^YPVEQIVEKVVQVTRQVTVPEIVQ^VP^VP^HEVIVERR^VP^VP^VERITHK^AVP^YPVEQIVEKIVQ^VP^VP^QY^QK^VP^VQ^VP^VP^V

ERIVT

RDVP^YPVEQIV

DKVV

ERQ^VP^VPT^VQ^VP^VPT^VQ^VP^YPV

EKIV

DRPVPHEVVRVV

ERRVEVP^YD^VP^VP^VVIETVQVPHEVIRTVEVPFPVEQIV

EKIV

EKIVH^VP^VP^VTRQEHVTRQVVQNTLQQTRAPPQQL^GTQVLPGRDLGVTSGIVRGGPAAYGAYGA

TPVAAAPVTYGSVAAPVTYGAAPVTYGAAPVTAYGAPPVYAGAPAVVTAPSSVT

RYGYGPPRPAYPAT^VP^VAYAAAPT^ELRNDDFDWPPVPSV^CGVQRQPQ

Euglena 80K articulin AAB23240.1

Predicted secondary structure

```
1 M S Y L A G A P V T A A A S Y P V V P G A A A Y P V G A S Y P A V A Q A S V L R A N A Q F A G Q R V
51 I R K E V V G N V E Y I P V Q H N H E V E L I E R E F I V P V E K I V Q R H V P V P V E R I V Q R R
101 V P V P R Q V P V P Q R V E I P V P V E R I Q H R Q V P Y P V E Q I V E K R I P V P V T Q T V E Q A
151 V E V P V P V H R R V I Q Q V P V P H A V V R E V I R H E P Y P V T K E V T R Q V P V E V P R E V V
201 R Q V T V D V P V Q V P Q H V Q V P Y P V E K V V H R Q V P Y P V E K V V Q R Q V P Y P V Q K I V E
251 R Q V Q V P Y E V L V P E R V E I P V P H E V I T H R D V P V P Q E V I R T V Q V P V P V E Q I V H
301 R D V P Y P V E Q I V E K V V Q V T R Q V T V P E I V Q V P V P H E V I V E R R V P V P V E R I T H
351 K A V P Y P V E Q I V E K I V Q V P V P Q Y Q K V P V Q V P V P V E R I V T R D V P Y P V E Q I V D
401 K V V E R Q V P V P T P V Q V P V P T P V Q V P Y P V E K I V D R P V P H E V V R V V E R R V E V P
451 Y D V P V P V I E T V Q V P H E V I R T V E V P F P V E Q I V E K I V E K I V H V P V P V T R Q E H
501 V T R Q V V Q N T L Q Q T R A P P Q Q L G T Q V L P G R D L G V T S G I V R G G P A A Y G A Y G A T
551 P V A A A P V T Y G S V A A P V T Y G A A A P V T Y G A A P V T A Y G A P P V Y A G A P A V V T A P
601 S S V T R Y G Y G P P R P A Y P A T V P V A Y A A A P T E L R N D D F D W P P V P S V C G V Q R Q P
651 Q
```

Articulin

Predicted secondary structure next slide

MSWVAAQTQQGAFALDAADGRIDGKYFGSNVGVSVPGAPVTYGAAAPVTSYA
 APAAVTSFAATAAPVTSFGAYGAYAPACPPCATGPRVVNDPLETRFVDVVKQV
 ETIRSVDVPVPHEVVRTVDVPEHYDVPVPHAVHVQVPYPV

DKFVDVPVPHTIQKIVETRVPYPVQQVVQRRV
 ERPYDVPVV
 ERVNVVPVEQVV
 ERRVPVPVEQIV
 ERVQVPV
 ERLV
 EKVQVHRQVPVPVRVEVPVPHEVIRTVDVPVPHEVVRTQDVPVPVEQIV
 EKVQVPVPVQKKVIQHVQVPYPVQKIV
 DRPVPYPV
 EKIVEQKVPYAVQKVI
 DRPVPYPVQKIV
 ERRVDVPVEVKVRQEVRVPYPVQKIV
 DRPEPYPV
 DKVVEVPQYPVPVQKVV
 ERRVEVPHVIQVREEVRVPYTV
 DKVV
 DRPVPYPVT
 KEVRYV
 DRPVPQPYEVRVPQPYEVKVPVEQIRY
 RDVPVPV
 ERIV
 EKVQVPVPRQVPVKQIQQVPVPV
 EKIVHVQRPYPVQKVVV
 REVVKHVQVPHEIIQRVEHVQHVVQPVEVIEEPVIQQIV

TINKRTIQGQPYVANTTTNVIGQEVRGAYAAGAYGATVTGPAVATGYGQSVTGYGQ
 SVTGYGYGAAAPVAVAGSQGGALSLDLADGRDLGRFYGAPVVSASPF

Predicted secondary structure

1 M S W V A A Q T Q Q G A F A L D A A D G R I D G K Y F G S N V G V S V P G A P V T Y G A A A P V T S 5
51 Y A A P A A V T S F A A T A A P V T S F G A Y G A Y A P A C P P C A T G P R V V N D P L E T R F V D 1
101 V V K Q V E T I R S V D V P V P H E V V R T V D V P E H Y D V P V P H A V H V Q V P Y P V D K F V D 1
151 V P V P H T I Q K I V E T R V P Y P V Q Q V V Q R R V E R P Y D V P V V E R V N V P Y P V E Q V V E 2
201 R R V P V P V E Q I V E R V V Q V P V E R L V E K V V Q V H R Q V P V P V R V E V P V P H E V I R T 2
251 V D V P V P H E V V R T Q D V P V P V E Q I V E K V V Q V P V P V Q K K V I Q H V Q V P Y P V Q K I 3
301 V D R P V P Y P V E K I V E Q K V P Y A V Q K V I D R P V P Y P V Q K I V E R R V D V P V E V K V R 3
351 Q E V R V P Y P V Q K I V D R P E P Y P V D K V V E V P Q P Y P V Q K V V E R R V E V P H V I Q V R 4
401 E E V R V P Y T V D K V V D R P V P Y P V T K E V V R Y V D R P V P Q P Y E V R V P Q P Y E V K V P 4
451 V E Q I R Y R D V P V P V E R I V E K V V Q V P V P R Q V P V K Q I Q Q V P V P V E K I V H V Q R P 5
501 Y P V Q K V V V R E V V K H V Q V P H E I I Q R V E H V Q H V V Q P V E V I E E P V I Q Q I V T I N 5
551 K R T I Q G Q P Y V A N T T T N V I G Q E V R G A Y A A G A Y G A T V T G P A V A T G Y G Q S V T G 6
601 Y G Q S V T G Y G Y G A A A P V A V A G S Q G G A L S L D L A D G R L D G R F Y G A P V V S A S P F 6
651

6

Euglena Gene.3890_GEFR01000334.1
(First third of sequence)

Articulin. Predicted secondary structure at end of sequence

...FFFELFISTMSKPVTTTSGRVRYDGGSQMRVLGGATSGSASSTPIFVSETSP
ITTTHYGGSVTRTYGASTVGVGGISSHYVGGTTVAGSSLAGSTRTHAVGGTAIA
GSSALAGSSLTGSAAFGTGGSVLVGGSALVGGASSVGGASVGGSRVRHIG
GEGPIYVDSGYTTSGYSNVCCAPAGVRSTFVGSRRREVLEAATVTQHI

EKVPVEIDVDIDTY
REIVEVPI
EKIVEVPYPV
EKIVQVPV
DRVVEVPV
ERLV
ERIVRVPEVRTVEVPVEHC
EKLVTVPEIL
EREVAI
EKIV
ERVVRVPEIHEIEVPV
EKLI
ERILRVPEIHEIEVPV
ERIV
EKVIRVPEIRQVEVPV
ERII
EKIIRVPEIHEIEVPV
EKRIETVVRVPEVHQIEVPV
EKII
ERVVRVPEVHQIEVPV
EKIV
ERVVQTHT
RDVIEVPV
EKVV
ERIVRVPQETLIEVPV
EKIV
ERVVRVPQHTTVEVPV
EKVI
ERVVQTHT
RDIIEVPV
EKVV
ERVVRVPQETIIEVPV
EKVV
ERIVTVPQHTTVEVPV
EKIV
ERVVQTHT
RDVIEVPV
EKVV
ERIVKVPQETIIEVPV
EKLI
ERVCRVPQETIVEVPV
EKII
ERVSQTHV
RDIVEVPV
EKLV
ERIVTVP
KETVIEVPV
EKII
ERIVKVPQETIVEVPV
EKLV
ERVQHTHT

Euglena Gene.3890_GEFR01000334.1
(Second third of sequence)

RDIEVPV
EKIV
ERVVKVPQETIIEVPV
EKVI
ERVVTVPQETI
REVPV
EKVV
ERVVQTHT
RDIEVPV
EKII
ERVVKVPQETIVEVPV
ERVI
EKVSTHT
RDIEVPI
EKII
ERIVKVPQETIIEVPV
EKI
ERVVTVHSTDVVEVPV
EKVI
ERIVRVPQETIIEVPV
EKII
ERVCEVQV
EKIVEVPVQQIVRVPVEVPVEHIV
ERVVKVPVETI
KEVPV
EKIVEIIRPVEVIHHIEVPV
ERIV
ERVVEVL
KEVPV
DRVV
ERLVEVPVTHTVEVPVVSVV
EKLVEVP
KEII
REVPV
EKIV
ERVVTNVDNYIEVPV
EKII
ERVVTVPV
EKRIEVPVEVVV
ERI
KEVQV
EKIVEVPV
EKII
EKIVKVPV
EKIVEVPVECVV
ERIVQVPVTQTVEVPV
ERIV
EKIVEVQVQQIVEVPV
EKVV
ERVVKVPVDHVIEVPV
ERVV
ERI
REVPV
EKIIQVPV
EKIKIVTV
DRIVEVPVEV

Euglena Gene.3890_GEFR01000334.1
(Final third of sequence)

EKIVTVQV
EKIVQVPV
EKIV
EKLVEV
EKIV
ERTVTV
ERVVEVPKVV
EKIV
EKIVEVP
REYI
REVTVTSI
KEVPV
EKIV
ERIVEVPITNVVEVPV
ERVV
ERVVEIPL
ERVV
ERIVEVPV
EKII
ERIEVPV
ERVVEHII EVPV
ERLV
ERIVEIPV
EKII
EKIIEVPI
EKII
ERVIEVEVPIEI
EKVVEI
EKIVEVPVWV
EKVVEV
EKIVEVSGGVEIV
EKIVEV
EKIVEVPVWV
EKVIET
EKII
EKGVV

TEIVPEVSTELYSTSYSTGSKGTSSSISASGGASYSVGGGYGVSGGASTGGGYAVS
GGTGGAGYGSIGTSGGTGYSSITTGGSGTGYSSIGGTGAGVGYSSATTGGAGAGA
GTGYSSISGSGAPGGAGGYSSTSYSSGTYGAGGGTYGAGAPGAGGYSYSSTSY
SGTGGSGSAPGRTGYTATSGSGTRSGSATRSGSGSRA

Predicted secondary structure

```

1 F F F E L F I S T M S K P V T T T S G R V R Y D G G S Q M R V L G G A T S G S A S S T P I F V S E T
51 S P I T T T H Y G G S V T R T Y G A S T V G V G G I S S H Y V G G T T V A G S S L A G S T R T H A V
101 G G T A I A G S S A L A G S S L T G S A A F G T G G S V L V G G S A L V G G A S S V V G G A S V V G
151 G S R V R H I G G E G P I Y V D S G Y T T S G Y S N V C C A P A G V R S T F V G S R R E V L E A A T
201 V T Q H I

1 E K V P V E I D V D I D T Y R E I V E V P I E K I V E V P Y P V E K I V Q V P V D R V V E V P V E R
51 L V E R I V R V P E V R T V E V P V E H C I E K L V T V P E I L E R E V A I E K I V E R V V R V P E
101 I H E I E V P V E K L I E R I L R V P E I H E I E V P V E R I V E K V I R V P E I R Q V E V P V E R
151 I I E K I I R V P E I H E I E V P V E K R I E T V V R V P E V H Q I E V P V E K I I E R V V R V P E
201 V H Q I E V P V E K I V E R V V Q T H T R D V I E V P V E K V V E R I V R V P Q E T L I E V P V E K
251 I V E R V V R V P Q H T T V E V P V E K V I E R V V Q T H T R D I I E V P V E K V V E R V V R V P Q
301 E T I I E V P V E K V V E R I V T V P Q H T T V E V P V E K I V E R V V Q T H T R D V I E V P V E K
351 V V E R I V K V P Q E T I I E V P V E K L I E R V C R V P Q E T I V E V P V E K I I E R V S Q T H V
401 R D I V E V P V E K L V E R I V T V P K E T V I E V P V E K I I E R I V K V P Q E T I V E V P V E K
451 L V E R V C Q T H T

1 R D I I E V P V E K I V E R V V K V P Q E T I I E V P V E K V I E R V V T V P Q E T I R E V P V E K
51 V V E R V V Q T H T R D I I E V P V E K I I E R V V K V P Q E T I V E V P V E R V I E K V C S T H T
101 R D I I E V P I E K I I E R I V K V P Q E T I I E V P V E K C I E R V V T V H S T D V V E V P V E K
151 V I E R I V R V P Q E T I I E V P V E K I I E R V C E V Q V E K I V E V P V Q Q I V R V P V E V P V
201 E H I V E R V V K V P V E T I K E V P V E K I V E I I R P V E V I H H I E V P V E R I V E R V V E V
251 L K E V P V D R V V E R L V E V P V T H T V E V P V V S V V E K L V E V P K E I I R E V P V E K I V
301 E R V V T V N V D N Y I E V P V E K I I E R V V T V P V E K R I E V P V E V V E R I K E V Q V E K
351 I V E V P V E K I I E K I V K V P V E K I V E V P V E C V V E R I V Q V P V T Q T V E V P V E R I V
401 E K I V E V Q V Q Q I V E V P V E K V V E R V V K V P V D H V I E V P V E R V V E R I R E V P V E K
451 I I Q V P V E K I K I V T V D R I V E V P V E V

1 E K I V T V Q V E K I V Q V P V E K I V E K L V E V E K I V E R T V T V E R V V E V P K V V E K I V
51 E K I V E V P R E Y I R E V T V T S I K E V P V E K I V E R I V E V P I T N V V E V P V E R V V E R
101 V V E I P L E R V V E R I V E V P V E K I I E R I I E V P V E R V V E H I I E V P V E R L V E R I V
151 E I P V E K I I E K I I E V P I E K I I E R V I E V E V P I E I E K V V E I E K I V E V P V V V E K
201 V V E V E K I V E V S G G V E I V E K I V E V E K I V E V P V V V E K V I E T E K I I E K G V V

1 T E I V P E V S T E L Y S T S Y S T G S K G T S S S I S A S G G A S Y S V G G G Y G V S G G A S T G
51 G G Y A V S G G T G G A G Y G S I G T S G G T G Y S S I T T G G S G T G Y S S I G G T G A G V G Y S
101 S A T T G G A G A G A G T G Y S S I S G S G A P G G A G G Y S S T S Y S S G T Y G A G G G T Y G A G
151 A P G A G G Y G S Y S S T S Y S G T G G S G S A P G R T G Y T A T S G S G T R S G S A T R S G S G S
201 R A

```

N.B Sequence submitted to PSPRED in 5 sectors. Longer sectors yield predictions of no secondary structure using both PSIPRED and YASPIN.

Euglena 15434_G

(First half of sequence)

Articulin

MTVEESRTFSQTPITTYQSGARAVSNTRILSTGATTLSSATSIVDGSRRYAR
 SAAVCCQPFDP^CNSVVTRQLQTRYLDPVETTHVEKVPMEIEVDVEGYRE
 VIEVPV

ERIVEVPV
 ERIRTVPMQVVEVPVNRHV
 ERLVVPEL
 KERQVFIPKTI
 ERVVKVPEV^YY^EVPVQRTV
 EKVVVP
 REVLVEVPV
 DKVTTV
 DRVV
 EKHVPVEHIV
 ERLVEQ
 ERVVETQVKVQQIVEVP
 REVIKN^YEV
 ERIV
 ERVVEVPQL
 KEVQVPV
 ERV^Y
 ERIVPVPK
 RDTVQIPV
 ERIH
 ERIVPVPQTNIVEVPV
 DKLI
 ERV^CTVPKQ^CTVEVQV
 EKVV
 DRIVEVPRTITQEVLIDMVV
 EKVVTVPS
 EKIEVPV
 EKVVESIVQIPKVVTHDVPV
 EKVI
 ERVVTVPHEQRIEIQRVV
 ERVVQVPV
 DRVV
 ERVVEVPHVVT
 REVPVESV
 ERVVTVPVNRVVEVPV
 DRY^V
 DRVVQVPV
 EKVVNRVVEIPKVIT
 REVELQQIV
 EKIVEIPTTKVVEVPVIRY^I
 EKIVEVPVT
 REVEVPVQKY^V
 EKIVEIPY^Y
 ERVVEVQVEVPV
 ERIV
 EKVVEVPFTV
 DRVVEISI
 EKVI
 EKIIEV
 ERRVEVPVEVIV
 EKVV

Euglena 15434_G
(Second half of sequence)

EKVVEIEVPVII
EKVVEI
EKIVEV
ERRVEVPFDTIVEV
ERVVEV
EKIVEIPIEVVKNV
ERVV
EKTIE

TSYDVTPVVEQKVSSFSASGGGSYSTSASYSTGGGYASGGGGYVSGGGGGFT
SGGGGYVSGGGYSSGGGTGGSAFIGGGGSGGSAFISGGGSGGSAFISGGGSG
GSAFISGGGGSGSGSRVTGGFSSSGALPGAAGVTGTTQLYSSQTSSLSGSTGA
WGLPKQY

Articulin

...AASPVDDLPGRRPGGGPELVAMAPALGSGVAGGQAVGDQGEVVDLMQIGSE
 QEIVEALEVESFDSRAVHQVALLEEDSPLWRSPVKPNGVV

ERIVPV
 ERLIEVPKVIHCVE
 ERIVEVPIQRIV
 ERLVPMEVTHI
 KEVPV
 ERIV
 EKLVEVPVVTI
 KEVPVETVV
 ERVEVTV
 EKVV
 EKLTV
 EKIVEVPTFVEVIHVKVIEVPV
 ERVI
 EKIV
 ERPVPVSVV
 ERLVEVPVQHVVQVPV
 EKIV
 ERIVEV
 EKPVTVTH
 EKIVEVPV
 DRIV
 EKVVTQEVPRIVEVPV
 ERLV
 ERLVEVPVVTI
 KEVPI
 ERVV
 ERVEVITV
 EKVV
 EKVPV
 EKIVEV
 EKVVRVVE
 DRIVEVPVAVRIQQVV
 ERPVVETRVETI
 REVPV
 REIVQIPVYEVI
 ERI
 KEVQVPYT
 REVPVEHV
 EKIVEFVEVPIYV
 DRIVEVPP

SVRSTSRAQQKISAHQRQCGASRRNGSGRW

Articulin

...YVRHYTDLQPTHREYVEMVPVEVSVQEEIIQDIEVPTFERRRRMVWEEFGG
 AAPTVPFRPVTPERPVPVPLQRPADTFAPVPVPHEGAQTVAPSERHGPV**G**
 PAEHTLQVVRVPQVQRVEVPV

ERVV
 EKVVTVPQVQRVEVPV
 ERLV
 EKVVKVPEVRIVEVPV
 DRVI
 ERRVEVPV
 DRVV
 ERVVKVPEVTTVEVPVETI
 KEKIVTVPVQRTIEVPI
 EKfV
 ERIVKIPEVRTVEIPV
 ERII
 EKIVQVPI
 EKIVEVPV
 DRII
 ERVIEVPVTKVV
 EKVIEVPV
 DRVVEVPV
 EKIV
 EKfVEVEVLV
 ERIV
 ERLVEVEVP
 KEKII
 ERIIEVPV
 ERVI
 EKVVEVYVPV
 ERVVEV
 ERIVEVPV
 ERIV
 ERIVEV

PYSEDLA**G**AQVPSEQMAEK**G**SSEEATRSAQAQSGSHPVSHVSSEVSSSALPPVGS
 SSWMPGSSYASTQPSQSVASSAQSAGAVDVSTTWGPGSGYRPAARSGPSLAGY
 VPSAGARSEPPGTAAGVDFRQYSAQLAPEAEQQRf

Articulin

...FFSPFPLSFLSFTLQSADSQRYMSYSGAYRADYPVRSTSPRLSRPYPFSRPA
 ERVVSPTRVVDGRYRSPSRLVSES^CFPAQSTVVETTSAC^IRHTQEIVE^{VPV}VAQ
 PVEVVRTVPMEVPTLEVVERARVQVQEHVIEVP

REHIV

EKIV

EKVVLKTIE^{VPV}EHIVEKIIE^{VPV}

ERLV

ERIV

ERIVEVPQ

DRLV

ERVI

EKIVE^{VPVE}Y^{YPV}

ERVV

EKIVEVLVEIPV

ERIV

ERIVE^{VPVE}^{VPV}

ERIV

EKIVHVTTEVPRVQRQ

ERIIE^{VPV}DHIVERIVE^{VPV}DRPIS^{AIQ}

ERLVPHTV

ERIVK^{YE}^{VPV}DRVIE^{VPV}EREVI

REVHV

DRPVIHQV

ERPVVHTVTRPVIHEV

ERPVIHEV

ERPVIHEV

ERPVVHRV

ERPV

A^{VPV}SVVRQEPAV^CYETAPYRQRYVNG^{TTTY}STRAPR

Articulin

MPIPKTARDVAAALDNADGVADGKGFHGLPIVVGDSRPSSRNIVYSSPRSTRYV
 SDEFIRTPGLPATRVGYPATRVGYPATRSRYMSGPVIHHHEPPTIVRHSTAPSA
 SRIIRSEGPLVRRERIDQATSRTHIEKRVVPVTRVIEESVEPVRRFVQPVERRV
 EVQPVQTVRIVEVPVERRVVDVQQVEHVVEVEVPVPVEVTKEVKVEVPVEK
 PVYVEHVVEKIVQVDVNHETQVPIEKIVETLRVEEVVHIQVVPVQQTVEKLVTV
 EVPHIIEVPV

EKVV
 EKIVEV
 EKII
 ERIVPV
 EKIV
 EKEVILKIVE
 EKIIIEVPV
 ERVKHV
 EKIVEVPKIIETIV
 EKIVEVPVKHVVEVPV
 EKIV
 EKLVQVPVMYT
 KEVPIETVV
 EKEIPGPV
 REV
 EKIVQVPV
 EKLVEV
 EKIIIEVPVQVPV
 EKIVEVPVEVPV
 ERLV
 EKLVEVPV
 EKIVEVMV
 EKVIEVPVEVPV
 ERII
 EKVVEVPVHV
 EKIVEIPLEI
 EKIV
 EKIIIEVEIPVEIPVEI
 ERIV
 EKLVEVPVVV
 EKVVEVEIPIEI
 EKIIIEI
 EKIVEIPIEI
 EKIVEI
 EKIIIEVEVPIEI
 EKVIEV
 EKIVEVPIEI
 EKVI
 EKTIEV

GTEPTTVTKYSSFSTSSAPIETKRFSSQGMSRSSTQSSSTTPGGEIQYSSSAGSAESTGT
 GTGWGTPGTNRSAWLGNMNN

Articulin

MTVTITPTMTIQQTPGRFTX

EREVG^YRKLH
 RDVLQEIHEVLV
 EKPVEVIQTEVME^{VPV}KHIE
 ERVVAETH
 ERIVE^{VPV}
 EKIV
 ERVV
 EKVVVKTVE^{VPV}
 ERIV
 EKIVE^{VPV}
 ERIVEV^YV
 EKIVEVP
 RDKIV
 ERIV
 EKFVEVPQIV
 ERIV
 EKIE^{VPV}D^YPV
 ERVV
 EKL^{VQ}^{VPVE}^{VPV}
 ERIV
 ERIVE^{VPVE}^{VPV}
 DRIV
 EKIVEVEVE^{VPVIQ}
 DRERVVEVPTDNI
 REKVVE^{VPV}
 ERLVEVPS
 EREVVKQVQ^{VPVE}EIPVEVEPEVDVVIV
 DKA^V
 EKT
 ERPKRITLADQ^YVPP^YSL
 DK^Y
 KERSS

ASYQYTAPEKSVEYVTHAPSEKRYTRSSSG^QVVHEYAAPAKSVEYVTS
 GPERGQVLREYTREYYEDPAHEG^DYRISSGRHGTA^YEAA^RALDAADGRLDG
 KYQGKDIYVNGEKVRS^GXA^{CC}LRRPQ^RLHKSQPAGEQQVAVAITRVQAXV
 TM^CERLSILPRSPSLPPQQGHLVLAGDLHLPLYTQKKQVRDR^VNKATGTGWV
 MAIGNQRQVKGQPSSAIRVLRMSSKKHSMGSAEPTGS^CRVEGENKHSEVR
 GRALESDNSPKHHHLPALTSNQAIFFRNEWLQSSWWANVVSGESL^C

Eutreptiella CAMPEP_0113741250
(First half of sequence)

Articulin

...GAPHSAPISSNCTIDNMSSRVTTGSASGGRVRYDGGSSQMRVRDPKISVTE
TGTPGNEVIVVSAERDVVPTGTVTSSTYQGGSTIVRERGGYPVTSVGHSYFRS
GSPVAVGGSSVARSAVRVGQVQAGGAYYRQGGRALYSQSVRRELLDEATE
TTMI

EKVPIEIDVEVDN^Y
REVVEIPV
EKIVEVP^YVPV
ERVVE^{VPV}VEHTV
EKIITVPEL
REIEIPV
ERIV
EKVVRVPQETLVE^{VPV}
EKVI
ERVVRVPEVTEIE^{VPV}
ERIV
EKVVRVPEVHEIE^{VPV}
EKII
ERVVNTHS
RDIIE^{VPV}
ERVI
ERVVEVEV
EKIVH^{VPV}
EKVVTVPHEVVNSV
EKIVTVEVETI
KE^{VPV}ESVVEVIRPVEVTRHVE^{VPV}EHIV
ERIVEVV
KE^{VPV}
ERV
REKIVQ^{VPV}TNTVE^{VPV}VSVV
EKLVEVP
REIV
RE^{VPV}
EKIV
ERIVQ^{VPV}DN^YIE^{VPV}
EKSV
EKLQIPV
EKTVEVPI
ERVV
EKLVE^{VPV}ETIVE^{VPV}
EKL
RESIKVVTV
EKLIE^{VPV}
ERVV
EKIIQ^{VPV}EHSVQVPI
ERIV
EKIVEVQVQQIIE^{VPV}
EKEV
ERVVK^{VPV}
EKVVEVPI
EKIV
ERIVQ^{VPV}
EKIIEVPI
EKVVEVIKVVQN

Eutreptiella CAMPEP_0113741250
(Second half of sequence)

ERIVEVPVEV
ERLVSVPV
EKIIEVPL
ERIV
EKAVPV
DRII
EKVVEV
ERIVEVPMIV
EKFV
DRTVTV
EKIV
EKIVEVPVTNVVEVPI
ERTV
EKVVEVPV
EKIV
ERIVEVPV
EKIVEV
EKIVENIIEIPI
ERVI
EKIIEVEIPIEI
EKIVEI
EKVVEVEIPIEI
EKIVEI
EKIVEVEIPIEI
EKIVEI
EKIVEIEVPVIV
EKVVEI
EKIVEVPIEI
EKL
ERTIVE

TAPTVETT VTTETVEVPRGAGYSSYSSSSSFSSGSRGAGGGGGGAAGGGSGGSTSYSTWSSQSGX

Eutreptiella CAMPEP_0113737346 (First half of sequence)

Articulin. Head domain a perfect repeat, followed by a bit of a third repeat, followed by a novel sequence (blue font)

...NDGLFAEPDHTLSQLALPGPPGGEAMAVATQRLAGKCLKATVSPKHLVFS
 KLIPALKRAC^CKRSDLLEEMADGLTTIPPSDWVLAKMAREWR^CWLPSGP
 GYFDYASQPVPWELVGGRHWPNSRGDFDERDTGHEGFSTQDFLNRPEAK
 AAGLTAAEVIALRLYTGPYIPINRSLRVNSGRFAVTQWALD^CAIGKLALAERE
 GLLLRGLRLLPKEEWQQQYED^CR^CCADDAMDLWISDPAYSSTTTDMAVATGT
 DFGGP^CTFVFHAQ^CDLSPADGLIGNAASVQWVSQYPDEVERLLPSNS^CFLSL
 PQGMRSELPEGMGDRKIFQFFSRFLWDYER^{CC}PPVVTEVEEYVVRINEVMW
 HVYRGLGSEKPPPEVTADAELTAFFLAPDPVGVDGMDGGLPEVLEEEM

NDGLFAEPDHTLSQLALPGPPGGEAMAVATQRLAGKCLKATVSPKHLVFSKLI
 PALKRAC^CKRSDLLEEMADGLTTIPPSDWVLAKMAREWR^CWLPSGPGYF
 DYASQPVPWELVGGRHWPNSRGDFDERDTGHEGFSTQDFLNRPEAKAAG
 LTAAEVIALRLYTGPYIPINRSLRVNSGRFAVTQWALD^CAIGKLALAEREGLL
 LRGLRLLPKEEWQQQYED^CR^CCADDAMDLWISDPAYSSTTTDMAVATGTDF
 GGP^CTFVFHAQ^CDLSPADGLIGNAASVQWVSQYPDEVERLLPSNS^CFLSLP
 QGMRSELPEGMGDRKIFQFFSRFLWDYER^{CC}PPVVTEVEEYVVRINEVMW
 HVYRGLGSEKPPPEVTADAELTAFFLAPDPVGVDGMDGGLPEVLEEEM

NDGLFAEPDHTLSQLALPGPPGGEAMA

ASSLN^GLPTF^CAVSQEDPVAVETRIEVPTVE^KVTAVQAEDIQEVPMVE^{VPV}

ERVIE^{VPV}
 EKII
 ERI
 ERIPVETIRE^{VPVEV}
 EKIVH^{VPVEHV}
 EKVVE^{VPV}
 ERVV
 EKVLQ^{VPVEV}
 EKVV
 EKVVE^{VPV}
 EKIV
 EKVFQ^{VPVEV}
 EKIIH^{VPVEHV}
 EKVVE^{VPV}
 ERVV
 EKVLQ^{VPVEV}
 EKIIH^{VPVENV}
 EKVVE^{VPV}
 ERVV
 EKVLQ^{VPVEV}

Eutreptiella CAMPEP_0113737346
(Second half of sequence)

EKVV
EKVVEVPV
EKVV
EKVEKVPVEV
EKIIHVPVEHV
EKVVEVPV
ERVV
EKVLQVPVEV
EKIIHVPVXXVV
EKVVEVPV
ERVV
EKVLQVPVEV
EKVV
EKVVEVPV
EKVV
EKV
EKVPVEV
EKIIHVPVEHV
EKVVEVPV
ERV

Articulin. Predicted secondary structure next slide

MAEARTGPSAMDALEDARRRATRQTRSAAEALDAVDGVMGKFFGRPIVAT
 GPSRVVGGRSVVDSIDAFDGRRRYRSSGGVHGSAAEALDAADGVMGDRFYGR
 PIVETRSPQRLRGVTAQALDAADGVIDGRFYGRPIVETRGATTRYDDHVEVRR
 VHGGTRHRSKVVEVPVYHHENRIPVEVQVPVQIPVEVPVRI

ERPLAI
 ERVVEVP
 REVIKHVEVLVEVPHEVKVPYAV
 EKVI
 ERIVENVVEVPGPTKV
 EKIVEV
 EKIVEVEVQV
 ERI
 REVKVPYPV
 ERVV
 EKL
 ERTYPVDRIVEVPV
 ERIVHDHVDVP
 RERIV
 ERIVEVPYPV
 EKIV
 EKISEVADVRVV
 EKVVQIPV
 EKVV
 EKVTTVPI
 EKIVKVPI
 EKPVIKLVH
 KEVPVEHIV
 EKIVEVPV
 EKIEVEVDVPV
 EKVV
 EKVTRIPVEIPI
 EKFEVPI
 EKIVEVPVEHV
 EKIVEVPV
 ERII
 EKLKVMVPVPGPPHENIY
 RDREVPI
 EKRVY
 RDRIW

GGRQRGYEMPPPPPPPRDVQWISVPEAYPHQTWEWDRWQNWWDGGRYGSNMQY
 DGGYNAGYPPAAPQSYSNMQYDGGYNAGYPPAAPQSYRYGA

Predicted secondary structure

```

1  M A E A R T G P S A M D A L E D A D R R T A T R Q T R S A A E A L D A V D G V M D G K F F G R P I V
51  A T G P S R V V G G R S V V D S I D A F D G R R Y R S S G G V H G S A A E A L D A A D G V M D G R F
101 Y G R P I V E T R S P Q R L R G V T A Q A L D A A D G V I D G R F Y G R P I V E T R G A T T R Y D D
151 H V E V R R V H G G T R H R S K V V E V P V Y H H E N R I P V E V Q V P V Q I P V E V P V R I E R P
201 L A I E R V V E V P R E V I K H V E V L V E V P H E V K V P Y A V E K V I E R I V E N V V E V P G P
251 T K V V E K I V E V E K I V E V E V Q V E R I R E V K V P Y P V E R V V E K L V E R T Y P V D R I V
301 E V P V E R I V H D H V D V P R E R I V E R I V E V P Y P V E K I V E K I S E V A D V R V V E K V V
351 Q I P V E K V V E K V T T V P I E K I V K V P I E K P V I K L V H K E V P V E H I V E K I V E V P V
401 E K I I E V E V D V P V E K V V E K V T R I P V E I P I E K F V E V P I E K I V E V P V E H V V E K
451 I V E V P V E R I I E K L V K V M V P V P G P P H E N I I Y R D R E V P I E K R V Y R D R I W G G R
501 Q R G Y E M P P P P P P P R D V Q W I S V P E A Y P H Q T W E N D R W Q N W W D G R Y G S N M Q Y
551 D G G Y N A G Y P P A A P Q S Y S N M Q Y D G G Y N A G Y P P A A P Q S Y R Y G A

```

Articulin. Predicted secondary structure next slide

...PQVTE

REVVV
EKVV
EKILRVPEIHEFQVPI
EKIV
EKIVEVPV
EKVVEVPV
KEVVVEVMEV
EKIV
EKLVVV
ERIV
ERPVEV
EKIIEVPKVV
ERVVHKMVEVPVEVT
REVPV
EKIV
ERIVEVPNIT
EKLVVV
EKIVEVPKVV
EKVVVL
EKIVEVPKVVVEQVVVV
EKIVEVPKVI
DRVV
EKSVVV
EKIVEVPKVV
EKLI
EKVVI
EKVIEVPKVV
EKL
EKIVV
EKVVEVV
KEVKS
KDLIEVPV
ERVV
EKVMEVSNTHVVEVPI
EKII
EKVVEIPIEVVVEHVVEVPV
EKVVKFPV
EKIVEVPVEIYI
EKII
EKVVEVEVPVEV
ERILEV
EKAVEV
EKLVEVPI
EKTIEVMKTV
EKIVEVPV
EKII
EKIV
EKVVEI
EKVVEVPV
EKVVEI
EKVTH

Eutreptiella CAMPEP_0113747498
(second half of sequence)

SITQIPPAIETTFSFERESQSEGRHQYADVSSGESVERYHSSMSSSTHGERY
QVVESMTTSSPSIVSGGATSPRVGQSIRVADTTNYGQSSSVMSTTSSERDGP
MEVVEGSPRRYFSSISPRGTSEYKSIKAGAI

```

1 P Q V T E R E V V V E K V V E K I L R V P E I H E F Q V P I E K I V E K I V E V P V E K V V E V P V
51 K E V V V E V M E V E K I V E K L V V V E R I V E R P V E V E K I I E V P K V V E R V V H K M V E V
101 P V E V T R E V P V E K I V E R I V E V P N I T E K L V V V E K I V E V P K V V E K V V V L E K I V
151 E V P K V V E Q V V V V E K I V E V P K V I D R V V E K S V V V E K I V E V P K V V E K L I E K V V
201 V I E K V I E V P K V V E K L V E K I V V V E K V V E V V K E V K S K D L I E V P V E R V V E K V M
251 E V S N T H V V E V P I E K I I E K V V E I P I E V V V E H V V E V P V E K V V K F P V E K I V E V
301 P V E I Y I E K I I E K V V E V E V P V E V E R I L E V E K A V E V E K L V E V P I E K T I E V M K
351 T V E K I V E V P V E K I I E K I V E K V V E I E K V V E V P V E K V V E I E K V T H S I T Q I P P
401 A I E T T F S F E R E S Q S E G R H Q Y A D V S S G E S V E R Y H S S M S S S T G E R Y Q V V E S M
451 T T S S P S I V S G G A T S P R V G Q S I R V A D T T N Y G Q S S S V M S T T S S E R D G P M E V V
501 E G S P R R Y F S S I S P R G T S E Y K S I K A G A I

```


MSIYVGEQRYSSPQQYAAERQRYVSEPRPLDHIRVADAGYKYGTGDSKVVQ**G**
 KRVRSQVVN**VPV**VHHDVR**VPV**EVEYPVE**VPV**VPYEVSIPREVIREVT**VPV**VERI
 VESVN

ERKVI
 KEIPVPQEFETIV
 EKVV
 EKRVSVP**G**PKY**Y**V
 EKIVEV
 EKIVEV
 ERPVWRI
 KEIPVP**Y**PV
 EKVV
 EKIVNVPKY**Y**I
 DRKVEVPQDRLVQVP
 RERE**Y**I
 REVV
 KEVRVP**Y**PV
 ERIV
 EKVVKVPQ
 ERFV
 EKKVPFIV
 EKVVEHV
 EKVEFT
 KEVHVTV
 EKEVIEEVEVR**VPV**
 EKVI
 EKRVPKII
 EKVVEI
 EKPRIVQKIV
 EKTVMHP**Y**
 EKIT
 EKTIEVPI
 ERHVH**VPV**
 ERVI
 EKQIEVPI
 ERVV
 EKVVQVPQEVD

GGDVDKVVWRDIPVFDKVVYNDYTESGHLFAYHNMPTPPPPLPNFKWVPYSYNALTLGT
 GRSKGGPPLARSLVPPRGPPVPARLPGGRVGTIGRLPLPLGDARXXGNLQNLNIIDRFIF
 VDPEGNVVDERADLWTERGLDPNIMHQKVSAAQAAGFHRSD**K**FWEHHFDELAHAAAQR
 PIVDTQHFDHPVELNYHGHVSPTGRAYYPXLLAGGTATGTRVPTSLPAVLTVSLRGAASAL
 RVWLPTDEDSGVAGALAGTGLLAAPLGPRAFDSAAPNMDEEAV**C**GGGAPPHPVVGQPEP
 REEGSX

Eutreptiella CAMPEP_0113750188
(First half of sequence)

Articulin

MHAAAPRTPPRQGGSVTRPGSGVXXXXXXXXPHPHQLSGGVP**C**REAGSKNASHDGMV
MQVVQESSQMGSRLAPSYNRGYLTSQEVAVARSRGSTMPSSQSSSLTRSDVIVVDE
GRYYDSGMLSHPGSRVRH**GGYS**

ERSVV
EKREIEIEVEVDNI
REIVEIPV
EKIIE**VPV**
ERVV
EKIVQVPQVT
ERDVLV
EKVV
ERIVQ**VPV**ETRVEVPMVQTI
EKVMAVQVENV
RE**VPV**
EK**Y**R
EKLVE**VPV**
ERVIEV**G**V
EKLV
ERVVQVHVDNI
RE**VPV**
EKVI
EKLQIPVETI
RE**VPV**
ERIV
EKVVEVR**VPV**VPV
ERVI
EKIVE**Y**PVQNI
REVVV
EKVV
ERI
RE**VPV**
ERLVE**VPV**
ERIV
ERIVEVPIHNI
KE**VPV**VRLV
EKVVQ**VPV**EIPV
ERVI
EKIVKVPVENI
RE**VPV**
EKIV
ERVVE**VPV**
ERVV
ERVV

Eutreptiella CAMPEP_0113750188
(Second half of sequence)

EKQVPVEV
EKIVPVENIV
EKLVTVQVEHI
KEVPV
EKIV
ERIVEVPV
ERVV
EKIV
EKVV
ERVVEIPIEI
EKIV
ERVEVVEVPIEI
EKII
ERTVEVEVEVPVIV
EKVVEI
EKIVEVPIEI
EKIV
EKTIVE

SVGSA^GSAPSASVTVTTETSSSSASGYSSGGSGYVAGSWSSASKEQAAPGGANGYS
AFSNSAHRHANMLAALDAADGKLDGQVFSADVKTILPGSAPTAYVPPPTTRDKL
AXEGQPQCQPIHRQPPVWLQQFRDPNPRIRVDS

Eutreptiella CAMPEP_0113750414 (First third of sequence)

Articulin

XVMSGATQHPSYATYPSLQAPTYVGLASSAHHAAAAYRGPYASLQNTVEAS
GIGALMPATISGPPPLLPTAQYQPPSGQEVGGQWMEEQSTVTPGNMQVESL
TEISTHASKAKKSNKVSSGSGRSGLVRYVEVPV

EKVV
EKVVEV
EKIV
EKAVPV
ERIV
EKVVEV
ERIV
EKVIPV
EKIVTIPQVTT
REEKIV
DRVVKVPQKHEVVVPVEHIV
ERVVQVPTTV
EKVVVV
EKVV
ERPVEIPVATEHIVVV
EKVV
ERVVRVPEVLEVPVEHIVKQI
EKVPEVRTVEVPV
ERII
ERLVEVPQVT
EKEVPV
ERIV
ERIVKVPEVHTVEVPV
ERLL
EKLVEVPLVV
EREVPV
ERLV
ERIQLVP
EVRTVEVPVEVV
RERIVKVPEVHTVEVPV
EKVV
ERLVEVPQIIE
KVVVNTV
ERIQTVPEVCTIEVPV
EKRIETINKVPMETI
KEVPV
EKII
ERIVEIPQLTTIEVPV
EKVV
ERVQQVTVQNI
KEVPV
EKXV
ERIRKVPMTI
KEVPV
ERIH
EKIVEVPQVLTIEVPV
EKVV
ERVQQVVVQNI
KEVPV
EKL
ERIQKVPVETI

Eutreptiella CAMPEP_0113750414
(Second third of sequence)

KEVPVENII
EKVVQVITQNV
REVPV
EKT
ERIMKVPVETV
REVPI
EKTI
ERLQKVPMETI
REVPV
EKII
EKIVEVPLVI
ERAIPV
EKVV
EKIVEVPQLTTI
ERPVETII
ERVVQVNTENI
REVPV
EKTI
ERVHKVPMETI
KEVLV
EKVV
EKIVEVPQVI
ERAVPVETIV
ERILEVQVQNI
KEVPV
EKVV
ERVVQVPVEII
KEVPV
EKVMTVEVAKIVHQPVEVI
KEVPI
EKVVKVPEVHV
KEVPV
EKVI
ERIVQVPVETV
KEIPV
EKIVNVEVAKIVNVPVEVL
KEVPV
ERIVKVPEIHV
REAV
EKLI
EKLVEVPV
EKII
EKHIYIDRHVEVI
KEVPV
EKIV
ERIVEVP
KETVKMVVV
EKVV
ERVVEVP
KETVRSVVV
EKVV
EKVVKVPQLHEIEVAVEQIY
EKIVEVPGETT
KEVPI
EKII
ERIVEVPV

Eutreptiella CAMPEP_0113750414
(Final third of sequence)

ERII
EKAVTV
EKIVEVPVRI
ERTVEVPVTEIVEV
EKLVIYVNVV

HHAEGPVQMDDRVSDFDWSCAFVDCSNHVVPGYRFCATHQVDHGLSRSANSL
NIPIPHVGVPEPHKRTEKDYQALADMIREDGERRRSANSSSISQAHSSASSFSSIT
RMSGATKTKSYEGIKSEPIIARLTTDDGTRTINIPMHRVVVLGSSSELSLEYRFATF
EPQHCSLSSTSAGIVLRDLSDSGVVFVNEKAVGKQVSVMLQDKDRVAIGKARVSF
VFHPT

Articulin

...PQYSVQPLQMKQAEYQPAIYSMPQMQQAPAPVQENRVEYKEVIKEVEVPV
 PVEVPVYQDEPEYYEVPEYVRVEVPVPHEVPRYVEVPQPYTV

EKVVEVPEPYEVV
 KEVTVPVNVTYEVEVETVNVVPY PVEQVV
 EKVVVVP
 KEQIV
 EKVVEVPV
 EKVV
 EKVVEVPQSVSVPYVVEIPVPYET
 EKII EVPKPYEVIKYIEVPTPVTKEV
 EKIV
 EK RIPNYVKV
 REEIKVPSRRPKLWIKNVPV
 ERV
 ERVEVPYKVHKYVDVRKPYEVIKKI
 EREV
 EKVVEVKVL
 KEKAVPHKI
 DKIVEIPLPYV
 EKIVDVPQPYTVKKII
 EKEVEVPYTVLV
 KEEVRKPYEVTKVVNRVPQKV
 DKVVLKYV
 EKPF AQHYEVRVPKYDYDVPVPYEDIKY
 RDVPFPV
 EKL
 DKVVEVPVP
 REVPFKVFNEIPVPV
 EKII EVPKPEPYDTVV
 EKEVIVYQ

HEPHEVIQEIELVEHIVHPIEIIIEEQVIKRKIKKLRKKVYGEVLC DAPFTGQPVDQAPPDGK
 SGKSGKSGKSKKGTGATNPQIESGPPLAFPAAYGMPGMPGMPGMPGMPLFGMPMPY
 GVPMAF

Articulin

XFFRLGVSMTTYEELKARREREKRLFEVPSDVLSDRTTERRVVQPAYVQDD
 MGALAIRARPWTSQSPSTGELSLPITLCPXSVLS

REVPIQVPVEVP

REVPI

DRVFKHT

EKIPVVR

RERKIRVPVTKIVKKPVTKYVDVPVEVIT

EKVLIKRKKKIVEVPV

ERIV

EKPV

EKIIEVKV

EKIVEVPV

ERVV

EKKVNVPYEQVV

EKVVNVPVEILVNKVV

EKVV

EKIVEVPVEQVV

ERTVEVTI

ERIIQKNIEIPV

ERVQENIITIPIIKII

EKVVHVPV

EKII

EREVPVQV

EKIV

EKYDT

SDPIYVDKTVNVVVEPTVHTTTTTYETEYNYDVDGDWGGYGDNTEVTTEVITTEY
 DGDWGGDAGVTTTTTEYTYGDDGGYGGDDGYGGY

MPEPE^GSPADLQAADDIEIDVSVETLVHVDQRQVVDKPFVSVIEHLEE^{VPVE}
QVLESVN

ERLVL
REVRVPI
ERVV
ERVV^{GVAIEVPY}
ERVV
ERVVE^{VPV}
EKVVP
KEIVKVI
EKIVE^{VPV}TRVLKTKKITPV
ERIV
EKRVEVPTETTVIRRVEVP^YTRIV
ERFEEVII
DKIV
EKRVEVFV
DKIV
EKIVHVPHDVII
EKVVEIEEEEEIVE^{VPV}

THHETTYLDVPQRITLQDQ^GYKPPKAVILNHRQPVPLPTQHQPMMHHHQPTQHPIT
PIMLQQPYPPPTPVTALHTDWLRDRVQAAEEKNHWLQGRRYEPLSHDVEAPVRRLS
APDPGPGSGHRMHGSTAVHEPPYAAPSSVSSRTRRSRSADQYAA

Eutreptiella CAMPEP_0113754010
(First half of sequence)

Articulin

...PEVTADAELTAFFLAPDPVGVGDGMDGGLPEVLEEEMNDGLFAEPDHTLSQL
ALPGPPGGEAMAVATQRLAGKLGKATVSPKHLVFSKLIPALKRACRSDLLEEM
EADGLTTIPPPSDWVLAKMAREWRWLPSPGGYFDYASQPVPWELVGGRHW
PNSRGDFDERDTGHEGFSTQDFLNRPEAKAAGLTAAEVIALRLYTGPYIPINR
SLRVNSGRFAVTQWALDCAIGKLALAEREGLLLRGLRLLPKEEWQQQYEDCR
CADDAMDWLWISDPAYSSTTTDMAVATGTDFGGPCTFVFHAQCDLSPADGLIG
NAASVQWVSQYPDEVERLLPSNSCFLSLPQGMRSSELPEGMGDRKIFQFFSRF
LWDYERCCPPVVTEVEEYVVRINEVMWHVYRGLGSEKPPPEVTADAELTAFFL
APDPVGVGDGMDGGLPEVLEEEMNDGLFAEPDHTLSQLALPGPPGGEAMAASS
LNGLPTFCAVSQEDPVAVETRIEVPTVEKVTAVQAEDIQEVPMVEVPV

ERVIEVPV
EKII
ERIE
RIPVETI
REVPVEV
EKIVHVPVEHVV
EKVVEVPV
ERVV
EKVLQVPVEV
EKVV
EKVVEVPV
EKIV
EKVFQVPVEV
EKIIHVPVEHVV
EKVVEVPV
ERVV
EKVLQVPVEV
DKVA
EKVVEDVV
EKVVEVPV
ERVV
EKVLQVP

Eutreptiella CAMPEP_0113754468
(Second half of sequence)

XKVLQVPVEV

EKIIHVPVENVV

EKVVEVPV

ERVV

EKVLQVPVEV

EKIIHVPVEHVV

EKVVEVPV

ERVV

EKVLQVPVEV

DKVA

EKVVEVPVEQAVI

RDVPV

EKTVMIVEGNMLDAAIGLAEEGKPFVFEAANGYGLTARHVEAGQQLKNKR
ADNDRWRLHRMSPGEFTVESLSAPGLRLGVGNQSGHGFKAVLVPEGDQRA
LLRLVPARDGSQQRVSFESVSQPGSLLNH^CNGLMWFFDRPANQHFSNDSSW
VLLSSEKENS^KTSRRGGTAQLRLEGKYLNLDESSGFTIRPRAGQTIVQPVTFT
SAQDVYPGHNPH^TWTLKGTNDXGGMGITSGWGLSLEGS^GFREYIRFENWKA
YRMYRVTPPPGPELHPRKVDADRAVRRPRPGRRQPPAKKKRERKREEGLSQ
GWX

Articulin

MAEARTGPSAMDALEDADRRTATRQTRSAAEALDAVDGVMGKFFGRPIVAT
 GPSRVVGGRSVVDSIDAFDGRRYRSSGGVHGSAAEALDAADGVMGDRFYGR
 PIVETRSPQRLRGVTAQALDAADGVIDGRFYGRPIVETRGATTRYDDHVEVRR
 VHGGTRHRSKVVVVPVYHHENRIPVEVQVPVQIPVEVPVRI

ERPLAI
 ERVVEVP
 REVIKHVEVLVEVPHEVKVPYAV
 EKVI
 ERIVENVVEVPGPTKV
 EKIVEV
 EKIVEVEVQV
 ERI
 REVKVPFLTXV
 ERVV
 EKL
 ERTYPV
 DRIVEVPV
 ERIVHDHVDVP
 RERIV
 ERIVEVPYPV
 EKIV
 EKISEVADVRVV
 EKVVQIPV
 EKVV
 EKVTTVPI
 EKIVKVPI
 EKPVIKLVH
 KEVPVEHIV
 EKIVEVPV
 EKIIIEVEVDVPV
 EKVV
 EKVTRIPVEIPI
 EKFEVPI
 EKIVEVPVEHV
 EKIVEVPV
 ERII
 EKLVKVMVPVGPPHENIY
 RDREVPI
 EKRVY
 RDRIW

GGRQRGYEMPPPPPPPRDQVQWISVPEAYPHQTEWDRWQNWWDGRYGSNMQY
 DGGYNAGYPPAAPQSYSNMQYDGGYNAGYPPAAPQSYRYGA

Articulin. Repeat in ABD domain

MSYEEYGGYGDAEGYGEATTESWTEGYGEGGGGEAYAEGYGEASGAAYGEASG
 AAEAYGEASYSAYGEASGAAEAYGEAYGDGGGGEETHHTDSWVEHSGGGDGGG
 AEQWDQHEWTEEHEWGEPPGGYTSEYKMEALPDVDLGGYSEFTPTTDTYDGQA
 GVISSSYQVEGEGQKVYTDWSYQEPTPIKEYVPYVSPEIEPEGEKVVEMIHGEPV
 LTETKILAENQLDEVIVSVDLLSEGKVSSEVETGRRIVSKQRVIENRVEQVVQVPVET
 FEEHLVVN

EKPKPVRVNV
 DRLVA
 KDV
 DKIIY
 KEVKVPIEVK
 KEVVVKQEVEVPV
 EKVI
 ERIRKVP
 ERII
 EKVV
 EKVV
 EKIVEVPV
 ERIVIEEEIFV
 EKII
 EKIIIVPVPKIV
 EKVIKVPV
 ERIV
 EKVIKVPV
 ERIV
 EKIVEVPV
 ERIV
 EKIMEVPVENII
 EKTVEVLVPSRKQHYVEQVVVE
 ERPVEYT
 RERAPVYVE

Repeat

ERAPVYVE
 ERAPVYVE
 ERAPVYVE
 ERAPVYVE
 ERRVE
 ERRAPVYMK
 ERAPVYVE
 ERAPVYVE
 ERAPVYVE
 ERRVE
 ERHPVYVE
 ERAPVYVE
 ERAPVYVE
 ERGPVYAK
 EREPIY

TSSTYGTRRYADQEPVYANETSASPIYVDERSSAGAYY

Articulin

XFFRLGVSMTTYEELKARREREREKRLFEVPSDVLSDRRTTERRVVQPAYVQDDM
 GALAIRARPWTSPSRDFVGHDKTVSVYGGTVVPSHNTVSTSVLS

REVPIQVPVEVP

REVPI

DRVFKHT

EKIPVPR

RERKIRVPVTKIVKKPVTKYVDVPVEVIT

EKVLIKRKKKIVEVPV

ERIV

EKPV

EKIIEVKV

EKIVEVPV

ERVV

EKKVNVPYEQVV

EKVVNVPVEILVNKVV

EKVV

EKIVEVPVEQVV

ERTVEVTI

ERIIQKNIEIPV

ERVQENIITIPKII

EKVHVPV

EKII

EREVPVQV

EKIV

EKYD

TSDPIYVDKTVNVVVEPTVHTTTTTYETEYNYDVDGDWGGYGDNTEVTTEVITTEY
 DGDWGGDAGVTTTTTEYTYGDDGGYGGDDGYGGY

Articulin

MTTTTGAARALDLADGVEDGKIFYGRRIVEGNSRVLTGSRVYTRSPSRSVVR
 TEGSVIRALDAADGVIDGKIFYGSRIVDGARLGYSTQTEYGSRIVDGSSQVYSSI
 HRPRYSSITRGERRVVDVSSEATYKQKRVEVPVVHHETRVVPEYEVPVHIQVE
 VPVRVVEVPYPV

ERIVEVPQDVVQ
 ERQVIVDVPVEVRVPYPV
 ERVRHVPV
 ERIIEVQGPTRFVENVVHV
 ERLVEVPIEVPVVRTV
 ERYPVEQIVEVMVEQPVPVEHIVQVPEEYIVEEIEVEQ
 ERLVEQVVEVPVPVEHIVEQRVPVPVPHIVEEIVEVPVENIV
 EKLIQIPV
 ERV
 REVPVEQIVRQARYV
 DRPYEVL
 REKVVHVPV
 ERVRHVHVDVPV
 ERIV
 EKIVEVPFEVPV
 EKVVQVPV
 ERVVQVPVEHV
 ERPV

TRYSAPTYAAPPLRPVHHNPAFALDAADGRIDGQYFGARIAQPYNPALALDAADG
 RIDGTYYGSTIAPQPYGPGYRPY

MSIYVGEQRYSSPQQYAAERQRYVSEPRPLDHIRVADAGYKYGTGDSKVVQ**G**
 KRVRSQVVN**VPV**VHHDVR**VPV**EVE**YPVE****VPV**EVPEVSIPREVIREVT**VPV**ERI
 VESVN

ERKVI
 KEIPVPQEFETIV
 EKVV
 EKRVSV**P****G**PK**Y**V
 EKIVEV
 EKIVEV
 ERPVWRI
 KEIPVP**Y**PV
 EKVV
 EKIVNVPK**Y**I
 DRKVEVPQ
 DRLVQVP
RERE**Y**I
 REVV
 KEVRVP**Y**PV
 ERIV
 EKVVKVPQ
 ERFV
 EKKVPFIV
 EKVVEHV
 EKVEFT
 KEVHVTV
EKEVIEEVEVR**VPV**
 EKVI
 EKRVPKII
 EKVVEI
 EKPRIVQKIV
 EKTVMHP**Y**
 EKIT
 EKTIEVPI
 ERHVH**VPV**
 ERVI
 EKQIEVPI
 ERVV
 EKVVQVPQEVD

GGDVDKVVWRDIPVFIDKVVYNDYTESGHLFAYHNMPTPPPPLPNFKWVPYSYNALTL
 GTGRSAKGPVPARLPGGRVGTIGRLPLPLGDARMLGNLGQLNIIDRFIFVDPEGNVVD
 ERADLWTERGLDPNIMHQKVSAAQAAGFHRSDKFEHHEFDELAHAAAQRPIVDTQH
 FDHPVELNYHGHVSPTGRAYYPXLLAGGTATGTRVPTSLPAVLTVSLRGAASALRVWL
 PTDEDSGVAGALAGTGLLAAPLGPRAFDSAAPNMDEEAV**C**GGGAPHPVGGQPEPRE
 EGSX

Articulin

MAEARTGPSAMDALEDADRRTATRQTRSAAEALDAVDGVMGKFFGRPIVAT
 GPSRVVGGRSVVDSIDAFDGRRYRSSGGVHGSAAEALDAADGVMGDRFYGIV
 ETRSPQRLRGVTAQALDAADGVIDGRFYGRPIVETRGATTRYDDHVEVRRVH
GGTRHRSKVV**VPV**YHHENRIPVEVQ**VPV**QIPVE**VPV**RIERPLAIERVVEVPRE
 VIKHVEVLVEVPHEVKVPYAV

EKVI
 ERIVENVVEVPGPTKVV
 EKIVEV
 EKIVEVEVQV
 ERI
 REVKVPY**VPV**
 ERVV
 EKL
 ERT**Y**VPV
 DRIVE**VPV**
 ERIVHDHVDVP
RERIV
 ERIVEVP**Y**VPV
 EKIV
 EKISEV**A**DVRVV
 EKVVQIPV
 EKVV
 EKVTTVPI
 EKIVKVPI
 EKPVIKLVH
 KE**VPV**EHIV
 EKIVE**VPV**
 EKIEVEVD**VPV**
 EKVV
 EKVTRIPVEIPI
 EKfVEVPI
 EKIVE**VPV**EHVV
 EKIVE**VPV**
 ERII
 EKLVKVM**VPV**

GPPHENIIYRDREVPIEKRVYRDRIWGGRQRGYEMPPPPPPPPRDVQWISVPEAY
 PHQTWEWDRWQNWWDGRYGSNMQYDGGYNAGYPPAAPQSYSNMQYDGGYN
 AGYPPAAPQSYRYGA

Articulin

XKVLQVPVEV

EKIIHVPVENVV

EKVVEVPV

ERVV

EKVLQVPVEV

EKIIHVPVEHVV

EKVVEVPV

ERVV

EKVLQVPVEV

DKVA

EKVVEVPVEQAVI

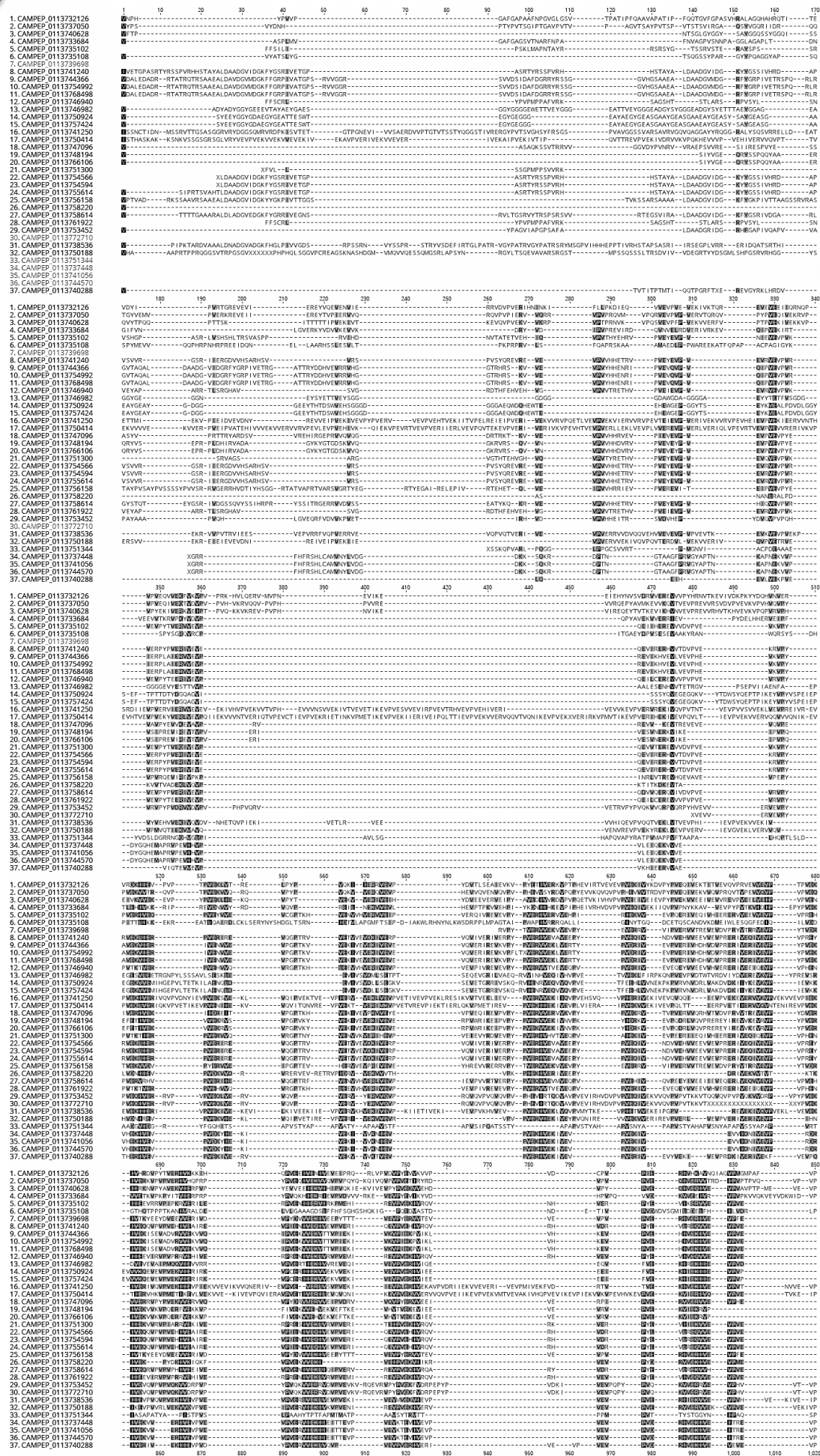
RDVPV

EKTVEMIVEGPNMLDAAIGLAEEGKPFVFEAANGYGLTARHVEAGQQLKNKR
ADNDRWRLHRMSPGEFTVESLSAPGLRLGVGNQSGHGFKAVLVPEGDQRA
LLRLVPARDGSQQRVSFESVSQPGSLLNH^CNGLMWFFDRPANQHFSNDSSW
VLLSSEKENS^KTSRRGGTAQLRLEGKYLNLDESSGFTIRPRAGQTIVQPVTFT
SAQDVYPGHNPHWTWKGTNDXGGMGITSGWGLSLEGSGFREYIRFENWKA
YRMYRVTPPPGPELHPRKVDADRAVRRPRPGRRQPPAKKKRERKREEGLSQ
GWX

Alignment of Eutreptiella articulins into two groups

Group 1

N-terminal half



Alignment of Eutreptiella articulans into two groups

Group 1 C-terminal half

Table showing sequence alignments for Group 1 C-terminal half. It consists of 37 rows of sequence identifiers (e.g., 1. CAMPEP_0113732126) and their corresponding amino acid sequences. The sequences are aligned in columns, with some gaps represented by dashes. The alignment is split into two main sections, each with its own column numbering (e.g., 1, 100, 200, 300, 400, 500, 600, 700, 800, 900, 1,000, 1,100, 1,200, 1,300, 1,400, 1,500, 1,600, 1,700, 1,800, 1,900).

Alignment of *Eutreptiella articulins* into two groups

Group 2

45

