

Swiss Institute of  
Bioinformatics

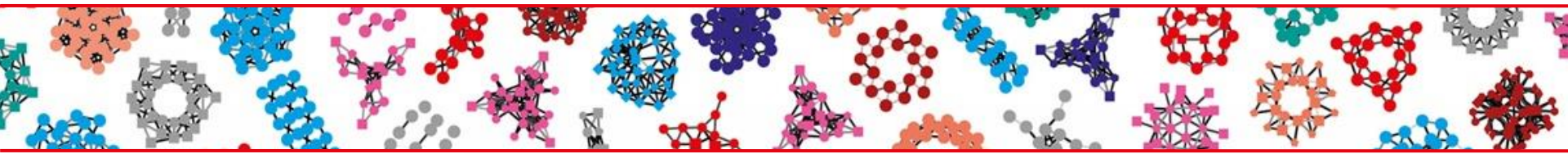
# Automatic annotation in UniProtKB

Protein sequence databases and sequence annotation

Streamed from Geneva, 14 October 2022



[www.sib.swiss](http://www.sib.swiss)



- **Why do we need predictive annotation tools?**
- Protein signatures for homology detection – A short primer
- Annotation rules for functional annotation
- HAMAP and PROSITE - automatic annotation in UniProtKB
- HAMAP and PROSITE - services for external users
- Practical exercises

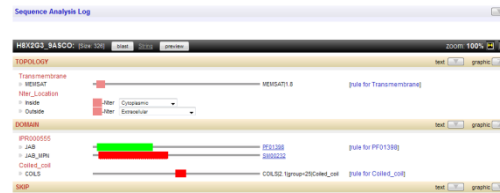
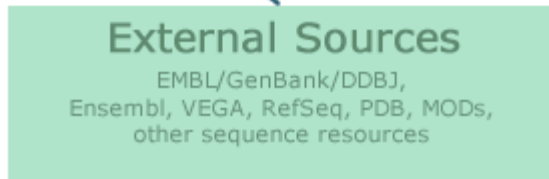
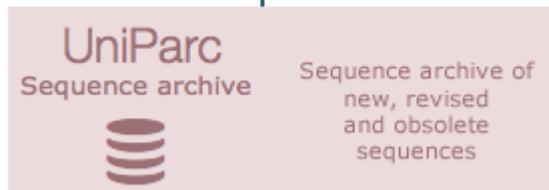
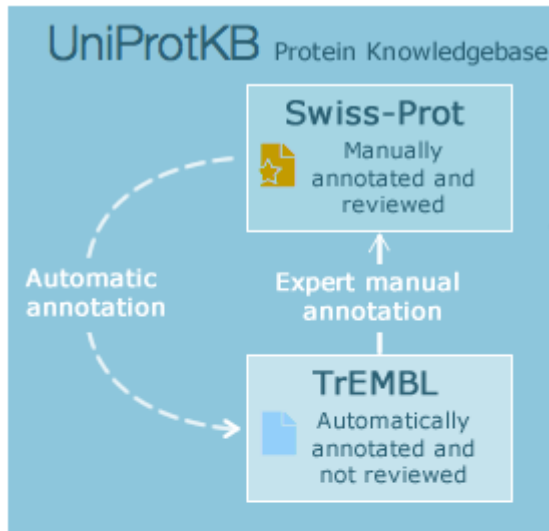
# The UniProt knowledgebase UniProtKB



Literature curation



Curator-supervised  
sequence analysis



- Motif database scans: Prosite, Pfam, Smart, etc.
- Prediction programs: TMHMM, coils, SignalP, TargetP, Repeat, etc.

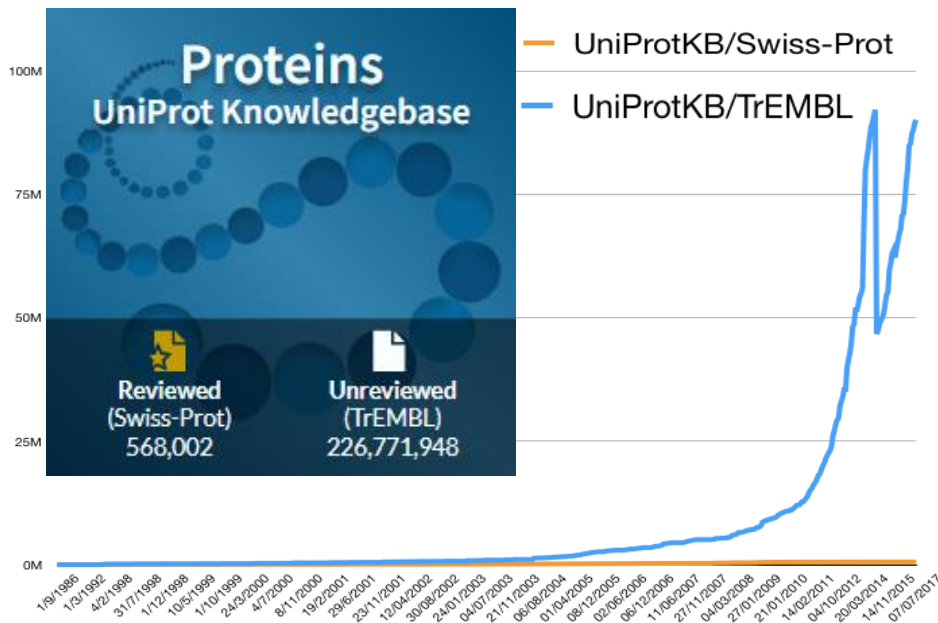
- Protein and gene names
- Function
- Catalytic activity
- Cofactors
- Subcellular location
- Expression
- Protein-protein interactions
- Domains/regions
- Catalytic residues
- PTMs
- Binding sites
- Splice variants

# Why do we need predictive annotation tools?

«Continuing advances in next generation sequencing mean that for every experimentally characterized protein, there are now many hundreds of proteins that will never be experimentally characterized in the laboratory !»

<http://nar.oxfordjournals.org/content/43/D1/D204>

UniProt release 2022\_03:



- Currently, reviewed entries make up less than 1% of UniProtKB
- Manual curation is time-intensive and published experimental data focuses on a rather limited number of model organisms
- Problem: The quality of the functional annotation attached to sequences submitted to UniProtKB is very variable (from original submitters).

# Annotation propagation in UniProtKB/Swiss-Prot

Characterized template entry




```

ADE_YEAST      -----MVSVEFLQELPKCEHHLHLEGTLEPDLDFPLAKRNDIILPE---GFPKSV
ADE_CANGA      -----MVPESFLLLEPKCEHHLHLEGTLEPDLDFPLAKRNNIQLPD---HFPQTF
ADE_KLULA      --MAKFECTDEVTNFLTELPKCEHHLHLEGTLEPELLFQLVERNGVQLPG---TFPKTV
ADE_CANAL      --MAQYECSEHMENFLRELPKCEHVVHLEGTLEPSELLFKLAKRNNITLPE---TFPKTV
ADE_ASPFU      MC-----QSPLHDFLHGLPKCEHVVHLEGCVTPELIFQLAKNNIQLFNPATHPAYASV
ADE_ASPOR      MC-----KSDLHDFLHGLPKCEHVVHLEGC LAPDLIFELAKRNNVSLFN---EPAYESI
ADE_EMENI      MCPNTPYQSQWHAFLHSLPKCEHVVHLEGCLEPPLIFSMARKNNVSLFN---RSSNPAYTSV
ADE_SCHFO      MS-----NLP IYNFIRKLPKCEHVVHLEGCLESPDLVFR LAKKNGITLPE---DDAAYTTP
ADE           -----PVPYVESA
ADE           -----SV
ADE           -----SI
ADE_RHORT      -----MAVDFAPLHALPKVELHLHI EGSLEPEMMVALAERNGLRLP-----SV
ADE_STRCO      -----MKRPYDALMPLPKAELHLHI EGTLEPELAFALAARNVSLP-----DE
ADE_BURPP      -MTT T V T P T P L A E K T A L P K A E L H I H I E G S L E P E L I F A L A E R N G V K L A Y D -----S I
ADE_BURXL      -M T T T V T S T P L A E K T V L P K A E L H I H I E G S L E P E L I F A L A E R N G V K L A Y D -----S I
ADE_CUPTR      ---M T I D A A L A E Q I R R T P K A E L H V H I E G T L E P E L I F R L A Q R N Q V A L P Y P -----S V
ADE_CUPNH      ---M T I D A A L A E Q I R R T P K A E L H V H I E G T L E P E L I F R L A Q R N Q V A L P Y P -----S V
ADE_CUPPJ      ---M T I D A A L A D K I R R T P K A E L H V H I E G T L E P E L I F R L A Q R N H N V L P Y P -----S V
ADE_RALME      ---M T I D A A L A D K I R R T P K A E L H V H I E G T L E P E R I F R L A Q R N N V K L A Y P -----D V
ADE_RALPJ      ---M P I S S A L A E R I A T S P K A E L H I H I E G S L E P E L M F A L A E R N G V K L P Y A -----S V
ADE_RALSO      ---M P I S P A L A E R I A T S P K A E L H I H I E G S L E P E L M F A L A E R N G V K L P Y A -----S V
ADE_GEOLS      ---M N L T N I P R Q A L P E L L C R M P K A E L H I H I E G S L E P E L I F A L A E R N R L Q L A Y P -----T I
ADE_GEOUR      ---M N F D C I P R E D L H G I L C H M P K A E L H I H I E G S L E P E L I F E L A T R N R I Q L P Y P -----T I
    
```

Propagation by sequence similarity


Uncharacterized protein family members in other species




MT2ID\_HUMAN


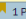

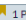

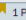
MT2ID\_MOUSE




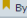

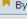
## Function<sup>i</sup>


Protein N-lysine methyltransferase that specifically trimethylates 'Lys-315' of VCP/p97; this modification may decrease VCP ATPase activity.  2 Publications


## Function<sup>i</sup>

Protein N-lysine methyltransferase that specifically trimethylates 'Lys-315' of VCP/p97; this modification may decrease VCP ATPase activity.  By Similarity

		S-adenosyl-L-methionine (UniProtKB   ChEBI  )  1 Publication
▶ Binding site	43	Manual assertion based on experiment (Inferred from experiment) <sup>i</sup> Ref.9
▶ Binding site	75-77	S-adenosyl-L-methionine (UniProtKB   ChEBI  )  1 Publication
▶ Binding site	96	S-adenosyl-L-methionine (UniProtKB   ChEBI  )  1 Publication

		S-adenosyl-L-methionine (UniProtKB   ChEBI  )  By Similarity
▶ Binding site	43	Manual assertion inferred from sequence similarity (Inferred from sequence or structural similarity) <sup>i</sup> Q9H867
▶ Binding site	75-77	S-adenosyl-L-methionine (UniProtKB   ChEBI  )  By Similarity
▶ Binding site	96	S-adenosyl-L-methionine (UniProtKB   ChEBI  )  By Similarity

 2 Publications

 By Similarity

# Microbial proteome annotation

## Statistics

Species	Species code	UniProtKB proteome identifier	Number of entries in UniProtKB		
			All	Swiss-Prot	TrEMBL
<i>Acaryochloris marina</i> (strain MBIC 11017)	ACAM1	UP000000268	8172	322	7850
<i>Accumulibacter phosphatis</i> (strain UW-1)	ACCPU	UP000001619	4438	1	4437
<i>Acetobacter pasteurianus</i> (strain NBRC 3283 / LMG 1513 / CCTM 1153)	ACEP3	UP000000948	2906	4	2902
<i>Acetobacter</i> sp. CAG:267	/	UP000017992	1538	0	1538
<i>Acetobacter</i> sp. CAG:977	/	UP000018259	1775	0	1775
<i>Acetobacterium woodii</i> (strain ATCC 29683 / DSM 1030 / JCM 2381 / KCTC 1655)	ACEWD	UP000007177	3445	4	3441
<b><i>Escherichia coli</i> (strain K12)</b>	<b>ECOLI</b>	<b>UP000000625</b>	<b>4305</b>	<b>4305</b>	<b>0</b>
<i>Acholeplasma laidlawii</i> (strain PG-8A)	ACHLI	UP000008558	1380	147	1233
<b><i>Bacillus subtilis</i> (strain 168)</b>	<b>BACSU</b>	<b>UP000001570</b>	<b>4197</b>	<b>4185</b>	<b>12</b>
<i>Acholeplasma</i> sp. CAG:878	/	UP000017930	1145	0	1145
<b><i>Mycobacterium tuberculosis</i> (strain ATCC 25618 / H37Rv)</b>	<b>MYCTU</b>	<b>UP000001584</b>	<b>3991</b>	<b>2081</b>	<b>1910</b>
<i>Acidaminococcus fermentans</i> (strain ATCC 25085 / DSM 20731 / VR4)	ACIFV	UP000001902	2016	12	2004
<i>Acidaminococcus intestini</i> (strain RyC-MR95)	ACIIR	UP000007093	2386	0	2386
<i>Acidaminococcus</i> sp. CAG:917	/	UP000018048	1276	0	1276
<i>Acidimicrobium ferrooxidans</i> (strain DSM 10331 / JCM 15462 / NBRC 103882 / ICP)	ACIFD	UP000000771	1935	8	1927
<i>Acidiphilium cryptum</i> (strain JF-5)	ACICJ	UP000000245	3521	262	3259
<i>Acidiphilium</i> sp. CAG:727	/	UP000018053	1479	0	1479

# Annotation propagation in UniProtKB/Swiss-Prot

Characterized template entry




```

ADE_YEAST      -----MVSVEFLQELPKCEHHLHLEGTLEPDLDFPLAKRNDIILPE---GFPKSV
ADE_CANGA      -----MVPESFLLLEPKCEHHLHLEGTLEPDLDFPLAKRNNIQLPD---HFPQTF
ADE_KLULA      --MAKFECTDEVTNFLTELPKCEHHLHLEGTLEPELLFQLVERNGVQLPG---TFPKTV
ADE_CANAL      --MAQYECSEHMENFLRELPKCEHVVHLEGTLEPELLFKLAKRNNITLPE---TFPKTV
ADE_ASPFU      MC-----QSPLHDFLHGLPKCEHVVHLEGCVTPELI FQLAKNNIQLNPATHPAYASV
ADE_ASPOR      MC-----KSDLHDFLHGLPKCEHVVHLEGC LAPDLI FELAKRNNVSLFN---EPAYESI
ADE_EMENI      MCPNTPYQSQWHAFLHSLPKCEHVVHLEGCLEPPLIFSMARKNNVSLFN---RSSNPAYTSV
ADE_SCHFO      MS-----NLP IYNFIRKLPKCEHVVHLEGCLESPDLVFR LAKNGITLPE---DDAAYTTP
ADE           -----PVPYESA
ADE           -----SV
ADE           -----SI
ADE_RHORT      -----MAVDFAPLHALPKVELHLHI EGSLEPEMMVALAERNGLRLP-----SV
ADE_STRCO      -----MKRPYDALMPLPKAELHLHI EGTLEPELAFALAARNGVSLP-----DE
ADE_BURFP      -MTTIVTPTPLAECTALAPKAELHLHI EGSLEPELIFALAERNGVKLAYD-----SI
ADE_BURXL      -MTTIVTSTPLAECTVLPKAELHLHI EGSLEPELIFALAERNGVKLAYD-----SI
ADE_CUPTR      ----MTIDAALAEQIRRTPKAELHVHI EGTLEPELIFRLAQRNOVALPYP-----SV
ADE_CUPNH      ----MTIDAALAEQIRRTPKAELHVHI EGTLEPELIFRLAQRNOVALPYP-----SV
ADE_CUPPJ      ----MTIDAALADKIRRTPKAELHVHI EGTLEPELIFRLAQRNNHVLPYP-----SV
ADE_RALME      ----MTIDAALADKIRRTPKAELHVHI EGTLEPERIFRLAQRNNVKLAYP-----DV
ADE_RALPJ      ----MPISSALAERIATSPKAELHLHI EGSLEPELMFALAERNGVKLYA-----SV
ADE_RALSO      ----MPISSALAERIATSPKAELHLHI EGSLEPELMFALAERNGVKLYA-----SV
ADE_GEOLS      -MNLINI PRQALPELLCRMPKAELHLHI EGSLEPELIFALAERNRLQLAYP-----TI
ADE_GEOUR      -MNFDCI PREDLHGILCHMPKAELHLHI EGSLEPELIFELATRNRILQLPYP-----TI
    
```

Propagation by sequence similarity


Uncharacterized protein family members in other species




MT2ID\_HUMAN


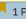

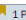

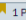
MT2ID\_MOUSE


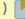

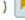

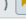
## Function<sup>i</sup>

Protein N-lysine methyltransferase that specifically trimethylates 'Lys-315' of VCP/p97; this modification may decrease VCP ATPase activity.  2 Publications


## Function<sup>i</sup>

Protein N-lysine methyltransferase that specifically trimethylates 'Lys-315' of VCP/p97; this modification may decrease VCP ATPase activity.  By Similarity

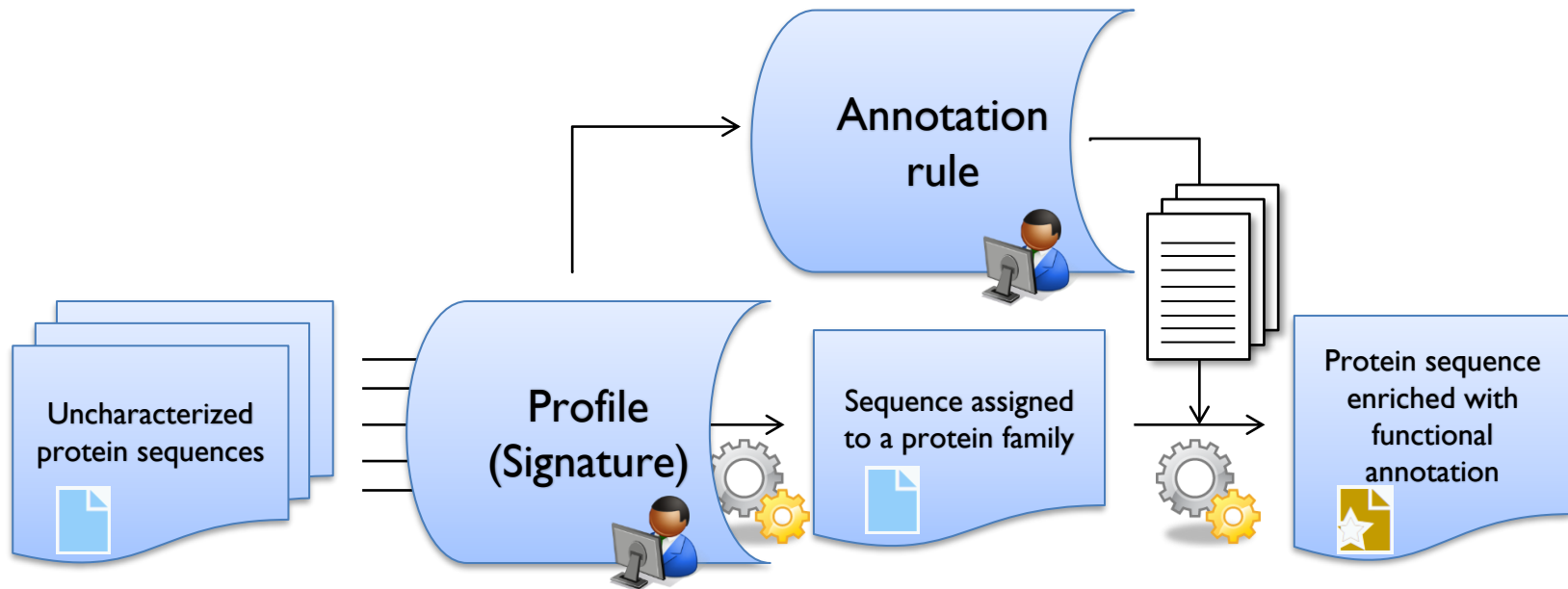
		S-adenosyl-L-methionine (UniProtKB   ChEBI  )  1 Publication
▶ Binding site	43	Manual assertion based on experiment (Inferred from experiment) <sup>i</sup> Ref.9
▶ Binding site	75-77	S-adenosyl-L-methionine (UniProtKB   ChEBI  )  1 Publication
▶ Binding site	96	S-adenosyl-L-methionine (UniProtKB   ChEBI  )  1 Publication

		S-adenosyl-L-methionine (UniProtKB   ChEBI  )  By Similarity
▶ Binding site	43	Manual assertion inferred from sequence similarity (Inferred from sequence or structural similarity) <sup>i</sup> Q9H867
▶ Binding site	75-77	S-adenosyl-L-methionine (UniProtKB   ChEBI  )  By Similarity
▶ Binding site	96	S-adenosyl-L-methionine (UniProtKB   ChEBI  )  By Similarity

 2 Publications

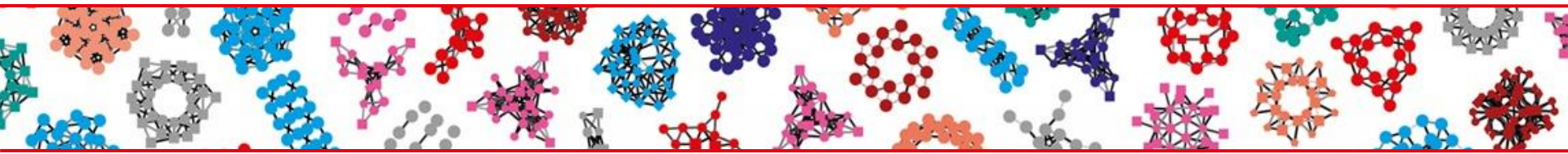
 By Similarity

# Automatic annotation of protein sequences



- Manually curated **signatures** – a protein matching a specific signature may be annotated according to the contents of the associated rule
- **Manually** created annotation **rules** that specify **annotations** AND the **conditions** under which they may be applied
- **HAMAP** signatures/rules for **full length protein sequences**, **PROSITE** signatures/rules mainly for **domains** and **sites**
- Common format and shared syntax





- Why do we need predictive annotation tools?
- **Protein signatures for homology detection – A short primer**
- Annotation rules for functional annotation
- HAMAP and PROSITE - automatic annotation in UniProtKB
- HAMAP and PROSITE - services for external users
- Practical exercises

# Types of sequence similarity

**Family** Groups of proteins that are conserved along the whole sequence, sharing a common evolutionary origin, as reflected in their related functions.



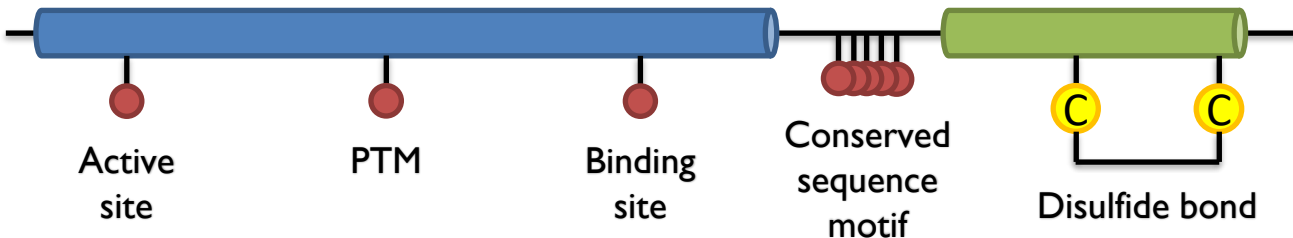
**Domains** Specific combination of secondary structures organized into a characteristic three dimensional structure or fold, that may exist in a variety of biological contexts.



**Repeats** Structural units typically repeated within a protein that assemble into a specific fold. Assemblies of repeats might also be thought of as domains.



**Sites & Motifs** Region of domains containing conserved active-site or binding residues or short conserved regions present outside domains that may adopt folded conformations only in association with their binding ligands



# Similarity identification with pairwise alignments

A popular way to identify similarities between proteins is to perform a pairwise alignment (Smith-Watermann, Needleman-Wunsch, BLAST, ...).

Normally, when the identity is higher than 40% this method gives good results.

```
>SCN2A_HUMAN_IQrepeat
EEVSAIIQRAYRRYLLKQKVKVSSIIYKK
      ↓ Blast
      ↓ Fasta

sp Q9UQD0 Sodium channel protein type 8 subunit alpha (Sodium channel protein 1980 AA
SCN8A_HUMAN type VIII subunit alpha) (Voltage-gated sodium channel align
subunit alpha Nav1.6) [SCN8A] [Homo sapiens (Human)]

Score = 36.3 bits (78), Expect = 0.025
Identities = 10/13 (76%), Positives = 12/13 (92%)

Query: 1 EEVSAIIQRAYR 13
          EEVSA++ QRAYR
Sbjct: 1895 EEVSAVVLQRAYR 1907
```

Only the N-ter of the query sequence matches and with a low score!



# Pairwise sequence alignments vs MSAs

Another weakness of the pairwise alignment is that no distinction is made between an amino acid at a crucial position (like an active site) and an amino acid with no critical role.

```

                *           20           *
SCN2A_HS : EEVSAIIIQRAYRRYLLKQKVKVSSSIYKK : 30
SCN8A_HS : EEVSAVVLQRAYRGHLARRGFICKKTTSNK : 30
          EEVSA666QRAYR   L 4           3   K

                *           20           *
SCN2A_HS : EEVSAIIIQRAYRRYLLKQKVKVSSSIYKK : 30
SCN2A_RN : EEVSAIVIQRAYRRYLLKQKVKVSSSIYKK : 30
SCN3A_HS : EEVSAAIQORNFRCYLLKQRLKNISSNYNK : 30
SCN8A_HS : EEVSAVVLQRAYRGHLARRGFICKKTTSNK : 30
SCN8A_MM : EEVSAVVLQRAYRGHLARRGFICRKITSNK : 30
IQGA1_HS_1 : NEGLITRLOARCRGYLVRQEFRRSMNFLKK : 30
IQGA1_HS_2 : QIPAITCIQSQWRGYKQKKAYQDRLAYLRS : 30
IQGA1_HS_3 : HKDEVVKIQSLARMHQARKRYRDRLOQYFRD : 30
IQGA1_HS_4 : HINDIIKIQAFIRANKARDDYKTLINAEDP : 30
IQGA1_MM_1 : NEGLITKLOACCRGYLVRQEFRRSMNFLKK : 30
          6Q   R   4

```

A multiple sequence alignment (MSA) gives a more general view of a conserved region by providing a better picture of the most conserved residues, which are usually essential for the protein function. It can help to identify subfamilies. An MSA contains more information than a pairwise alignment and several tools have been developed to extract this information.

# Extracting Information from MSAs

---

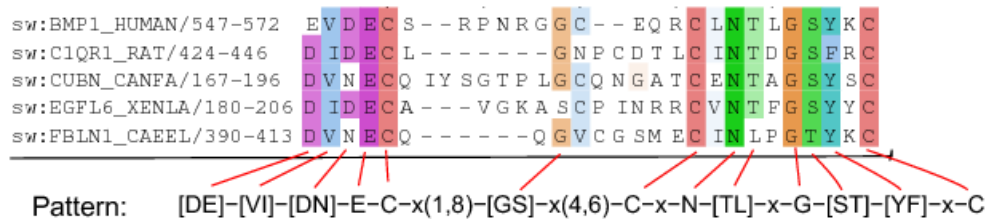
PROSITE patterns use a special syntax to describe the consensus of all the sequences present in the multiple alignment using a single expression.

Used to describe small functional regions:

- Enzyme catalytic sites;
- Prosthetic group attachment sites (heme, PLP, biotin, etc.);
- Amino acids involved in binding a metal ion;
- Cysteines involved in disulfide bonds;
- Regions involved in binding a molecule (ATP, calcium, DNA, etc.) or a protein.

Excellent tool to annotate active sites in combination with profiles (ProRules).

# How to build a PROSITE pattern

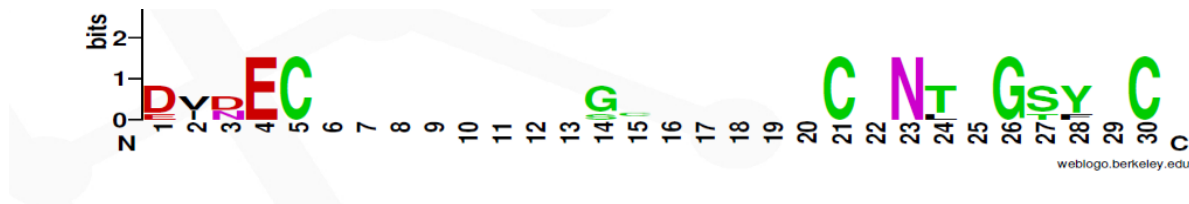


Collect sequences known to contain the signature and produce a multiple sequence alignment of the region of interest.

Build a pattern.

- By hand
- You can use automatic methods (e.g. <http://web.expasy.org/pratt/>) or a sequence logo to guide you

Example using a sequence logo (<http://weblogo.berkeley.edu/logo.cgi>):



# Tricks to build a PROSITE pattern

```
sw:BMP1_HUMAN/547-572  EVD E C S -- R P N R G G C -- E Q R C L N T L G S Y K C
sw:C1QR1_RAT/424-446   D I D E C L ----- G N P C D T L C I N T D G S F R C
sw:CUBN_CANFA/167-196 D V N E C Q I Y S G T P L G C Q N G A T C E N T A G S Y S C
sw:EGFL6_XENLA/180-206 D I D E C A --- V G K A S C P I N R R C V N T F G S Y Y C
sw:FBLN1_CAEEL/390-413 D V N E C Q ----- Q G V C G S M E C I N L P G T Y K C
```

Pattern: [DE]-[VI]-[DN]-E-C-x(1,8)-[GS]-x(4,6)-C-x-N-[TL]-x-G-[ST]-[YF]-x-C

- For the construction of the pattern, it is useful to consider residues and regions proved/thought to be important to the biological function of that group of proteins (e.g. enzyme catalytic sites, etc.).
- A first pattern is built from the MSA of the most conserved residues. It is used to scan the database.
- If it picks up too many false positives, it is modified to make it more stringent.
- The difficulty resides in achieving a pattern which does not pick up too many false positives yet does not miss too many sequences (false negatives).
- In some cases this result can not be achieved and an optimal sequence pattern can not be built.



# How to estimate the quality of a pattern

- We can not estimate the quality of a match with a pattern: PATTERNS don't produce a score, they match or not!
- But we can estimate the quality of the pattern.
- Two parameters can be computed to estimate the quality of a pattern: **precision** and **recall**.

**False positives** = known false hits.

**False negatives** = known missed hits.

**Precision** = true hits / (true hits + false positives).

"how useful the search results are"

Precision = 1  $\Rightarrow$  no false positive.

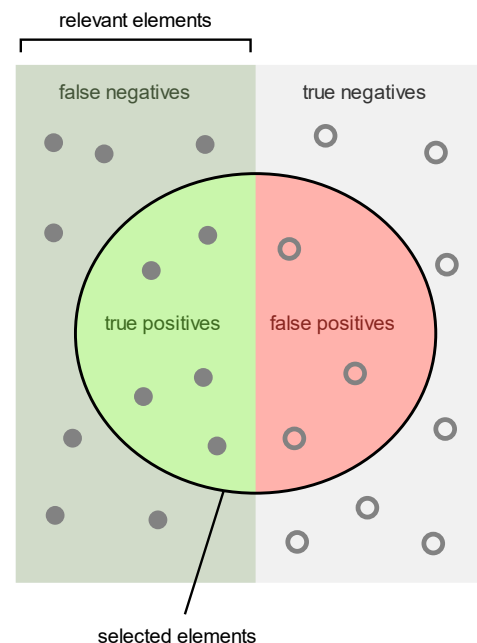
Precision = 0.8  $\Rightarrow$  20% false positives.

**Recall** = true hits / (true hits + false negatives).

"how complete the results are"

Recall = 1  $\Rightarrow$  No missed hits.

Recall = 0.8  $\Rightarrow$  20% missed hits.



- To obtain these measures we require a well annotated protein databases (PROSITE uses UniProtKB/Swiss-Prot).

# PROSITE patterns: example of a pattern entry

[Home](#) | [ScanProsite](#) | [ProRule](#) | [Documents](#) | [Downloads](#) | [Links](#) | [Funding](#)



Entry: **PS00972**

## General information about the entry

Entry name <a href="#">[info]</a>	USP_1
Accession <a href="#">[info]</a>	PS00972
Entry type <a href="#">[info]</a>	PATTERN
Date <a href="#">[info]</a>	JUN-1994 (CREATED); DEC-2013 (DATA UPDATE); APR-2015 (INFO UPDATE).
PROSITE Doc. <a href="#">[info]</a>	PDOC00750
Associated ProRule <a href="#">[info]</a>	PRU10092 <span style="color: red;">Active site</span>

## Name and characterization of the entry

Description <a href="#">[info]</a>	Ubiquitin specific protease (USP) domain signature 1.
Pattern <a href="#">[info]</a>	G-[LIVMFY]-x(1,3)-[AGCY]-[NASHOG]-x-C-[FYWC]-[LIVMFCA]-[NSTAD]-[SACV]-x-[LIVMSF]-[QF].

## Numerical results [\[info\]](#)

Numerical results for UniProtKB/Swiss-Prot release **2015\_06** which contains **548'586** sequence entries.

Total number of hits	282 in 282 different sequences	Number of true positives
Number of true positive hits	282 in 282 different sequences	
Number of 'unknown' hits	0	
Number of false positive hits	0	Number of false positives
Number of false negative sequences	28	Number of false negatives
Number of 'partial' sequences	1	
Precision (true positives / (true positives + false positives))	100.00 %	
Recall (true positives / (true positives + false negatives))	90.97 %	

## Comments [\[info\]](#)

Taxonomic range <a href="#">[info]</a>	Eukaryotes, Eukaryotic viruses
Maximum number of repetitions <a href="#">[info]</a>	1
Site <a href="#">[info]</a>	active_site at position 7
Version <a href="#">[info]</a>	1

# Limitations of PROSITE pattern

```
sw:BMP1_HUMAN/547-572  EVDECS--RPNRGGC--EQRCLNTLGSYKC
sw:ClQR1_RAT/424-446   DIDECL-----GNPCDTLCLINTDGSFRC
sw:CUBN_CANFA/167-196 DVNECQIYSGTPLGCGQNGATCENTAGSYSC
sw:EGFL6_XENLA/180-206 DIDECA---VGKASCPINRRCVNTFGSYYC
sw:FBLN1_CAEEL/390-413 DVNECQ-----QGVCGSMELINLPGTYKC
```

Pattern: [DE]-[VI]-[DN]-E-C-x(1,8)-[GS]-x(4,6)-C-x-N-[TL]-x-G-[ST]-[YF]-x-C

- OK to detect and annotate very conserved regions, but poor gap models
- residues at one position are considered equivalent in their frequencies
- if a symbol is not present at one position, this will exclude variants that have not yet been observed from being detected
- no score of the match is produced (you match or not)

# Profiles

```

A S T A M P V
A T S L M V T
S S S L M L T
A T P A M S S
A T A L L S A
A T A L L S A
  
```

Profiles turn a multiple sequence alignment into a position-specific scoring system suitable for searching databases for remotely homologous sequences

- Sequence weighting: correct sampling bias.
- Residue counts: get the frequency of each residue at each position of the MSA.
- Pseudo-counts: avoid frequencies of 0  $\Rightarrow$  avoid exclusion of residues.
- Build the final scoring matrix: used to build and score alignments.

w=0.2	A	S	T	A	M	P	V
w=0.2	A	T	S	L	M	V	T
w=0.2	S	S	S	L	M	L	T
w=0.1	A	T	P	A	M	S	S
w=0.1	A	T	A	L	L	S	A
	A	T	A	L	L	S	A

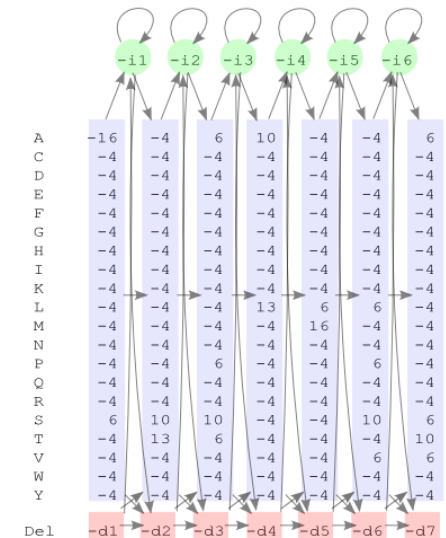
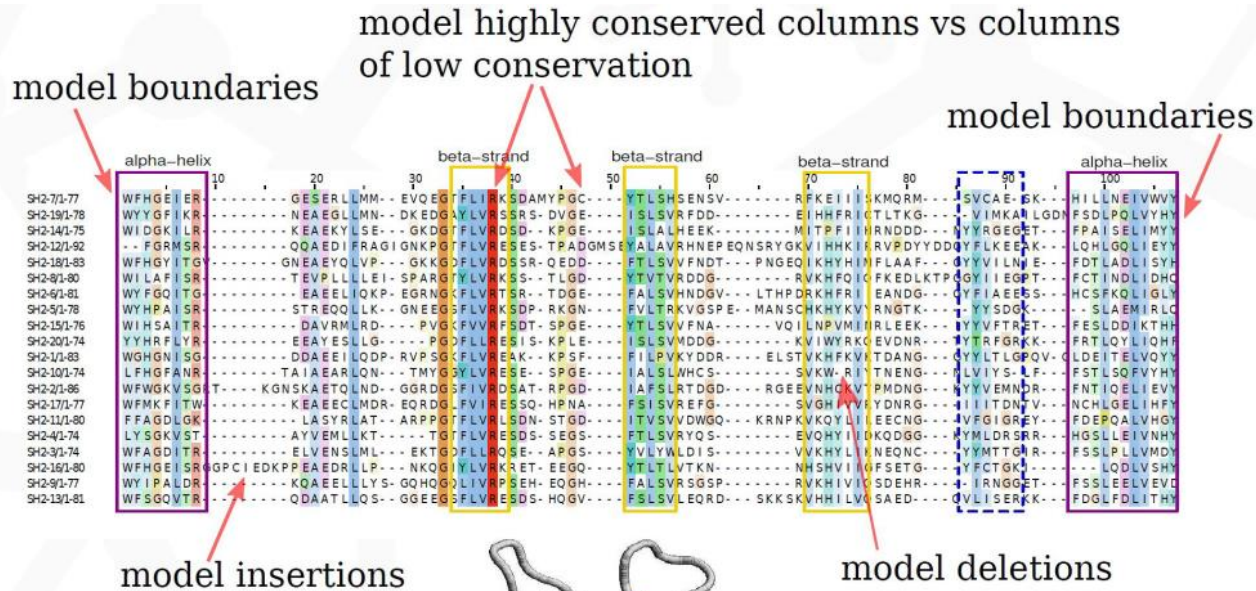
w=0.2	A	S	T	A	M	P	V	V
w=0.2	A	S	T	A	M	P	V	T
w=0.2	A	S	T	A	M	P	V	S
w=0.1	A	S	T	A	M	P	V	S
w=0.1	A	S	T	A	M	P	V	S

		1	2	3	4	5	6	7
A		0.8	0.0	0.2	0.4	0.0	0.0	0.2
C		0.0	0.0	0.0	0.0	0.0	0.0	0.0
D		0.0	0.0	0.0	0.0	0.0	0.0	0.0
E		0.0	0.0	0.0	0.0	0.0	0.0	0.0
F		0.0	0.0	0.0	0.0	0.0	0.0	0.0
G		0.0	0.0	0.0	0.0	0.0	0.0	0.0
H		0.0	0.0	0.0	0.0	0.0	0.0	0.0
I		0.0	0.0	0.0	0.0	0.0	0.0	0.0
K		0.0	0.0	0.0	0.0	0.0	0.0	0.0
L		0.0	0.0	0.0	0.6	0.2	0.2	0.0
M		0.0	0.0	0.0	0.0	0.8	0.0	0.0
N		0.0	0.0	0.0	0.0	0.0	0.0	0.0
P		0.0	0.0	0.2	0.0	0.0	0.2	0.0
Q		0.0	0.0	0.0	0.0	0.0	0.0	0.0
R		0.0	0.0	0.0	0.0	0.0	0.0	0.0
S		0.2	0.4	0.4	0.0	0.0	0.4	0.2
T		0.0	0.6	0.2	0.0	0.0	0.0	0.4
V		0.0	0.0	0.0	0.0	0.0	0.2	0.2
W		0.0	0.0	0.0	0.0	0.0	0.0	0.0
Y		0.0	0.0	0.0	0.0	0.0	0.0	0.0

		1	2	3	4	5	6	7
A		16	-4	6	10	-4	-4	6
C		-4	-4	-4	-4	-4	-4	-4
D		-4	-4	-4	-4	-4	-4	-4
E		-4	-4	-4	-4	-4	-4	-4
F		-4	-4	-4	-4	-4	-4	-4
G		-4	-4	-4	-4	-4	-4	-4
H		-4	-4	-4	-4	-4	-4	-4
I		-4	-4	-4	-4	-4	-4	-4
K		-4	-4	-4	-4	-4	-4	-4
L		-4	-4	-4	13	6	6	-4
M		-4	-4	-4	-4	16	-4	-4
N		-4	-4	-4	-4	-4	-4	-4
P		-4	-4	6	-4	-4	6	-4
Q		-4	-4	-4	-4	-4	-4	-4
R		-4	-4	-4	-4	-4	-4	-4
S		6	10	10	-4	-4	10	6
T		-4	13	6	-4	-4	-4	10
V		-4	-4	-4	-4	-4	6	6
W		-4	-4	-4	-4	-4	-4	-4
Y		-4	-4	-4	-4	-4	-4	-4

# Modeling of profiles



# Scoring: aligning a sequence to a profile

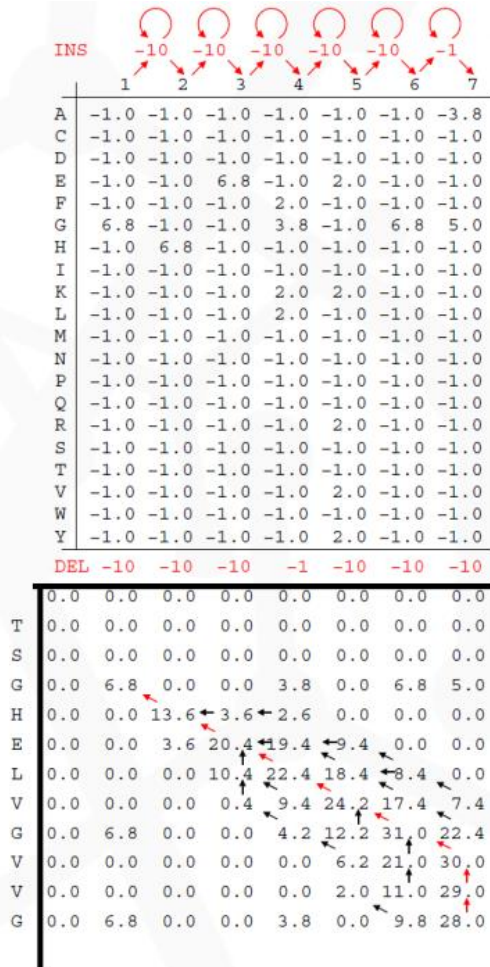
The similarity of new sequences to an existing profile can be tested by comparing each new sequence to the profile using a modification of the Smith/Waterman algorithm

The comparison of a sequence symbol to any row of the profile defines a specific value or "profile comparison value."

The best alignments of a sequence to a profile are found by aligning the symbols of the sequence to the profile in such a way that the sum of the profile comparison values minus the gap penalties is maximal

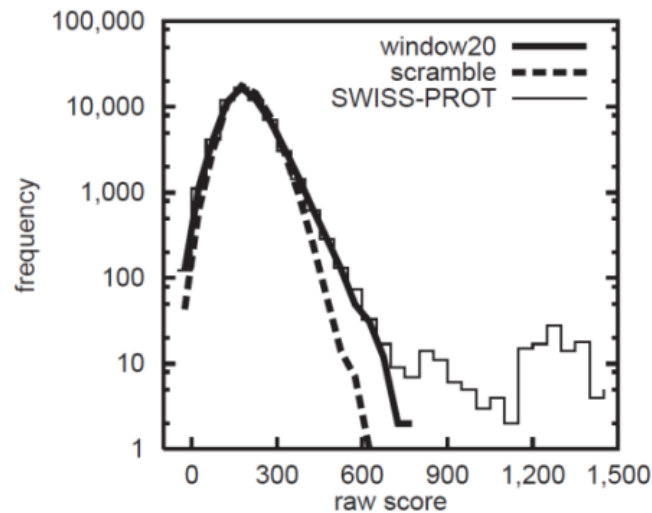
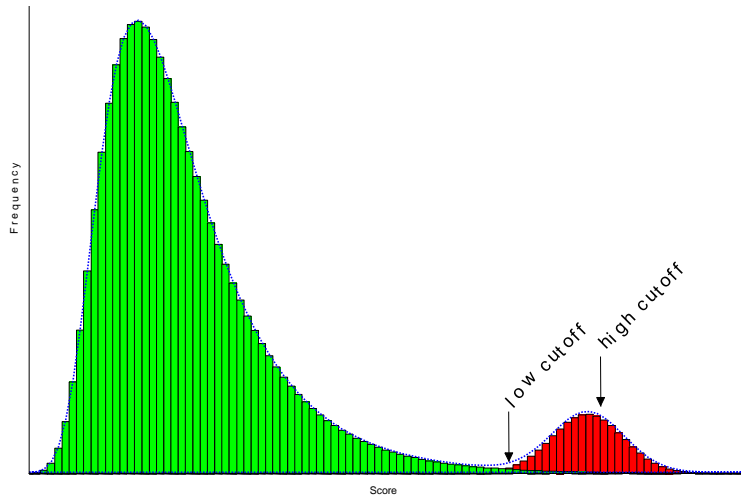
ALIGNMENT: score 28

	1	2	3	4	5	6	-	-	7
T S	G	H	E	L	V	G	V	V	G



# Profiles: interpretation of the score

- How do I interpret the score produced by a profile? Which is the lowest score I consider to produce a true match?
- Only biological arguments tell you if a match is true or not.
- However, a statistical analysis can help us decide if a match is statistically significant (true positive) or not (false positive).



- The score distribution of a profile on unrelated sequences is approximated by an Extreme Value Distribution (EVD) (green bars).
- This property permits to calculate the E-value: the number of matches that we expect to occur by chance with a score  $\geq$  a given cut off.

# Summary about patterns and profiles

---

## Patterns

- model a multiple sequence alignment using a compact string
- suited to model short and well conserved motifs
- good to describe functional residues
- easy to build, but not producing a score

## Profiles

- model a multiple sequence alignment using a numerical matrix representing the position-specific distribution of the residues
- suited to model protein domains and gapped motifs
- excellent technology to detect distant homologies
- matches produce a score that can be interpreted using statistical methods

Profiles and pattern can be used together (rules) to produce precise annotation

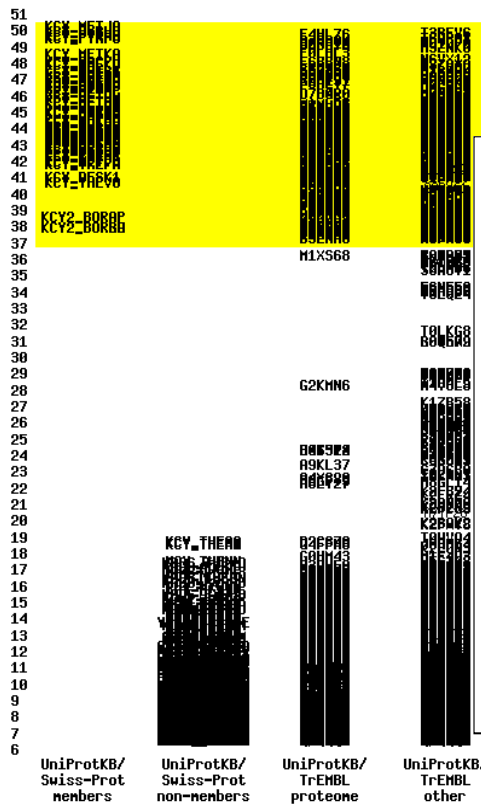




# Validation of a family profile



Score distribution of profile MF\_00239 matches in all kingdoms of UniProtKB



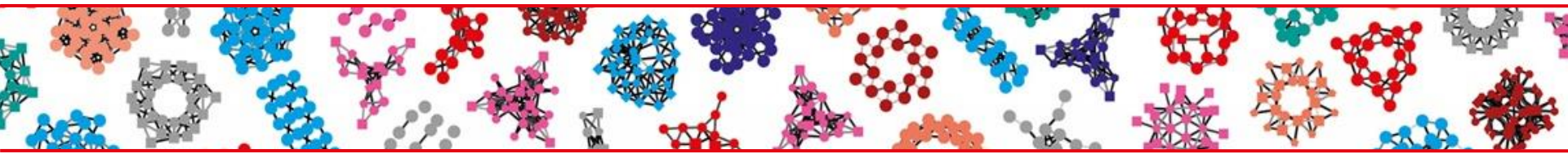
List of UniProtKB/Swiss-Prot true positive matches for the profile MF\_00239 (i.e. HAMAP MF\_00239 UniProtKB/Swiss-Prot members)

ac	id	kingdom	score	score_diff	description	organism
<a href="#">Q58071</a>	KCY_METIA	Archaea	50.213	+12.972	Cytidylate kinase	Methanocaldococcus jannaschii (strain ATCC 43067 / DSM 2661 / JAL-1 / JCM 10045 / NBRC 100440)

```

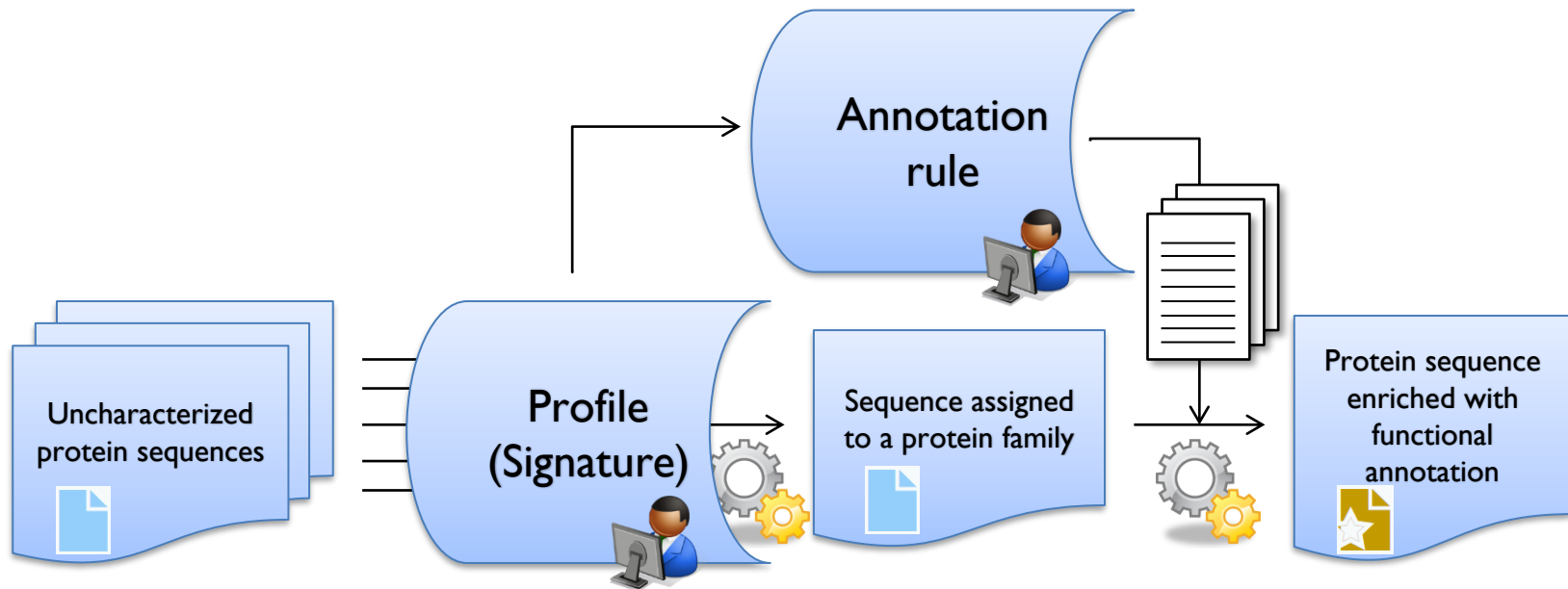
ID Cytidyl_kinase_type2; MATRIX.
AC MF_00239;
DT DEC-2013 (DATA UPDATE).
DE Cytidylate kinase [cmk].
CC /VERSION=4;
MA /GENERAL_SPEC: ALPHABET='ACDEFGHIKLMNPQRSTVWY'; LENGTH=191; LOG_BASE=1.071779; P0=0.9972;
MA P= 7.552363, 1.698108, 5.303439, 6.320015, 4.078187, 6.844419, 2.240667,
MA /DISJOINT: DEFINITION=PROTECT; N_SCORE=37.241;
MA /NORMALIZATION: MODE=1; FUNCTION=LINEAR; R1=17.3045; R2=0.0080365; TEXT='-LogE';
MA /CUT_OFF: LEVEL=1; SCORE=25.4; N_SCORE=37.241; MODE=1; TEXT='!';
MA /CUT_OFF: LEVEL=0; SCORE=25.4; N_SCORE=37.241; MODE=1; TEXT='?';
MA /CUT_OFF: LEVEL=-1; SCORE=-100; N_SCORE=8.5; MODE=1; TEXT='??';
MA /DEFAULT: B0=*; B1=*; E0=*; E1=*; MM=; II=*;
MA /I: B0=-85; B1=-85; BD=-49;
MA /M: SY='M'; M=-25,-24,-49,-42,8,-41,-30,-3,-10,-4,46,-37,-41,-35,-8,-32,-6,6,-28,-25; M0=-
MA /I: MM=0; MI=-100; MD=-111; IM=-9; II=-11; DM=-7; DD=-14;
MA /M: SY='R'; M=-24,-25,-39,-33,-2,-39,-27,20,0,-1,-17,-3,-39,-4,21,3,-24,11,-29,-25; M0=-12;
MA /I: MM=0; MI=-100; MD=-111; IM=-9; II=-11; DM=-7; DD=-14;
MA /M: SY='I'; M=-44,-39,-71,-68,-43,-70,-70,39,-68,-13,-30,-66,-66,-66,-69,-64,-44,4,-62,-57;
MA /I: MM=0; MI=-100; MD=-111; IM=-9; II=-11; DM=-7; DD=-14;
MA /M: SY='T'; M=10,11,-59,-53,-34,-51,-43,3,-50,-29,5,-49,-52,-47,-49,-43,31,15,-41,-37; M0=-
MA /I: MM=0; MI=-101; MD=-111; IM=-9; II=-11; DM=-7; DD=-14;
MA /M: SY='I'; M=-44,-39,-71,-68,-44,-70,-71,35,-68,-4,-31,-67,-67,-67,-70,-64,-44,20,-63,-58;
M0=-44;
    
```

<a href="#">Q5IMV8</a>	KCY_NAIPD	Archaea	47.159	+9.918	kinase	2160)
<a href="#">C6A187</a>	KCY_THESM	Archaea	47.159	+9.918	Cytidylate kinase	Thermococcus sibiricus (strain MM 739 / DSM 12597)



- Why do we need predictive annotation tools?
- Protein signatures for homology detection – A short primer
- **Annotation rules for functional annotation**
- HAMAP and PROSITE - automatic annotation in UniProtKB
- HAMAP and PROSITE - services for external users
- Practical exercises (afternoon)

# Automatic annotation of protein sequences



- Manually curated **signatures** – a protein matching a specific signature may be annotated according to the contents of the associated rule
- **Manually** created annotation **rules** that specify **annotations** AND the **conditions** under which they may be applied
- **HAMAP** signatures/rules for **full length protein sequences**, **PROSITE** signatures/rules mainly for **domains** and **sites**
- Common format and shared syntax

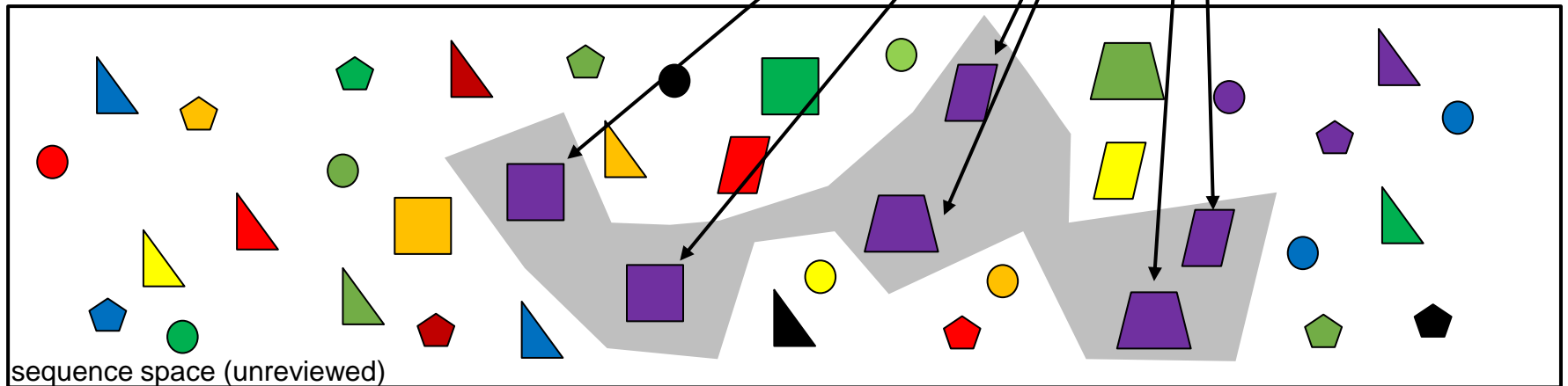
# A fictitious rule

If a protein meets these conditions...

- color: purple
- has four edges: yes

... then these annotations are applied

- is a purple quadrilateral



# A fictitious rule

If a protein meets these conditions...

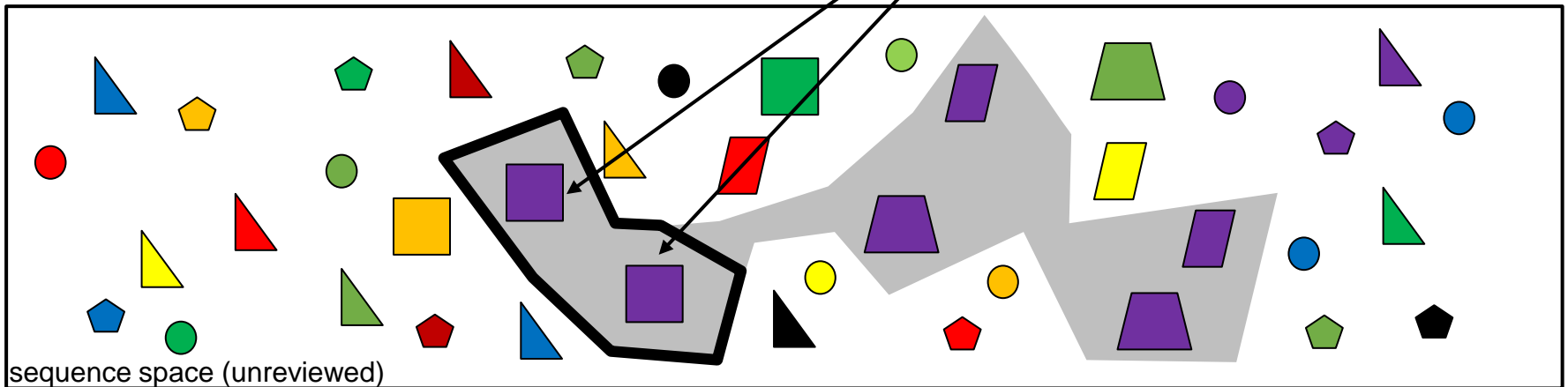
- color: purple
  - has four edges
- and:
- all sides of same length
  - all angles are 90 degrees

... then these annotations are applied

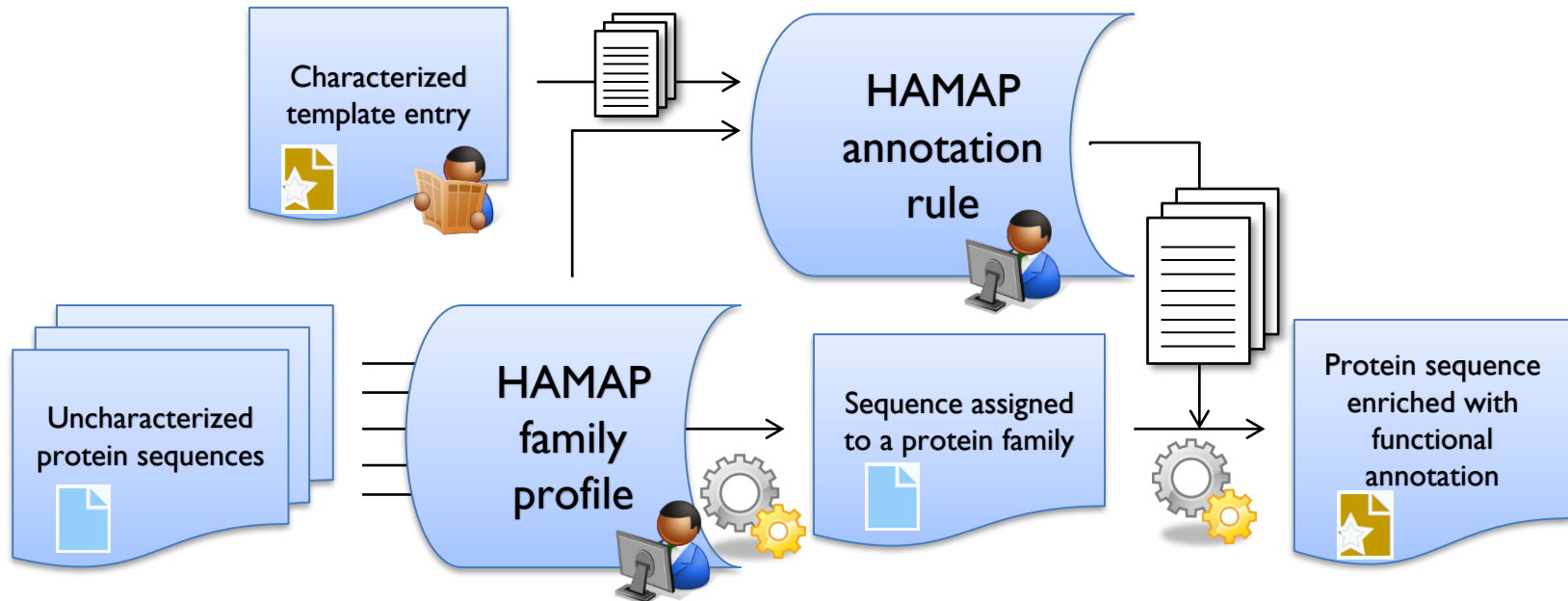
- is a purple quadrilateral

**CONDITIONS**

**ANNOTATIONS**



# HAMAP annotation rules



- Rule creation starts with the manual curation of UniProtKB/Swiss-Prot records (=template entries).
- Protein names, gene names, functional annotation, GO terms, keywords and sequence features from the templates are selected and combined to build the rule containing the annotation to be propagated.
- Annotations may be subject to control statements that limit their propagation to only those sequences satisfying one or more conditions.



### Annotation rule MF\_00041

[Send feedback](#)

#### General rule information

Accession	MF_00041
Date	1-JUN-2001 (Created) 20-NOV-2019 (Last updated, Version 1.0)
Profile	Cys_tRNA_synth
Taxa	Bacteria Archaea
Template	P21888 (SYC_ECOLI)
Triggered by	HAMAP; MF_00041 ( <a href="#">Get profile general information and statistics</a> )

Annotation applies only to sequences from specified taxa

Annotation applies to sequences matching specified profile

#### Propagated annotation [?]

##### Identifier, protein and gene names [?]

Identifier	MF_00041
RecName:	Full=Cysteine--tRNA ligase; EC 6.1.1.16;
AltName:	Full=CysteinyI-tRNA synthetase; Short=CysRS;
Gene name	cysS

A meaningful recommended name which can be safely propagated to orthologs; synonyms

**CONDITIONS**

**ANNOTATIONS**





## Annotation rule MF\_00041

[Send feedback](#)

### Propagated annotation [?]

#### Comments [?]

Catalytic activity	<a href="#">RHEA:17773</a> : ATP + L-cysteine + tRNA(Cys) = AMP + diphosphate + L-cysteinyl-tRNA(Cys) <a href="#">EC 6.1.1.16</a>
case <FTGroup:1> Cofactor	<a href="#">Zn(2+)</a> <u>Note</u> : Binds 1 zinc ion per subunit.
end case case <OC:Bacteria> Subunit	Monomer.
end case Subcellular location	Cytoplasm.
Similarity	Belongs to the class-I aminoacyl-tRNA synthetase family.



Protein function, interactions, location, pathways, family membership and others – Swiss-Prot syntax

#### Keywords [?]

[Cytoplasm](#)  
[Aminoacyl-tRNA synthetase](#)  
[Protein biosynthesis](#)  
[Ligase](#)  
[ATP-binding](#)  
[Nucleotide-binding](#)  
case <FTGroup:1>  
[Metal-binding](#)  
[Zinc](#)  
end case  
case <FT:8>  
[Phosphoprotein](#)  
end case



Summary in the form of keywords

Case statements restrict annotation to proteins satisfying these conditions



## Annotation rule MF\_01962

[Send feedback](#)

Characterized template

Reference coordinates

Feature descriptions

Conditions specify when a feature should be defined as present: PROSITE pattern syntax

A group of related features

### Features [?]

From: ADE\_PSEAE (Q916Y4)

Key	From	To	Description	Tag	Condition	FTGroup
ACT_SITE	197	197	Proton donor		E	
METAL	14	14	Zinc; catalytic		H	1
METAL	16	16	Zinc; catalytic		H	1
METAL	194	194	Zinc; catalytic		H	1
METAL	275	275	Zinc; catalytic		D	1
BINDING	276	276	Substrate		D	
SITE	218	218	Important for catalytic activity		H	

# Annotation transfer

- Sequence features are propagated by aligning target sequences to the profile and translating feature positions from the template to the target sequence.

```
FT From: ADE_PSEAE (Q9I6Y4)
FT METAL 14
FT /note="Zinc; catalytic"
FT Group: 1; Condition: H
FT METAL 16
FT /note="Zinc; catalytic"
FT Group: 1; Condition: H
```

```
ID C4ZN08_THASP Unreviewed; 343 AA.
AC C4ZN08;
```

```
FT METAL 16
FT /note="Zinc; catalytic"
FT /evidence="ECO:0000256|HAMAP-Rule:MF_01962"
FT METAL 18
FT /note="Zinc; catalytic"
FT /evidence="ECO:0000256|HAMAP-Rule:MF_01962"
```

CLUSTAL W (1.83) multiple sequence alignment template=ADE\_PSEAE profile\_method=hmmbuild

```
C4ZN08_THASP -----MELEAYVRALPKAEIHHIEGTLPEPEMMFALARRNGVALPWA-----SV
ADE_YEAST -----MVSVEFLQELPKCEHHHLEGTLPEPDLFFPLAKRNDIILPE----GFPKSV
ADE_CANGA -----MVPESFLELPLKCEHHHLEGTLPEPDLFFPLAKRNNIQLPD----HFPTQ
ADE_KLULA ---MAKFECTDEVTFNLTLPKCEHHHLEGTLPEPELLFQLVERNGVQLPG----TFPKTV
ADE_CANAL ---MAQYECSEHMENFLRELPKCEHHHLEGTLPEPELLFQLVERNGVQLPG----TFPKTV
ADE_ASPFU MC-----QSPLDHDFLHGLPKCEHHHLEGCVTPELIFQLAEKNNIQLPNPATHPAYASV
ADE_ASPOR MC-----KSDLHDFLHGLPKCEHHHLEGCLAPDLIFELAKRNNVSLPN---EPAYESI
ADE_EMENI MCPPNTPYQSQWHAFHLHSLPKCEHHHLEGCLEPPLIFSMARKNNVSLPSPSSNPAYTSV
ADE_SCHPO MS-----NLPIYNFIRKLPKCEHHHLEGCLSPDLVFRLAKKNGITLPS--DDAAYTTP
ADE_GIBZE MC-----KSRVHSFLQALPKVECHHIEGTLPEPELLFTLAEKNGIELPN---DPVYESA
ADE_CAUCR -MTDASFAPSASAEFVRGLPKAEIHHIEGSLEPELMFELAQRNGITLPPA-----SV
ADE_CAUCN -MTDASFAPSASAEFVRGLPKAEIHHIEGSLEPELMFELAQRNGITLPPA-----SV
ADE_SPHAL -MPDGFASHEERAAFIAGLPKAEIHHIEGSLEPELLFEFARRNVAIPFA-----SI
ADE_RHORT -----MAVDPAFLHALPKVEIHHIEGSLEPEMMVALAERNGLRRLPYA-----SV
ADE_STRCO -----MKRPYDALMPLPKAEIHHIEGTLPELAFALAARNGVSLPYA-----DE
ADE_BURPP -MTTTTPTVPTPLAEKTALAPKAEIHHIEGSLEPELIFALAERNVVKLAYD-----SI
ADE_BURXL -MTTTTPTVSTPLAEKTVLAPKAEIHHIEGSLEPELIFALAERNVVKLAYD-----SI
ADE_CUPTR -----MTIDAALAEQIRRTPKAEIHHIEGTLPELIFRLAQRNQVALPYP-----SV
ADE_CUPNH -----MTIDAALAEQIRRTPKAEIHHIEGTLPELIFRLAQRNQVALPYP-----SV
ADE_CUPPJ -----MTIDAALADKIRRTPKAEIHHIEGTLPELIFRLAQRNHVNLVYP-----SV
ADE_RALME -----MTIDAALADKIRRTPKAEIHHIEGTLPEPERIFRLAQRNNVKLAYP-----DV
ADE_RALPJ -----MPISSALAERIATSPKAEIHHIEGSLEPELMFALAERNVVKLAYP-----SV
ADE_RALSO -----MPISSALAERIATSPKAEIHHIEGSLEPELMFALAERNVVKLAYP-----SV
ADE_GEOLS -MNLTNIPRQALPELLCRMPKAEIHHIEGSLEPELIFALAERNRLQLAYP-----TI
ADE_GEOUR -MNFDCIPREDLHGILCHMPKAEIHHIEGSLEPELIFELATRNRQLPYP-----TI
ADE_RHOFD -MTIKPVSQERLPELLRTIPKAEIHHIEGSLEPELMFALAQRNGVSIYP-----DV
ADE_ZYMMO -----MNNLIKFTAALPKAEIHHIEGSLEPELMFELAKRNVTLFPF-----DV
ADE_ACIAA -----MNQSELIRALPKAEIHHIEGTFPELMFEIQRNHIDIPYK-----SV
ADE_PSYCK -----MIDLIKRLPKAEIHHIEGSLEPELMFRLAKKNQIEIPYK-----DI
ADE_VIBPA -----MNAFVQGLPKVEIHHIEGSLEPELMFKLAKRNGIDIPYS-----SP
ADE_PSEAE -----MYEWNALPKAEIHHIEGTLPEPELLFALAERNRIALPWN-----DV
```



Search HAMAP

Search

## Annotation rule MF\_01962

[Send feedback](#)

### Propagated annotation [?]

case <FTGroup:1>

Cofactor

[Zn\(2+\)](#)

Note: Binds 1 zinc ion per subunit.

end case

Similarity

Belongs to the metallo-dependent hydrolases superfamily. Adenosine and AMP deaminases family. Adenine deaminase type 2 subfamily.

#### Keywords [?]

[Hydrolase](#)

[Nucleotide metabolism](#)

case <FTGroup:1>

[Metal-binding](#)

[Zinc](#)

end case

#### Gene Ontology [?]

[GO:0000034](#); Molecular function: adenine deaminase activity.

[GO:0006146](#); Biological process: adenine catabolic process.

[GO:0043103](#); Biological process: hypoxanthine salvage.

case <FTGroup:1>

[GO:0008270](#); Molecular function: zinc ion binding.

end case

Annotation applied if, and only if, ALL members of the specified group (here, FTGroup I, are present).

# Annotation output

- Many checks are performed in order to prevent the propagation of wrong annotation:

```
DR   HAMAP; MF_01962; Adenine_deaminase; 1.
KW   Hydrolase; Metal-binding; Nucleotide metabolism; Zinc.
FT   CHAIN           1..343
FT                               /note="Adenine deaminase"
FT   ACT_SITE       199
FT                               /note="Proton donor"
FT                               /evidence="ECO:0000255|HAMAP-Rule:MF_01962"
FT   METAL          16
FT                               /note="Zinc; catalytic"
FT                               /evidence="ECO:0000255|HAMAP-Rule:MF_01962"
FT   METAL          18
FT                               /note="Zinc; catalytic"
FT                               /evidence="ECO:0000255|HAMAP-Rule:MF_01962"
FT   METAL          196
FT                               /note="Zinc; catalytic"
FT                               /evidence="ECO:0000255|HAMAP-Rule:MF_01962"
FT   METAL          277
FT                               /note="Zinc; catalytic"
FT                               /evidence="ECO:0000255|HAMAP-Rule:MF_01962"
FT   BINDING        278
FT                               /note="Substrate"
FT                               /evidence="ECO:0000255|HAMAP-Rule:MF_01962"
FT   SITE           220
FT                               /note="Important for catalytic activity"
FT                               /evidence="ECO:0000255|HAMAP-Rule:MF_01962"
**
** ##### INTERNAL SECTION #####
**EV ECO:0000255; HAMAP-Rule:MF_01962; XXX; 02-MAR-2020.
**HA FAM; Method MF_01962; ADE; Trusted match; 85.781 (+32).
**HA FAM; Method MF_00540; ADD; Weak match; 28.113 (-21.7).
**HA SAM; Annotated by HAMAP 3.68.17; MF_01962.9; MF_01962; 02-MAR-2020 12:42:26.
SQ   SEQUENCE   343 AA;  37571 MW;  E2DA4A77193E5C4B CRC64;
      MELEAYVRAL PKAELHLHIE GTLEPEMMFA LARRNGVALP WASVEAVRAA YAFTDLQSFL
      DLYYAGAAVL VREQDFELA FAYFERAHAD GVVHAELFFD PQTHTARGVA LETVLDGLER
      ACVEARARWG IGSRLILCFL RHLSEEEGFA TLQQALPHLS RIDGVGLDSS ERGHPPAKFA
      RLFARCRELG LHVVAHAGEE GPPAYIVDAL DLLKAERIDH GVRCTEDPAL VGRLVREQVP
      LTVCP LSNVK LCVFPDLARH NLGQLFAAGL KVTINSDDPA YFGGYVAKNY VDTARALGLG
      RAELRRIARN SLEASFVSAA ERAPWLARLD ALGEDCEGEG GAA
```

//

# Annotation output validation

```
ID F1VZU6_9BURK Unreviewed; 343 AA.
AC F1VZU6;
..
DR HAMAP; MF_01962; Adenine_deaminase; 1.
KW Hydrolase; Nucleotide metabolism.
FT CHAIN 1..302
FT /note="Adenine deaminase"
FT ACT_SITE 163
FT /note="Proton donor"
FT /evidence="ECO:0000255|HAMAP-Rule:MF_01962"
FT BINDING 242
FT /note="Substrate"
FT /evidence="ECO:0000255|HAMAP-Rule:MF_01962"
FT SITE 184
FT /note="Important for catalytic activity"
FT /evidence="ECO:0000255|HAMAP-Rule:MF_01962"
**
** ##### INTERNAL SECTION #####
**EV ECO:0000255; HAMAP-Rule:MF_01962; XXX; 02-MAR-2020.
**HA SAM; Annotated by HAMAP 3.68.17; MF_01962.9; MF_01962; 02-MAR-2020 11:28:48.
**HW HAMAP: WARNING(0) for [MF_01962]: Conditions for mandatory group [1] not met (in one or more instances of the
```

```
**HW HAMAP: WARNING(0) for [MF_01962]: Conditions for mandatory group [1]
not met (in one or more instances of the trigger) (group condition: ...
```

```
catalytic|__FT_5,__FTGRP_1,__HAS_PATTERN_COND,D)))
```

```
**HW HAMAP: WARNING(0) for [MF_01962]: Protein size (302) smaller than
specified minimum (315)
```

```
MALAFESVDA LRAAYAFIDL QSFLLDIYAG ASVLLTEQDF YDMTRAYLLR AQQDQVRHTE
LEFDDQSTLVA PQLVDFEVEN STVALQDAG QSLVHTQALTE LQSLRHLGSD PAFSTLFLAL
```

```
**HW HAMAP: WARNING(0) for [MF_01962]: Trigger [MF_01962] (with MaxNterGaps:
17): N-ter too short by 29 AA
```

```
PA
```

```
//
```

# Annotation output validation

CLUSTAL W (1.83) multiple sequence alignment template=ADE\_PSEAE

```
ADE_YEAST      -----MVSVEFLQELPKCEHHLHLEGTLEPDLLFPLAKRNDIILPE----GFPKSV
ADE_CANGA      -----MVPESFLLLELPKCEHHLHLEGTLEPDLLFPLAKRNNIQLPD----HFPQTP
ADE_KLULA      --MAKFECTDEVTNFLTLPKCEHHLHLEGTLEPELLFQLVERNGVQLPG----TFPKTV
ADE_CANAL      --MAQYECSEHMENFLRELPKCEHVVHLEGTLEPSLLFKLAKRNNITLPE----TFPKTV
ADE_ASPFU      MC-----QSPLHDFLHGLPKCEHVVHLEGCVTPELIFQLAEKNNIQLPNPATHPAYASV
ADE_ASPOR      MC-----KSDLHDFLHGLPKCEHVVHLEGCCLAPDLIFELAKRNNVSLPN---EPAYESI
ADE_EMENI      MCPPNTPYQSQWHAFLHSLPKCEHVVHLEGCLEPPLIFSMARKNNVSLPSPSSNPAYTSV
ADE_SCHPO      MS-----NLPIYNFIRKLPKCEHVVHLEGCCLSPDLVFRLLAKKNGITLPS--DDAAYTTP
ADE_GIBZE      MC-----KSRVHSFLQALPKVEQHLHIEGTLEPELLFTLAEKNGIELPN---DPVYESA
ADE_CUPTR      -----MTIDAALAEQIRRTPKAELHVVHIEGTLEPELIFRLAQRNQVALPYP-----SV
ADE_CUPNH      -----MTIDAALAEQIRRTPKAELHVVHIEGTLEPELIFRLAQRNQVALPYP-----SV
ADE_CUPPJ      -----MTIDAALADKIRRTPKAELHVVHIEGTLEPELIFRLAQRNHVNL PYP-----SV
ADE_RALME      -----MTIDAALADKIRRTPKAELHVVHIEGTLEPERIFRLAQRNNVKLAYP-----DV
ADE_RALPJ      -----MPISSALAERIATSPKAEHLHIHIEGSLEPELMFALAERNGVKLPYA-----SV
ADE_RALSO      -----MPISPALAERIATSPKAEHLHIHIEGSLEPELMFALAERNGVKLPYA-----SV
F1VZU6_9BURK  -----MALAFE-----SV
ADE_BURPP      -MTTIVTPTPLAECTALAPKAEHLHIHIEGSLEPELIFALAERNGVKLAYD-----SI
ADE_BURXL      -MTTIVTSTPLAECTVLAPKAEHLHIHIEGSLEPELIFALAERNGVKLAYD-----SI
ADE_GEOLS      -MNLTNIPRQALPELLCRMPKAEHLHIHIEGSLEPELIFALAERNRLQLAYP-----TI
ADE_GEOUR      -MNFDCIPREDLHGILCHMPKAEHLHIHIEGSLEPELIFELATRNNRIQLPYP-----TI
ADE_RHOFD      -MTIKVVSQERLELLRTIPKAEHLHIHIEGSLEPELMFALAQRNGVSI PYP-----DV
ADE_ZYMMO      -----MNNLIKFIAALPKAEHLHLHIEGSLEPELMFELAKRNVKTL PFP-----DV
ADE_CAUCR      -MTDASFAPSASAEFVRGLPKAEHLHMHIEGSLEPELMFELAQRNGITLPFA-----SV
ADE_CAUCN      -MTDASFAPSASAEFVRGLPKAEHLHMHIEGSLEPELMFELAQRNGITLPFA-----SV
```

# Annotation output validation

```
ID G9PG85_9ACTO          Unreviewed;      332 AA.
AC G9PG85;
..
DR HAMAP; MF_01962; Adenine_deaminase; 1.
KW Hydrolase; Nucleotide metabolism.
FT CHAIN                1..332
FT                      /note="Adenine deaminase"
FT BINDING              278
FT                      /note="Substrate"
FT                      /evidence="ECO:0000255|HAMAP-Rule:MF_01962"
FT SITE                 222
FT                      /note="Important for catalytic activity"
FT                      /evidence="ECO:0000255|HAMAP-Rule:MF_01962"
**
** ##### INTERNAL SECTION #####
**EV ECO:0000255; HAMAP-Rule:MF_01962; XXX; 02-MAR-2020.
**HA SAM; Annotated by HAMAP 3.68.17; MF_01962.9; MF_01962; 02-MAR-2020 11:36:03.
**HW HAMAP: WARNING(0) for [MF_01962]: Conditions for mandatory group [1] not met (in one or more instances of the
```

```
**HW HAMAP: WARNING(0) for [MF_01962]: Mandatory Positional (pattern [E])
not matched (1 or more times for
Positional|ACTIVE_SITE|{197//}|{197//}|Proton
donor|__FT_1,__HAS_PATTERN_COND)
```

```
SQ SEQUENCE 332 AA; 37491 MW; 701E1C335EA545DB CRC64;
MAHEALVQAL PKVELHVHIE GTLEPELKFK LAQRNGIKLP FSSEQEVKDS YTFNDLASFL
DAYYDGMNVL LASEDFYDLA MAYRKAASQ GLRYAEIFFD PQAHTCRGVS FHTVISGLRR
AQLAEQTLG VFSQFIMCFL RDFQPEYAMA TLLESLPYKQ WIVGVGLDSD ETDHEPAKFA
AVFARARREG YQLTMHCDVD IKDSIEHIRQ VIEDIQVQRV DHGTNVVENP ALVKYLVEHR
IGLTSCPISN LWISESDKVE LVKQLVADGV LVTINSDDPA YFGGYIGDNF QRVADHDGVD
EDFLRALVAN AIEISWAPLA VKQRLREELA AV
```

```
//
```



# Annotation output validation

```
ADE_RALME      ---RDRFVGVGLDSSERGNPPEKFARVFARAKE--LGLHLVAHAGEE--GPPQYVTDALD
ADE_RALPJ      --PANRIIGVGLDSSERGNPPEKFARVFARCKE--LGLRLVAHAGEE--GPAQYVIDALD
ADE_RALSO      --PANRIIGVGLDSSERGNPPEKFARVFARCKA--LGLRLVAHAGEE--GPAQYVIDALD
ADE_GEOLS      ---RDKFIGVGLDSGEKGNPPEKFSRVFARCRE--LGLRLVAHAGEE--GTAEYIWHALD
ADE_GEOUR      ---RDKFIGVGLDSSERGNPPEKFTRVFARCRE--LGLRLVAHAGEE--GSAEYISHSLD
ADE_RHOFD      ---LDKLVGVGLASSEMGHPPEKFARVFARARE--LGLRLVAHAGEE--GPPAYIWSALD
ADE_ZYMMO      ---LDKIAGVGLDSSEVGNPPSKFRHVFAEARQ--KGLKLVAHAGEE--GDASYIKEALD
ADE_MARMS      ---LKWIDGIGLDSSEVGHPPPEKFLRVFEACKN--LGLKVTAHAGEE--GPPDYVWQAIE
ADE_RHOPB      ---KDQIIAIGMGAELGNPPAKFARFFKAARD--RGFRTTVHAGEE--GPAAYVREALE
ADE_STRCO      ---LDRITGVGLDSAIEVGHPPVKFREVEAAAA--LGLRRVAHAGEE--GPPAYVVEALD
G9PG85_9ACTO   ---KQWIVGVGLDSDETDHEPAKFAAVFARARR--EGYQLTMHCDVDIKDSIEHIRQVIE
ADE_RHIEC      ---NPLITGFNLAGEERMGRVADYIRAFDIARD--AGLGLTIHAGEV--CGAFSVADALD
ADE_RHIE6      ---NPLITGFNLAGEERMGRVADYSRAFDIARD--AGLGLTIHAGEV--CGAFSVADALD
ADE_RHILW      ---NPLITGFNLAGEERMGRVADYARAFDIARD--AGLGLTIHAGEV--CGAFSVADALD
ADE_RHIL3      ---NPLITGFNLAGEERMGRVADYARAFDIARD--AGLGLTIHAGEV--CGAFSVADALD
ADE_AGRRK      ---NPLISGFNMAGEERMGRVADYARAFDIARE--AGLGITIHAGEV--CGAFSVADAVE
ADE_AGR5T      ---HPLVTGFNMAGEERMGRVADYARAFDIARD--AGLGLTIHAGEV--CGPESVADALD
ADE_RHIME      ---HPLVTGFNLAGEERMHSVAEFSRAFDIVRD--AGLGLTIHAGEE--SGAFSVRDALD
ADE_RHISN      ---HPLVTGFNLAGEERMHSVAEFARAFDIVRD--AGLGLTIHAGEE--SGAFSVRDALD
ADE_SINMW      ---HPLVTGFNLAGEERMHSVAEFSRAFDIVRD--AGLGLTIHAGEE--SGAFSVRDALD
ADE_AGRVS      ---HPLITGFNMAGEERMNRVADYAPAFDIARE--AGYGITIHAGEE--CGAFSVRDALD
ADE_MESSB      ---HPLVTGFVGMAGDERAGHPRDFAYAFDIARE--AGLGISIHAGEE--GGAESVEAALD
ADE_RHILO      ---NPLVTGFGVAGDERVGEEMEDYVRAFEIARE--AGLGITIHAGEE--TGWETVQAALD
```

```
. . . . * : *.. :
```

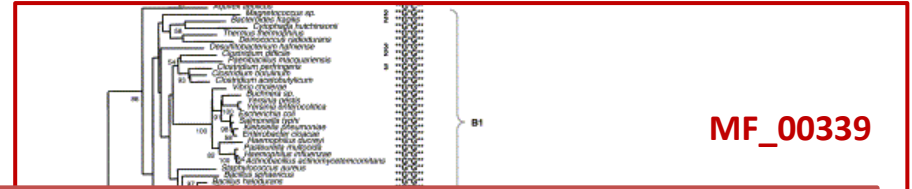
# Controlling annotations using cases and conditions



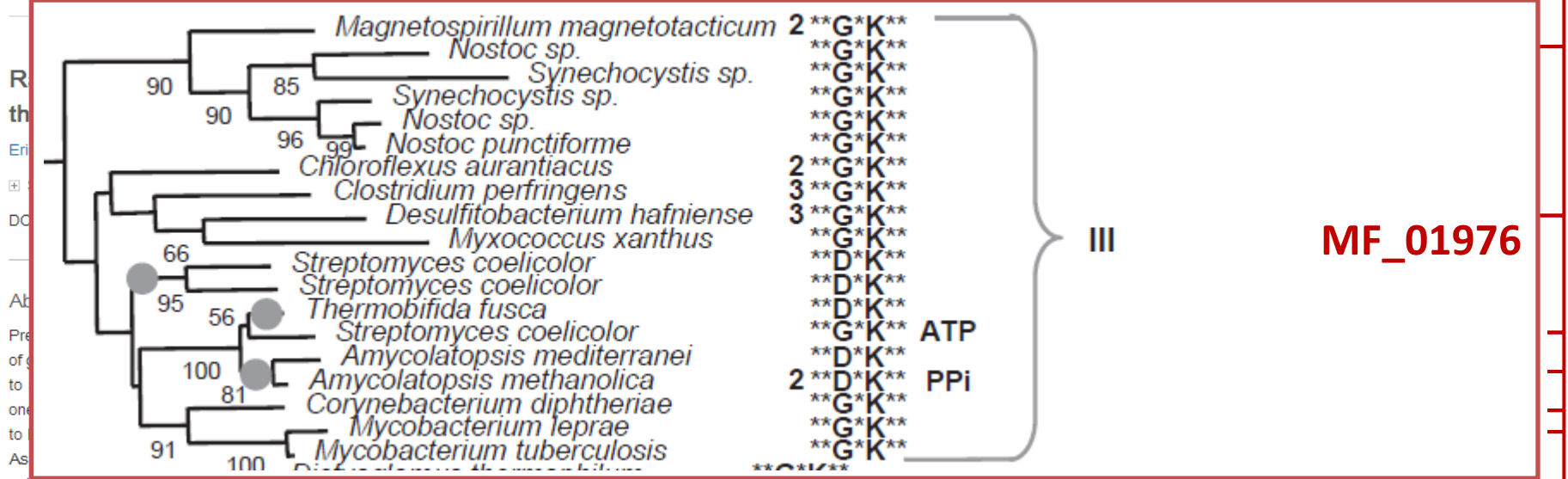
Gene

Volume 318, 30 October 2003, Pages 185–191

GENE



MF\_00339



MF\_01976

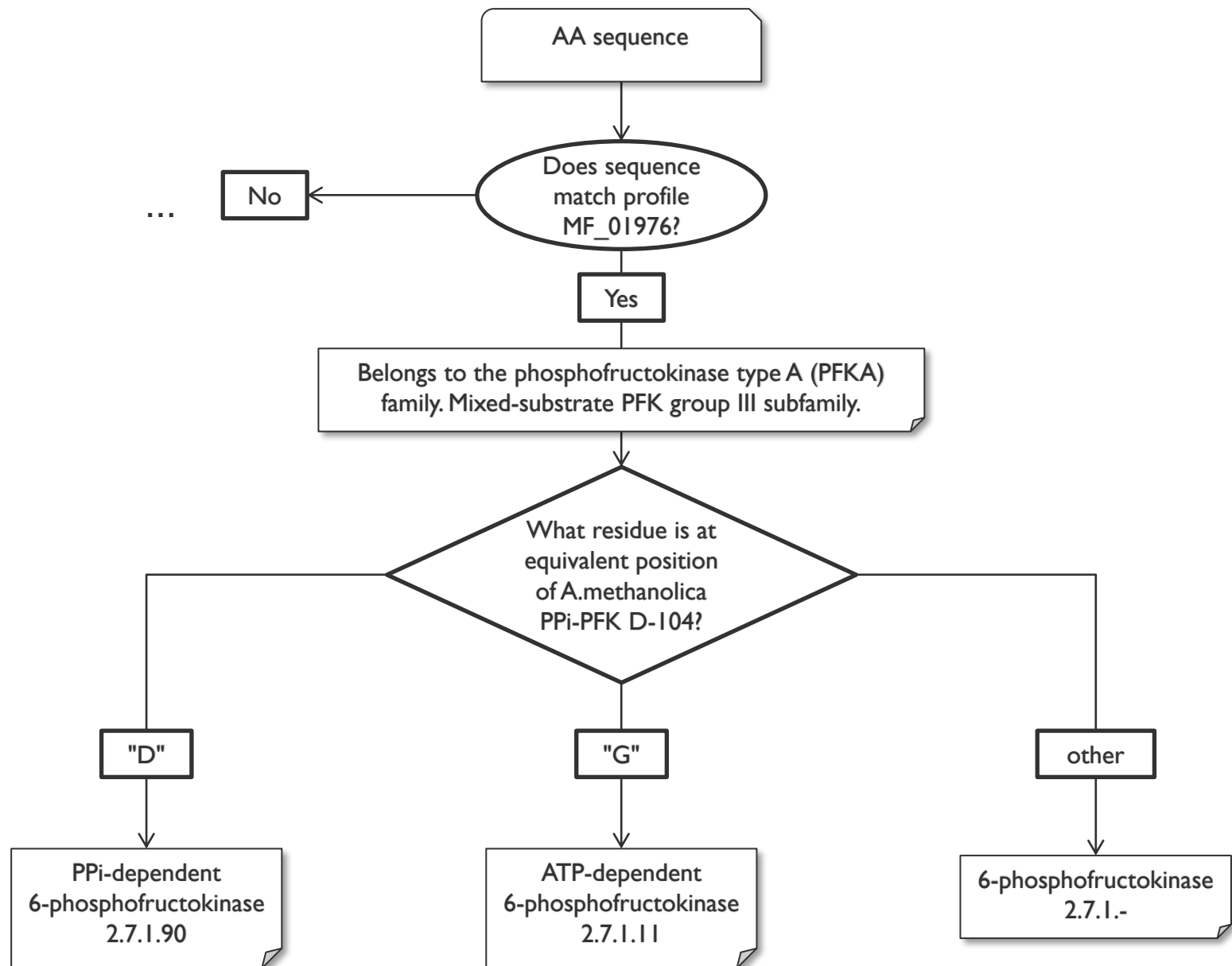
existence of numerous sequences from distantly related species carrying atypical combinations of amino acids. Several adaptive changes of phospho-donors, probably requiring a single mutation at position 104, have likely occurred independently in many lineages. The analysis of this gene suggests the existence of a high rate of both HGT and substitution in its active sites. These rampant HGT events and flexibility in phospho-donor use illustrate the importance of tinkering in molecular evolution.



MF\_01980



MF\_01981









## Annotation rule MF\_01976

[Send feedback](#)

### Features [\[?\]](#)

From: PFP\_AMYPE ([Q59126](#))

Key	From	To	Description	Tag	Condition	FTGroup
 SITE (Optional)	104	104	Important for catalytic activity and substrate specificity; stabilizes the transition state when the phosphoryl donor is P <sub>i</sub> ; prevents ATP from binding by mimicking the alpha-phosphate group of ATP	PPI	D  <span style="border: 1px solid black; padding: 2px;">"D"</span>	
<span style="border: 1px solid black; padding: 2px;">other</span>						
 SITE (Optional)	104	104	Important for substrate specificity; cannot use P <sub>i</sub> as phosphoryl donor	ATP	G  <span style="border: 1px solid black; padding: 2px;">"G"</span>	
REGION	125	127	Substrate binding		T-x-D	
REGION	169	171	Substrate binding		M-G- [RH]	
REGION	271	274	Substrate binding		[HY]-x (2) -R	
ACT_SITE	127	127	Proton acceptor		D	

## Identifier, protein and gene names [?]

"D"

case <FTTag:PPI>

Identifier	PFP
Protein name	RecName: Full=Pyrophosphate--fructose 6-phosphate 1-phosphotransferase; EC=2.7.1.90; AltName: Full=6-phosphofructokinase, pyrophosphate dependent; AltName: Full=PPI-dependent phosphofructokinase; Short=PPI-PFK; AltName: Full=Pyrophosphate-dependent 6-phosphofructose-1-kinase;
Gene name	pfp

"G"

else case <FTTag:ATP>

Identifier	PFKA
Protein name	RecName: Full=ATP-dependent 6-phosphofructokinase; Short=ATP-PFK; Short=Phosphofructokinase; EC=2.7.1.11; AltName: Full=Phosphohexokinase;
Gene name	pfkA

other

else

Identifier	PFP
Protein name	RecName: Full=6-phosphofructokinase; EC=2.7.1.-;

end case

## Annotation rule MF\_03125

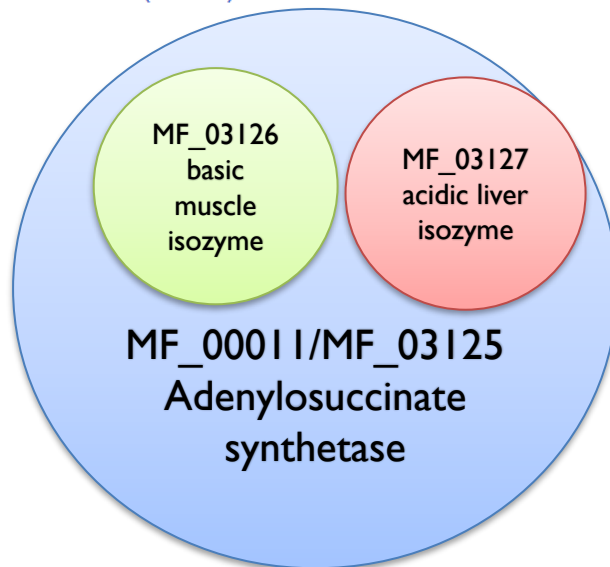
[Send feedback](#)

### Additional information [?]

Size range	411-710 amino acids
Related rules	<a href="#">MF_03126</a> (PURA1 <b>supersedes the current rule</b> ); <a href="#">MF_03127</a> (PURA2 <b>supersedes the current rule</b> )
Fusion	None

[View rule in raw text format \(no links\)](#)

We can indicate related HAMAP rules and specify their precedence





### Overview

[Proteins matched](#) (31019)

[Domain architectures](#) (41)

[Pathways & interactions](#)

[Species](#)

[Structures](#)

[Literature](#) (7)

[Cross-references](#) (5)

## F Family

# Adenylosuccinate synthetase (IPR001114)

*Short name: Adenylosuccinate\_synthetase*

## Overlapping homologous superfamilies i

[P-loop containing nucleoside triphosphate hydrolase](#) (IPR027417)

## Family relationships

**Adenylosuccinate synthetase (IPR001114)**

**Adenylosuccinate synthetase isozyme 1, chordates (IPR027509)**

**Adenylosuccinate synthetase isozyme 2, chordates (IPR027529)**

## Description

Adenylosuccinate synthetase ([EC:6.3.4.4](#)) plays an important role in purine biosynthesis, by catalysing the GTP-dependent conversion of IMP and aspartic acid to AMP. IMP and L-aspartate are conjugated in a two-step reaction accompanied by the hydrolysis of GTP to GDP in the presence of Mg<sup>2+</sup>. In the first step, the  $\gamma$ -phosphate group of GTP is transferred to the 6-oxygen atom of IMP. An aspartate then displaces this 6-phosphate group to form the product adenylosuccinate. Adenylosuccinate synthetase has been characterised from

[Add your annotation](#)

### Contributing signatures

Signatures from InterPro member databases are used to construct an entry.

**TIGRFAMs** i

[TIGR00184](#) (purA)

**Pfam** i

[PF00709](#)

(Adenylsucc\_synt)

**CDD** i

[cd03108](#) (AdSS)

**PANTHER** i

[PTHR11846](#)

(PTHR11846)

**SMART** i

[SM00788](#)

(Adenylsucc\_synt)

**HAMAP** i

[MF\\_00011](#)

(Adenylosucc\_synt)



### Overview

[Proteins matched](#) (93)

[Domain architectures](#) (1)

[Pathways & interactions](#)

[Species](#)

[Structures](#)

[Literature](#) (1)

[Cross-references](#) (1)

## Family

# Adenylosuccinate synthetase isozyme 1, chordates (IPR027509)

*Short name: AdSS\_1\_vert*

## Overlapping homologous superfamilies ⓘ

ⓘ [P-loop containing nucleoside triphosphate hydrolase \(IPR027417\)](#)

## Family relationships

↳ [Adenylosuccinate synthetase \(IPR001114\)](#)

↳ **[Adenylosuccinate synthetase isozyme 1, chordates \(IPR027509\)](#)**

## Description

This entry represents the adenylosuccinate synthetase isozyme 1 (AdSS1) ([EC:6.3.4.4](#)) from vertebrates. It is a component of the purine nucleotide cycle (PNC), which interconverts IMP and AMP to regulate the nucleotide levels in various tissues, and which contributes to glycolysis and ammoniagenesis. It catalyses the first committed step in the

Add your annotation

### Contributing signatures

Signatures from InterPro member databases are used to construct an entry.

■ **HAMAP ⓘ**

[MF\\_03126](#)

(Adenylosucc\_synth\_vert\_basic)



# Types of sequence similarity

**Family** Groups of proteins that are conserved along the whole sequence, sharing a common evolutionary origin, as reflected in their related functions.



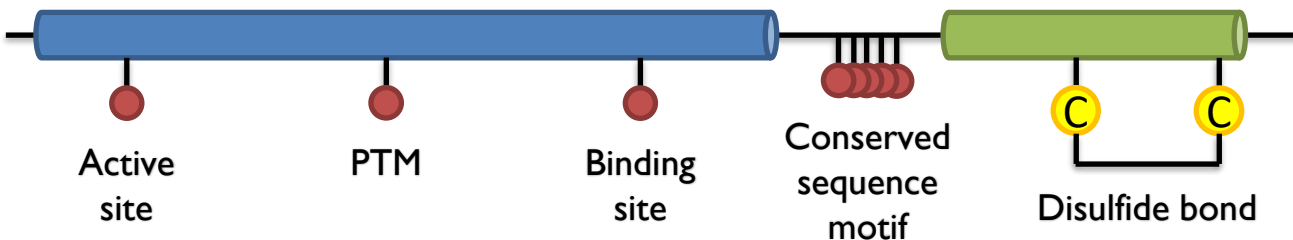
**Domains** Specific combination of secondary structures organized into a characteristic three dimensional structure or fold, that may exist in a variety of biological contexts.



**Repeats** Structural units typically repeated within a protein that assemble into a specific fold. Assemblies of repeats might also be thought of as domains.



**Sites & Motifs** Region of domains containing conserved active-site or binding residues or short conserved regions present outside domains that may adopt folded conformations only in association with their binding ligands



# ProRules

prosite.expasy.org/unirule/PRU00297



PROSITE

Home | Contact



annotation rule: PRU00297

Home ScanProsite ProRule Documents Downloads Links Funding

## General rule information

Accession	PRU00297
Dates	12-DEC-2003 (Created) 2-DEC-2010 (Last updated, Version 05)
Data class	Domain
Predictors	PROSITE; PS50873; PEROXIDASE_4
Name	Plant heme peroxidase
Function	Removal of H(2)O(2), oxidation of toxic reductants

Different classes of signature - domains, patterns, families

Generic function of a domain

## Propagated annotations

Conditional annotations (again)

## Comments

case <FTGroup:1> and <FTGroup:2> and <FTGroup:3> and <FTGroup:4>

FUNCTION Removal of H(2)O(2), oxidation of toxic reductants, biosynthesis and degradation of lignin, suberization, auxin catabolism, response to environmental stresses such as wounding, pathogen attack and oxidative stress. These functions might be dependent on each isozyme/isoform in each plant tissue.

CATALYTIC ACTIVITY Donor + H(2)O(2) = oxidized donor + 2 H(2)O.

# ProRules

← → ↻ [prosite.expasy.org/unirule/PRU00297](https://prosite.expasy.org/unirule/PRU00297) ☆



PROSITE

[Home](#) | [Contact](#)

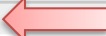


annotation rule: PRU00297

[Home](#) [ScanProsite](#) **[ProRule](#)** [Documents](#) [Downloads](#) [Links](#) [Funding](#)

Propagated annotation

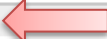
## Gene Ontology

case <Feature:PS50873:168=H>   
GO:0020037; Molecular function: heme binding.  
GO:0005506; Molecular function: iron ion binding.

## Keywords

Iron  
Heme  
end case

Coordinates refer to positions in profile PS50783 - no template sequence here

case <Feature:PS50873:1=Q>   
Pyrrolidone carboxylic acid  
end case

# ProRules

← → ↻ [prosite.expasy.org/unirule/PRU00297](http://prosite.expasy.org/unirule/PRU00297) ☆



PROSITE

Home | [Contact](#)

- Home
- ScanProsite
- ProRule**
- Documents
- Downloads
- Links
- Funding



Profile identifier

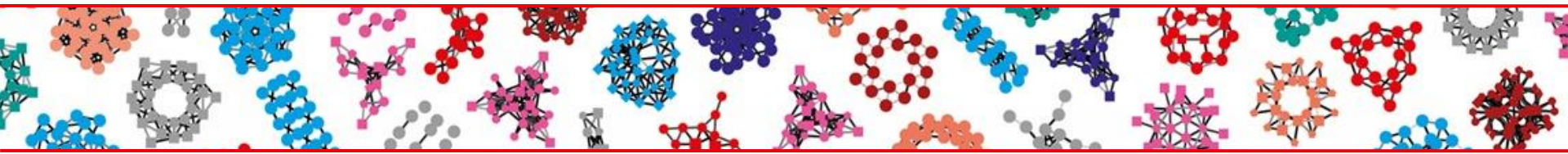
Positions in profile

annotation

## Features

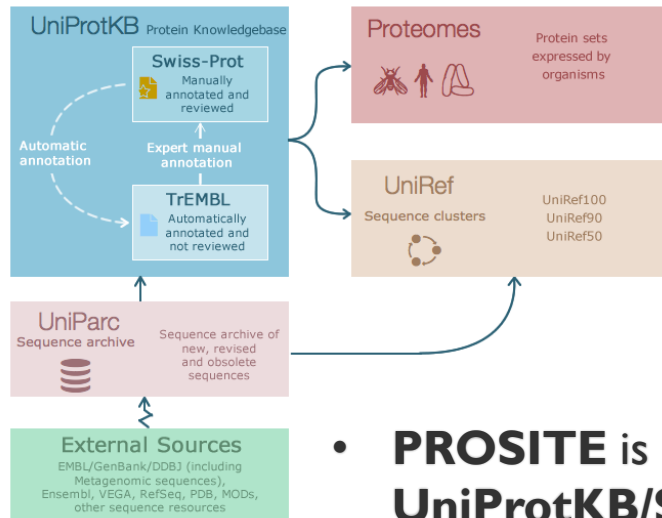
From: [PS50873](#)

Key	From	To	Description	Tag	Condition	FTGroup
METAL	43	43	Calcium #1 ( <i>By similarity</i> )		[DE]	1
METAL	46	46	Calcium #1; via carbonyl oxygen ( <i>By similarity</i> )		[VI]	1
METAL	48	48	Calcium #1; via carbonyl oxygen ( <i>By similarity</i> )		G	1
METAL	50	50	Calcium #1 ( <i>By similarity</i> )		[DE]	1
METAL	52	52	Calcium #1 ( <i>By similarity</i> )		[ST]	1



- Why do we need predictive annotation tools?
- Protein signatures for homology detection – A short primer
- Annotation rules for functional annotation
- **HAMAP and PROSITE - automatic annotation in UniProtKB**
- HAMAP and PROSITE - services for external users
- Practical exercises (afternoon)

# PROSITE annotation in UniProtKB/Swiss-Prot



- **PROSITE** is routinely used in the manual annotation in **UniProtKB/Swiss-Prot**
- It forms part of an integrated annotation tool '**Anabelle**' which presents PROSITE results in the context of other predictors for evaluation by Swiss-Prot curators
- PROSITE matches are then **evaluated** during curation and subsequently integrated into UniProtKB/Swiss-Prot records
- This curated information is used to calculate the reliability of patterns and profiles – indicated on PROSITE pages

# UniProtKB/Swiss-Prot annotation system

Function

Names & Taxonomy

Subcellular Location

Disease & Variants

PTM/Processing

Expression

Interaction

Structure

Family & Domains

Sequence

Similar Proteins

## Q8N3Z0 · PRS35\_HUMAN

Inactive serine protease 35 · Homo sapiens (Human) · Gene: PRSS35 (C6orf158) · 413 amino acids · Evidence at protein level · Annotation score: 3/5

[Entry](#) [Feature viewer](#) [Publications](#) [External links](#) [History](#)

[BLAST](#) [Align](#) [Download](#) [Add](#) [Add a publication](#) [Entry feedback](#)

### Sequence<sup>i</sup>

Sequence status<sup>i</sup> | Complete

Sequence processing<sup>i</sup> | The displayed sequence is further processed into a mature form.

[Tools](#) [Download](#) [Add](#) [Highlight](#) [Copy sequence](#)

Length 413

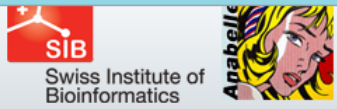
Mass (Da) 47,098

Last updated 2007-09-11 v2

Checksum<sup>i</sup> 818D9C951BD2D6C1

MENMLLWLI <sup>10</sup> F	FTPGWTLIDG <sup>20</sup>	SEMEWDFMWH <sup>30</sup>	LRKVPRIVSE <sup>40</sup>	RTFHLTSPAF <sup>50</sup>	EADAKMMVNT <sup>60</sup>	VCGIECQKEL <sup>70</sup>	PTPLSELED <sup>80</sup>	YLSYETVFEN <sup>90</sup>	GRTLTRVKV <sup>100</sup>	QDLVLEPTQN <sup>110</sup>	ITTKGVSVRR <sup>120</sup>	KRQVYGTDSR <sup>130</sup>	FSILDKRFLT <sup>140</sup>
NFPFSTAVKL <sup>150</sup>	STGCSGILIS <sup>160</sup>	PQHVLTA <sup>170</sup> AHC	VHDGKDYVKG <sup>180</sup>	SKKLRVGLLK <sup>190</sup>	MRNKSGGKRR <sup>200</sup>	RGSKRSRREA <sup>210</sup>	SGGQREGTR <sup>220</sup>	EHLRERAKGG <sup>230</sup>	RRRKSGRGQ <sup>240</sup>	RIAEGRPSFQ <sup>250</sup>	WTRVKNTHIP <sup>260</sup>	KGWARGMGD <sup>270</sup>	ATLDYDYALL <sup>280</sup>
ELKRAHKKKY <sup>290</sup>	MELGISPTIK <sup>300</sup>	KMPGGMIHFS <sup>310</sup>	GFDNDRADQL <sup>320</sup>	VYRFCVSDE <sup>330</sup>	SNDLLYQYCD <sup>340</sup>	AESGSTGSGV <sup>350</sup>	YLRLKDPDKK <sup>360</sup>	NWKRKIIAVY <sup>370</sup>	SGHQWVDVHG <sup>380</sup>	VQKDYNVAVR <sup>390</sup>	ITPLKYAQIC <sup>400</sup>	LWIHGNDANC <sup>410</sup>	AYG

back

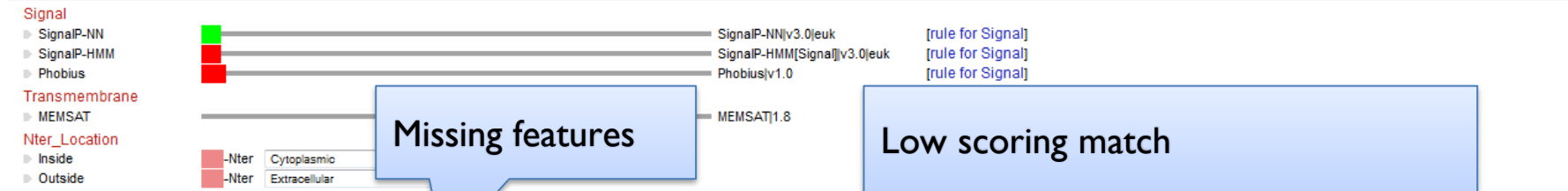


color codes quick help tutorial anabelle full technical doc

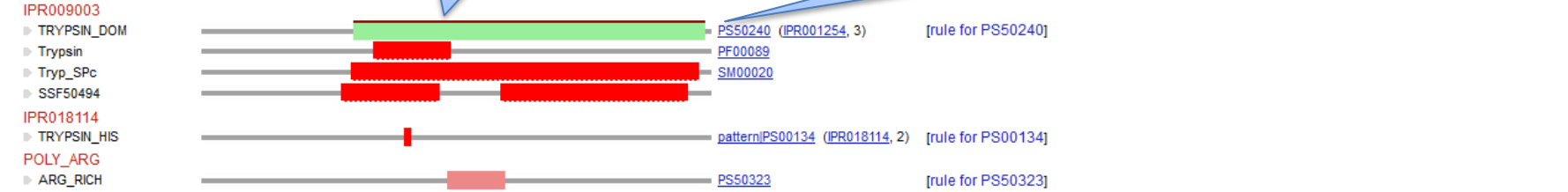
### Sequence Analysis Log

PRSS35\_HUMAN: [Size: 413] [blast](#) [String](#) [HoverProt links](#) [Family](#) [Alignment](#) [Tree](#) [preview](#) Zoom: 100%

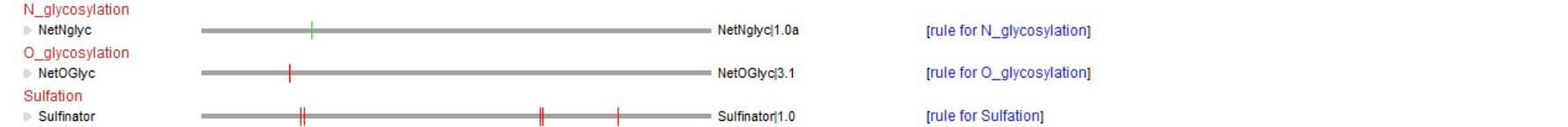
#### TOPOLOGY



#### DOMAIN

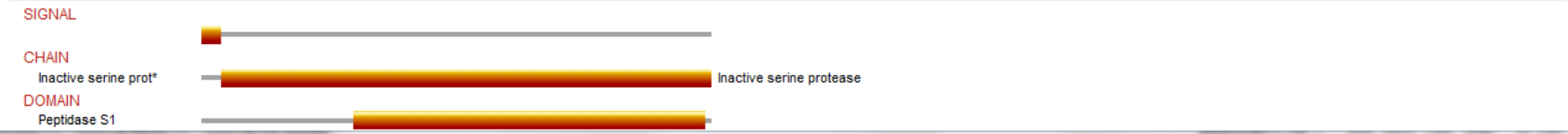


#### SITE



#### SKIP

#### CURRENT\_FT





- Function
- Names & Taxonomy
- Subcellular Location
- Disease & Variants
- PTM/Processing
- Expression
- Interaction
- Structure
- Family & Domains
- Sequence
- Similar Proteins


## Q8N3Z0 · PRS35\_HUMAN

Inactive serine protease 35 · Homo sapiens (Human) · Gene: PRSS35 (C6orf158) · 413 amino acids · Evidence at protein level · Annotation score

[Entry](#) [Feature viewer](#) [Publications](#) [External links](#) [History](#)

BLAST [Align](#) [Download](#) [Add](#) [Add a publication](#) [Entry feedback](#)

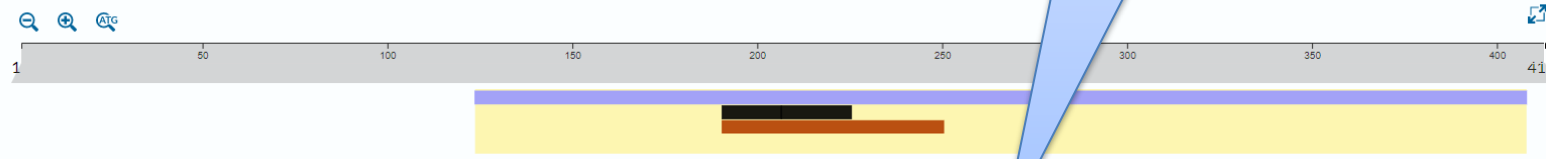
### Function<sup>i</sup>


**Caution**  
 Although related to peptidase S1 family, lacks the conserved active Ser residue in position 346 which is replaced by a Thr, suggesting that it has no protease activity.  Curated

### Family & Domains<sup>i</sup>

#### Features

Showing features for domain<sup>i</sup>, compositional bias<sup>i</sup>, region<sup>i</sup>.



TYPE	ID	POSITION(S)	DESCRIPTION	
Domain				
▶ Domain		124-408	Peptidase S1	BLAST  Add

### Sequence Similarities<sup>i</sup>

Belongs to the peptidase S1 family.  Curated

The domain is nevertheless annotated following manual evaluation by the curator...

[Expand table](#)

## General information about the entry

Entry name <a href="#">[info]</a>	TRYPSIN_DOM
Accession <a href="#">[info]</a>	PS50240
Entry type <a href="#">[info]</a>	MATRIX
Date <a href="#">[info]</a>	DEC-2001 (CREATED); OCT-2013 (DATA UPDATE); APR-2015 (INFO UPDATE).
PROSITE Doc. <a href="#">[info]</a>	<a href="#">PDOC00124</a>
Associated ProRule <a href="#">[info]</a>	<a href="#">PRU00274</a>

## Name and characterization of the entry

Description <a href="#">[info]</a>	Serine proteases, trypsin domain profile.
Matrix / Profile <a href="#">[info]</a>	<pre> /GENERAL_SPEC: ALPHABET='ABCDEFGHIJKLMNPQRSTUVWXYZ'; LENGTH=234; /DISJOINT: DEFINITION=PROTECT; N1=6; N2=229; /NORMALIZATION: MODE=1; FUNCTION=LINEAR; R1=0.0169; R2=0.00836256; TEXT='NScore'; /NORMALIZATION: MODE=-1; FUNCTION=LINEAR; R1=6057.56402; R2=28.29509; TEXT='HScore'; /CUT_OFF: LEVEL=0; SCORE=1134; N_SCORE=9.5; H_SCORE=34749; MODE=1; TEXT='!'; /CUT_OFF: LEVEL=-1; SCORE=775; N_SCORE=6.5; H_SCORE=27986; MODE=1; TEXT='?'; /DEFAULT: M0=-9; D=-20; I=-20; B1=-60; E1=-60; MI=-105; MD=-105; IM=-105; DM=-105; ... </pre> <p style="text-align: right;"><a href="#">» More</a></p>

## Numerical results [\[info\]](#)

Numerical results for UniProtKB/Swiss-Prot release **2015\_06** which contains **548'586** sequence entries.

Total number of hits	759 in <a href="#">739 different sequences</a>
Number of true positive hits	759 in <a href="#">739 different sequences</a>
Number of 'unknown' hits	
Number of false positive hits	
Number of false negative sequences	<b>17</b>
Number of 'partial' sequences	<b>80</b>
Precision (true positives / (true positives + false positives))	100.00 %
Recall (true positives / (true positives + false negatives))	97.81 %

**...and tagged as false negative**

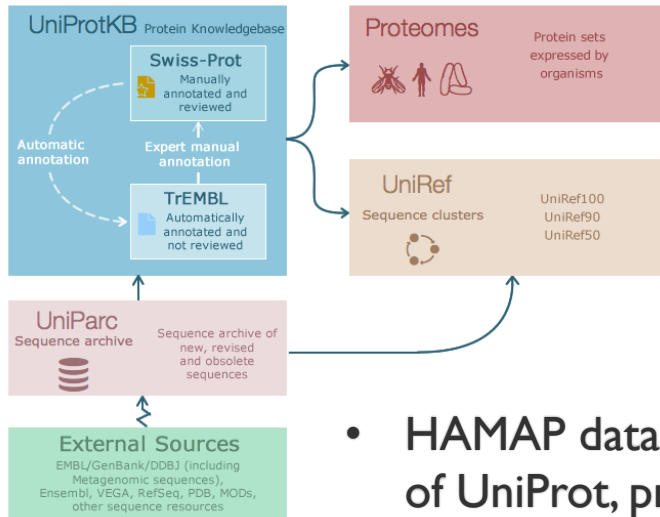
Feature key <a href="#">[info]</a>	DOMAIN
Feature description <a href="#">[info]</a>	Peptidase S1
Version <a href="#">[info]</a>	2

### Cross-references [\[info\]](#)

UniProtKB/Swiss-Prot True positive sequences	<p>739 sequences</p> <p>ACH1_LONAC (<a href="#">P23604</a>), ACH2_LONAC (<a href="#">P23605</a>), ACRO_HUMAN (<a href="#">P10323</a>),  <a href="#">» More</a></p> <ul style="list-style-type: none"> <li>Retrieve an alignment of UniProtKB/Swiss-Prot true positive hits:  <a href="#">[Clustal format, color, condensed view]</a> <a href="#">[Clustal format, color]</a>  <a href="#">[Clustal format, plain text]</a> <a href="#">[Fasta format]</a></li> <li>Retrieve the sequence logo from the alignment</li> </ul>
UniProtKB/Swiss-Prot False negative sequences	<p>17 sequences</p> <p>HTR1B_DANRE (<a href="#">A9JRB3</a>), HTRA1_BOVIN (<a href="#">F1N152</a>), HTRA1_XENLA (<a href="#">A6YFB5</a>),  HTRA1_XENTR (<a href="#">A4IHA1</a>), PRS23_BOVIN (<a href="#">Q1LZE9</a>), PRS23_HUMAN (<a href="#">O95084</a>),  PRS23_MACMU (<a href="#">Q1WK23</a>), PRS23_MOUSE (<a href="#">Q9E0M6</a>), PRS23_RAT (<a href="#">Q6AY61</a>),  PRS35_BOVIN (<a href="#">Q5E9X0</a>), PRS35_HUMAN (<a href="#">Q8N3Z0</a>), PRS35_MACMU (<a href="#">Q1WK24</a>),  PRS35_MOUSE (<a href="#">Q8C0F9</a>), PRS35_RAT (<a href="#">Q5E212</a>), Y3671_MYCTO (<a href="#">P9WHR8</a>),  Y3671_MYCTU (<a href="#">P9WHR9</a>), YH05_SCHPO (<a href="#">O74325</a>)  <a href="#">Less «</a></p>
UniProtKB/Swiss-Prot 'Partial' sequences	<p>80 sequences</p> <p>ACRO_CAPHI (<a href="#">P10626</a>), CATG_RAT (<a href="#">P17977</a>), CBP_MESMA (<a href="#">P0C8M2</a>),  <a href="#">» More</a></p>
PDB <a href="#">[Detailed view]</a>	<p>1685 PDB</p> <p>1A0H; 1A0J; 1A0L; 1A2C; 1A3B; 1A3E; 1A46; 1A4W; 1A5G; 1A5H; 1A5I; 1A61;  <a href="#">» More</a></p>

[View entry in original PROSITE format](#)  
[View entry in raw text format \(no links\)](#)  
[Direct ScanProsite submission](#)

# HAMAP annotation in UniProtKB



- HAMAP data are incorporated by InterPro and the UniRule pipeline of UniProt, providing annotation of UniProtKB/Swiss-Prot quality for millions of unreviewed protein sequences in UniProtKB/TrEMBL
- Each month, new sequences entering UniProtKB/TrEMBL are scanned against the collection of HAMAP profiles
- Matching sequences are evaluated against corresponding rule conditions and annotated accordingly
- Results are made available in the corresponding UniProtKB/TrEMBL records (search for 'source:HAMAP')

- Function
- Names & Taxonomy
- Subcellular Location
- Phenotypes & Variants
- PTM/Processing
- Expression
- Interaction
- Structure
- Family & Domains
- Sequence
- Similar Proteins

## Q63RY2 · Q63RY2\_BURPS

Adenine deaminase · Burkholderia pseudomallei (strain K96243) · EC:3.5.4.2 · Gene: add · 341 amino acids · Inferred from homology · Annotation score: 3/5

Entry Feature viewer Publications External links History

BLAST Align Download Add Add a publication Entry feedback

### Family & Domains<sup>1</sup>

#### Features

Showing features for domain<sup>1</sup>.

TYPE	ID	POSITION(S)	DESCRIPTION
Domain		13-336	Adenine deaminase domain

**Cross-reference links back to HAMAP profile**

#### Sequence Similarities<sup>1</sup>

Belongs to the metallo-dependent hydrolases superfamily, Adenosine and AMP deaminases family, Adenine deaminase subfamily. 1 Automatic Annotation

#### Phylogenomic databases

OMA	DERLMQR <a href="#">↗</a>	eggNOG	COG1816 <a href="#">↗</a> Bacteria
-----	---------------------------	--------	------------------------------------

#### Family and domain databases

CDD	cd01320 <a href="#">↗</a> ADA 1 hit	PANTHER	PTHR43114 <a href="#">↗</a> PTHR43114 1 hit
HAMAP	MF_01962 <a href="#">↗</a> Adenine_deaminase 1 hit	Pfam	View protein in Pfam <a href="#">↗</a> PF00962 <a href="#">↗</a> A_deaminase 1 hit
InterPro	View protein in InterPro <a href="#">↗</a> IPR001365 <a href="#">↗</a> A_deaminase_dom IPR028892 <a href="#">↗</a> ADE IPR006330 <a href="#">↗</a> Ado/ade_deaminase IPR032466 <a href="#">↗</a> Metal_Hydrolase	SUPFAM	SSF51556 <a href="#">↗</a> SSF51556 1 hit
		TIGRFAMs	TIGR01430 <a href="#">↗</a> aden_deam 1 hit
		MobiDB	Search... <a href="#">↗</a>
		ProtoNet	Search... <a href="#">↗</a>



Search HAMAP

## Family profile MF\_01962

### General profile information [?]

Accession	MF_01962	<a href="#">[View profile]</a>
Entry name	Adenine_deaminase	<a href="#">[View seed alignment]</a>
Entry type	MATRIX	
Date	DEC-2013 (DATA UPDATE).	
Version	4	
Description	Adenine deaminase.	
Taxonomic range	Archaea, Bacteria, Eukaryota	
InterPro	<a href="#">IPR028892</a> ADE	
Associated rules	<a href="#">MF_01962</a> name: Adenine_deaminase scope: Bacteria <a href="#">MF_03145</a> name: Adenine_deaminase_euk scope: Eukaryota	

### Statistics [?]

<b>Number of hits in UniProtKB</b>	1,867	<a href="#">[Graphical view of score distribution]</a>
• Number of hits in UniProtKB/Swiss-Prot	70	
• Number of hits in UniProtKB/TrEMBL	1,797	
<b>Taxonomic distribution of hits in UniProtKB</b>		<a href="#">[View taxonomic distribution of UniProtKB matches]</a>
• Archaea	4	<a href="#">[Taxonomic distribution in UniProtKB complete proteomes]</a>
• Bacteria	1,515	
• Eukaryota	347	
• unclassified sequences	1	

- Function
- Names & Taxonomy
- Subcellular Location
- Phenotypes & Variants
- PTM/Processing
- Expression
- Interaction
- Structure
- Family & Domains
- Sequence
- Similar Proteins

## Q63RY2 · Q63RY2\_BURPS

Adenine deaminase · Burkholderia pseudomallei (strain K96243) · EC:3.5.4.2 · Gene: add · 341 amino acids · Inferred from homology · Annotation score: 3/5

Entry Feature viewer Publications External links History

BLAST Align Download Add Add a publication Entry feedback

### Function<sup>i</sup>

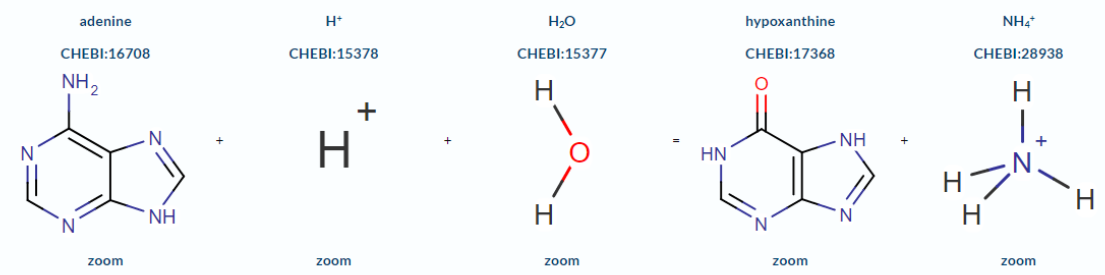
Catalyzes the hydrolytic deamination of adenine to hypoxanthine. Plays an important role in the purine salvage pathway and in nitrogen catabolism. 1 Automatic Annotation

Automatic assertion according to rules (Automatically inferred from sequence model)<sup>1</sup>  
 UniRule HAMAP-Rule: MF\_01962

### Catalytic Activity

adenine + H<sup>+</sup> + H<sub>2</sub>O = hypoxanthine + NH<sub>4</sub><sup>+</sup> 1 Automatic Annotation  
 EC:3.5.4.2 (UniProtKB | ENZYME | Rhea )  
 Source: Rhea 23688

Evidence attribution links back to annotation rule



### Cofactor

Zn(2+) (UniProtKB | Rhea | CHEBI:29105 ) 1 Automatic Annotation

## UniRule - UR000089071

[Download](#) | [View proteins](#)Source Rule: [MF\\_01962](#)The annotation and conditions in this rule are derived from the following entries: [Q916Y4](#) and [P53909](#)

IF		THEN	
HAMAP signature	<a href="#">MF_01962</a>	protein name	Adenine deaminase[...]
taxon	Bacteria	catalytic activity	adenine + H <sup>+</sup> + H <sub>2</sub> O = hypoxanthine + NH <sub>4</sub> <sup>+</sup> EC:3.5.4.2 ( <a href="#">UniProtKB</a>   <a href="#">ENZYME</a>   <a href="#">Rhea</a> ) Source: <a href="#">Rhea 23688</a>
fragmented	NO		<a href="#">View Rhea reaction</a>
		function	Catalyzes the hydrolytic deamination of adenine to hypoxanthine. Plays an important role in the purine salvage pathway and in nitrogen catabolism
		similarity	Belongs to the metallo-dependent hydrolases superfamily, Adenosine deaminases family, Adenine deaminase type 2 subfamily
		keyword	Hydrolase Nucleotide metabolism
		GO term	<a href="#">GO:0000034</a> <a href="#">GO:0006146</a> <a href="#">GO:0043103</a>
ADDITIONALLY			
IF		THEN	
positional features	14: H 194: H 16: H 275: D	cofactor	Zn(2+) ( <a href="#">UniProtKB</a>   <a href="#">PDB</a> :29105 ) Binds 1 zinc ion per subunit
		keyword	Metal-binding Zinc
		GO term	<a href="#">GO:0008270</a>

CONDITIONS

ANNOTATIONS



Status

Unreviewed (TrEMBL) (7,912)

# UniProtKB 7,912 results

BLAST Align Map IDs Download Add View: Cards Table Customize columns Share

Taxonomy

Filter by taxonomy

Proteins with

3D structure (1)

Active site (7,325)

Binding site (7,909)

Catalytic activity (7,912)

Chain (18)

More items

Protein existence

Homology (7,910)

Protein level (1)

Transcript level (1)

Annotation score

3 (5,744)

2 (2,168)

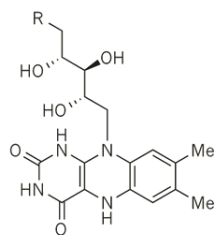
Sequence length

201 - 400 (7,900)

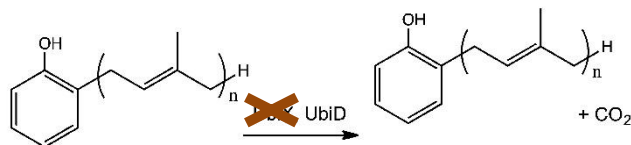
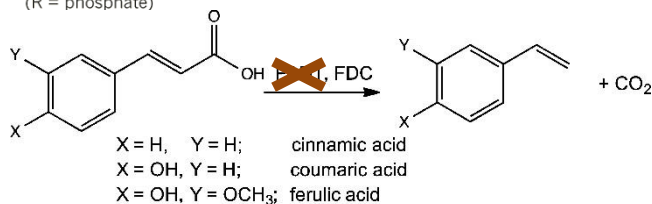
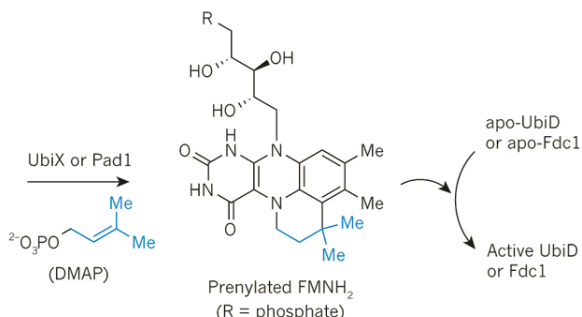
401 - 600 (10)

601 - 800 (2)

Entry	Entry Name	Protein Names	Gene Names	Organism
<input type="checkbox"/> A0A8F3VRU4	A0A8F3VRU4_ACIP1	Adenine deaminase, ADE, 3.5.4.2, Adenine aminohydrolase, AAH	add_1, OCUAc17_25480	Acinetobacter pittii (Acinetobacter genomosp. 3)
<input type="checkbox"/> A0A0Q0C003	A0A0Q0C003_PSEAJ	Adenine deaminase, ADE, 3.5.4.2, Adenine aminohydrolase, AAH	ALO60_02379	Pseudomonas amygdali pv. tabaci (Pseudomonas syringae pv. ta)
<input type="checkbox"/> A0A1C2VBD9	A0A1C2VBD9_ACIP1	Adenine deaminase, ADE, 3.5.4.2, Adenine aminohydrolase, AAH	B9X68_08750	Acinetobacter pittii (Acinetobacter genomosp. 3)
<input type="checkbox"/> A0A1C4GHH0	A0A1C4GHH0_9NOCA	Adenine deaminase, ADE, 3.5.4.2, Adenine aminohydrolase, AAH	RQCS_38680	Rhodococcus qingshengii
<input type="checkbox"/> A0A4Y3IV08	A0A4Y3IV08_ACILW	Adenine deaminase, ADE, 3.5.4.2, Adenine aminohydrolase, AAH	add, AL1T_00220	Acinetobacter lwoffii
<input type="checkbox"/> A0A6F9YQ48	A0A6F9YQ48_9LACO	Adenine deaminase, ADE, 3.5.4.2, Adenine aminohydrolase, AAH	add, PTL465_16780	Ligilactobacillus agilis
<input type="checkbox"/> A0A7V1K9E9	A0A7V1K9E9_9GAMM	Adenine deaminase, ADE, 3.5.4.2, Adenine aminohydrolase, AAH	ENI21_09280	Halopseudomonas sabulinigri
<input type="checkbox"/> A0A7V6CCL5	A0A7V6CCL5_9GAMM	Adenine deaminase, ADE, 3.5.4.2, Adenine aminohydrolase, AAH	ENK16_07150	Chromatiales bacterium
<input type="checkbox"/> C8ZXX9	C8ZXX9_ENTGE	Adenine deaminase, ADE, 3.5.4.2, Adenine aminohydrolase, AAH	EGBG_00546	Enterococcus gallinarum (strain EG2)
<input type="checkbox"/> C9YHA2	C9YHA2_CURXX	Adenine deaminase, ADE, 3.5.4.2, Adenine aminohydrolase, AAH	add, Csp_B21520	Curvibacter symbiont subsp. Hydra magnipapillata
<input type="checkbox"/> D0SMM8	D0SMM8_ACIJU	Adenine deaminase, ADE, 3.5.4.2, Adenine aminohydrolase, AAH	add, HMPREF0026_00561	Acinetobacter junii SH205
<input type="checkbox"/> D0STI9	D0STI9_ACILW	Adenine deaminase, ADE, 3.5.4.2, Adenine aminohydrolase, AAH	add, HMPREF0017_00613	Acinetobacter lwoffii SH145
<input type="checkbox"/> D6S693	D6S693_LACJE	Adenine deaminase, ADE, 3.5.4.2, Adenine aminohydrolase, AAH	add, HMPREF0526_11415	Lactobacillus jensenii JV-V16
<input type="checkbox"/> E2XKS5	E2XKS5_PSEFL	Adenine deaminase, ADE, 3.5.4.2, Adenine aminohydrolase, AAH	add, PFWH6_0677	Pseudomonas fluorescens WH6
<input type="checkbox"/> G2ZZZ4	G2ZZZ4_9RAL5	Adenine deaminase, ADE, 3.5.4.2, Adenine aminohydrolase, AAH	add, RALSY_10442	Ralstonia syzygii R24
<input type="checkbox"/> J2MNG7	J2MNG7_PSEFL	Adenine deaminase, ADE, 3.5.4.2, Adenine aminohydrolase, AAH	add, PflQ2_5014	Pseudomonas fluorescens Q2-87
<input type="checkbox"/> S5IWQ1	S5IWQ1_VIBPH	Adenine deaminase, ADE, 3.5.4.2, Adenine aminohydrolase, AAH	M634_20200	Vibrio parahaemolyticus O1:Kuk str. FDA_R31



Riboflavin: R = OH  
 FMNH<sub>2</sub>: R = phosphate



**nature**

International weekly journal of science

Home | News & Comment | Research | Careers & Jobs | Current Issues

Archive | Volume 522 | Issue 7557 | Letters | Article

NATURE | LETTER

日本語要約

## UbiX is a flavin prenyltransferase required for bacterial ubiquinone biosynthesis

Mark D. White, Karl A. P. Payne, Karl Fisher, [Stephen A. Mars](#),  
 J. W. Rattray, Drupad K. Trivedi, Royston Goodacre, Stephen  
 Scrutton, Sam Hay & David Leys

Affiliations | Contributions | Corresponding author

Nature 522, 502–506 (25 June 2015) | doi:10.1038/nature14559

```

DT 05-NOV-2015 (Rel. 126, Created)
DT 24-MAR-2016 (Rel. 128, Last updated, Version 4)
XX
DE Klebsiella pneumoniae 3-octaprenyl-4-hydroxybenzoate carboxy-lyase
XX
.
FT CDS CP012753.1:5179010..5179594
FT /codon_start=1
FT /transl_table=11
FT /locus_tag="AN966_25510"
FT /product="3-octaprenyl-4-hydroxybenzoate carboxy-lyase"
FT /EC_number="4.1.1.-"
FT /note="catalyzes the formation of 2-octaprenylphenol from
FT 3-octaprenyl-4-hydroxybenzoate; Derived by automated
FT computational analysis using gene prediction method:
FT Protein Homology."
FT /inference="EXISTENCE:similar to AA
FT sequence:RefSeq:WP_017899515.1"
FT /protein_id="ALH88064.1"
FT /translation="MKRLIIGISGASGAIYGVRLQVLRDVPDIETHLILSQARQTLA
FT METDFTVREVQALADVVDHARDIAASISSGSFKTAGMVILPCSMKTLGSIHVSYDGLL
FT TRAADVVLKERRPLVLCVRETPFHLGHLRLLVQAAELGAVIMPPVPAFYHRPQSLDDVI
FT NQTVNRVLDQFDISLEQDLFTRWQGSQDCTK"
XX
  
```

Search CiteSpace  
 Enter search text / DOI  
 ACS Chem. Biol.

Issue Multimedia Submission & Review Open Access

< Previous Article

## UbiX Catalyze Formation of a Novel 3-Octaprenyl-4-Hydroxybenzoate Decarboxylase and 4-Hydroxybenzoate

ert, Xiaoxia Nina Lin<sup>†</sup>, and E. Neil G. Marsh<sup>†§</sup>  
 engineering, and <sup>§</sup>Department of Biological Chemistry, University of

Protein | **Flavin prenyltransferase UbiX**

Gene | **pad1\_1**

Organism | *Klebsiella pneumoniae*

Status | Unreviewed - Annotation score: - Protein inferred from homology<sup>i</sup>

# A0A0C2MA92\_KLEPN

## Function<sup>i</sup>

Flavin prenyltransferase that catalyzes the synthesis of the prenylated FMN cofactor (prenyl-FMN) for 4-hydroxy-3-polyprenylbenzoic acid decarboxylase UbiD. The prenyltransferase is metal-independent and links a dimethylallyl moiety from dimethylallyl monophosphate (DMAP) to the flavin N5 and C6 atoms of FMN. [UniRule annotation](#) ▼

### Catalytic activity<sup>i</sup>

Dimethylallyl phosphate + FMNH<sub>2</sub> = prenylated FMNH<sub>2</sub> + phosphate. [UniRule annotation](#) ▼

### Sites

Automatic assertion according to rules<sup>i</sup>






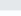

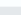

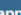

HAMAP-Rule:MF\_01984

Feature key	Position(s)	Length	Graphical view	Feature identifier	Actions
Binding site <sup>i</sup>	37 – 37				
Binding site <sup>i</sup>	123 – 123	1		FMN <a href="#">UniRule annotation</a> ▼	
Binding site <sup>i</sup>	153 – 153	1		DMAP <a href="#">UniRule annotation</a> ▼	

- Proteins with
- 3D structure (8)
- Active site (1)
- Beta strand (5)
- Binding site (18,436)
- Biotechnological use (1)
- More items
- Protein existence
- Homology (18,426)
- Protein level (9)
- Transcript level (1)
- Annotation score
- 5 (1)
- 4 (1)
- 3 (32)
- 2 (18,402)
- Sequence length
- 1 - 200 (10,266)
- 201 - 400 (8,148)
- 401 - 600 (11)
- 601 - 800 (7)
- >= 801 (4)

## UniProtKB 18,436 results

BLAST Align Map IDs  Download  Add View: Cards  Table   Customize columns  Share

Entry	Entry Name	Protein Names	Gene Names	Organism
<input type="checkbox"/> P33751	 PAD1_YEAST	Flavin prenyltransferase PAD1, mitochondrial, 2.5.1.129, Phenylacrylic acid decarboxylase 1, PAD	PAD1, POF1, YDR538W	Saccharomyces cerevisiae (strain ATCC 204508 / S288c) (Baker's yeast)
<input type="checkbox"/> P0AG03	 UBIX_ECOLI	Flavin prenyltransferase UbiX, 2.5.1.129	ubiX, defF, b2311, JW2308	Escherichia coli (strain K12)
<input type="checkbox"/> Q9HX08	 UBIX_PSEAE	Flavin prenyltransferase UbiX, 2.5.1.129	ubiX, PA4019	Pseudomonas aeruginosa (strain ATCC 15692 / DSM 22644 / CIP 104116 / JCM 14847 / LMG 12228 / 1C / PRS 101 / PAO1)
<input type="checkbox"/> P0AG04	 UBIX_ECO57	Flavin prenyltransferase UbiX, 2.5.1.129	ubiX, Z3573, ECs3195	Escherichia coli O157:H7
<input type="checkbox"/> Q9V030	 UBIX_PYRAB	Flavin prenyltransferase UbiX, 2.5.1.129	ubiX, PYRAB09680, PAB0651	Pyrococcus abyssi (strain GE5 / Orsay)
<input type="checkbox"/> Q9ZD09	 UBIX_RICPR	Flavin prenyltransferase UbiX, 2.5.1.129	ubiX, RP541	Rickettsia prowazekii (strain Madrid E)
<input type="checkbox"/> Q9ZJE3	 UBIX_HELPJ	Flavin prenyltransferase UbiX, 2.5.1.129	ubiX, jhp_1369	Helicobacter pylori (strain J99 / ATCC 700824) (Campylobacter pylori J99)
<input type="checkbox"/> P57767	 UBIX_THAAR	Flavin prenyltransferase UbiX, 2.5.1.129	ubiX	Thauera aromatica
<input type="checkbox"/> Q9JXP4	 UBIX_NEIMB	Flavin prenyltransferase UbiX, 2.5.1.129	ubiX, NMB1945	Neisseria meningitidis serogroup B (strain MC58)
<input type="checkbox"/> Q9PKH2	 UBIX_CHLMU	Flavin prenyltransferase UbiX, 2.5.1.129	ubiX, TC_0493	Chlamydia muridarum (strain MoPn / Nigg)
<input type="checkbox"/> Q9Y8K8	 UBIX_SACS2	Flavin prenyltransferase UbiX, 2.5.1.129	ubiX, SSO0437, C41_048	Saccharolobus solfataricus (strain ATCC 35092 / DSM 1617 / JCM 11322 / P2) (Sulfolobus solfataricus)
<input type="checkbox"/> O26011	 UBIX_HELPY	Flavin prenyltransferase UbiX, 2.5.1.129	ubiX, HP_1476	Helicobacter pylori (strain ATCC 700392 / 26695) (Campylobacter pylori)

Feedback

Help

## A0A0C2MA92 · A0A0C2MA92\_KLEPN

Flavin prenyltransferase UbiX · *Klebsiella pneumoniae* · EC:2.5.1.129 · Gene: pad1\_2 (pad1\_1, ubiX) · 194 amino acids · Inferred from homology · Annotation score: 2/5

Entry Feature viewer Publications External links History

BLAST Align  Download  Add Add a publication Entry feedback

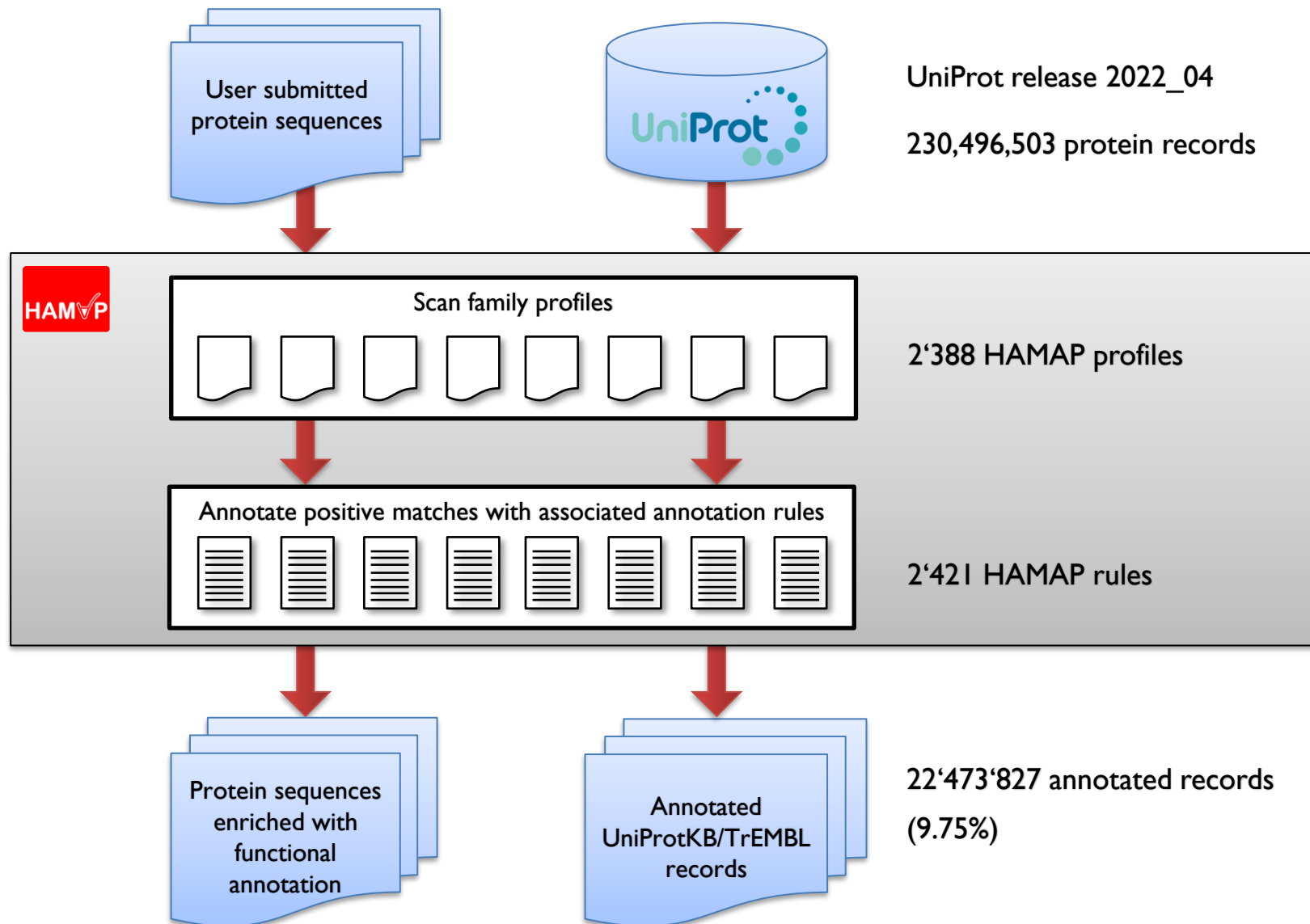
### Function<sup>i</sup>

Flavin prenyltransferase that catalyzes the synthesis of the prenylated FMN cofactor (prenyl-FMN) for 4-hydroxy-3-polyprenylbenzoic acid decarboxylase UbiD. The prenyltransferase is metal-independent and links a dimethylallyl moiety from dimethylallyl monophosphate (DMAP) to the flavin N5 and C6 atoms of FMN. 1 Automatic Annotation

Automatic assertion according to rules (Automatically inferred from sequence model)<sup>i</sup>  
 UniRule HAMAP-Rule: MF\_01984

Feedback

# Current status





## Reference proteomes: archaea

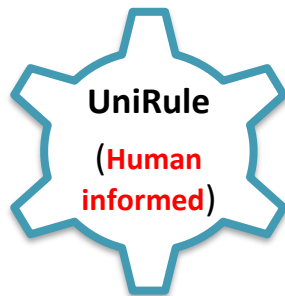
 

To retrieve information from a particular organism (description, genome reference paper(s), taxonomy or UniProtKB entries), click on the UniProtKB proteome identifier.

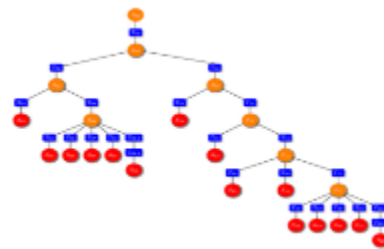
### Statistics

Species	Species code	UniProtKB proteome identifier	Number of entries in UniProtKB			HAMAP coverage in UniProtKB
			All	Swiss-Prot	TrEMBL	
Acidianus hospitalis (strain W1)	ACIHW	<a href="#">UP000008458</a>	2329	0	2329	14%
Acidilobus saccharovorans (strain DSM 16705 / JCM 18335 / VKM B-2471 / 345-15)	ACIS3	<a href="#">UP000000346</a>	1499	2	1497	16%
Escherichia coli (strain K12)	ECOLI	<a href="#">UP000000625</a>	4305	4305	0	26%
Aeropyrum pernix (strain ATCC 70633 / DSM 11037 / JCM 3026 / NBRC 100130 / K1)	AERPE	<a href="#">UP000002316</a>	1766	367	1399	20%
Buchnera aphidicola subsp. Acyrthosiphon pisum (strain APS) (Acyrthosiphon pisum symbiotic bacterium)	BUCAI	<a href="#">UP000001806</a>	572	572	0	69%
Archaeoglobus sulfaticallidus PM70-1	/	<a href="#">UP000013307</a>	2213	0	2213	18%
Archaeoglobus veneficus (strain DSM 11195 / SNP6)	ARCVS	<a href="#">UP000008136</a>	2065	0	2065	20%
archaeon GW2011_AR10	/	<a href="#">UP000031777</a>	1339	0	1339	14%
archaeon GW2011_AR15	/	<a href="#">UP000031776</a>	1308	0	1308	13%
archaeon GW2011_AR20	/	<a href="#">UP000031765</a>	1010	0	1010	11%
Caldisphaera lagunensis (strain DSM 15908 / JCM 11604 / IC-154)	CALLD	<a href="#">UP000010469</a>	1477	0	1477	17%

# Automatic annotation in UniProtKB

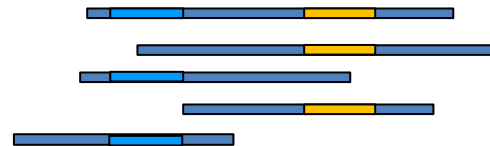
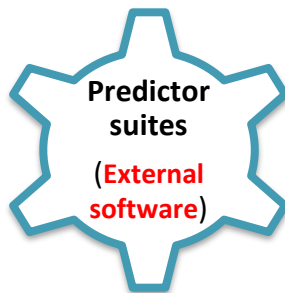
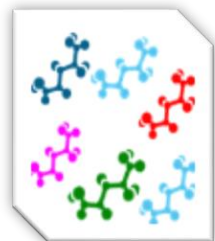


Rules are created by curators.  
(HAMAP, PIR name and site rules, and RuleBase rules)

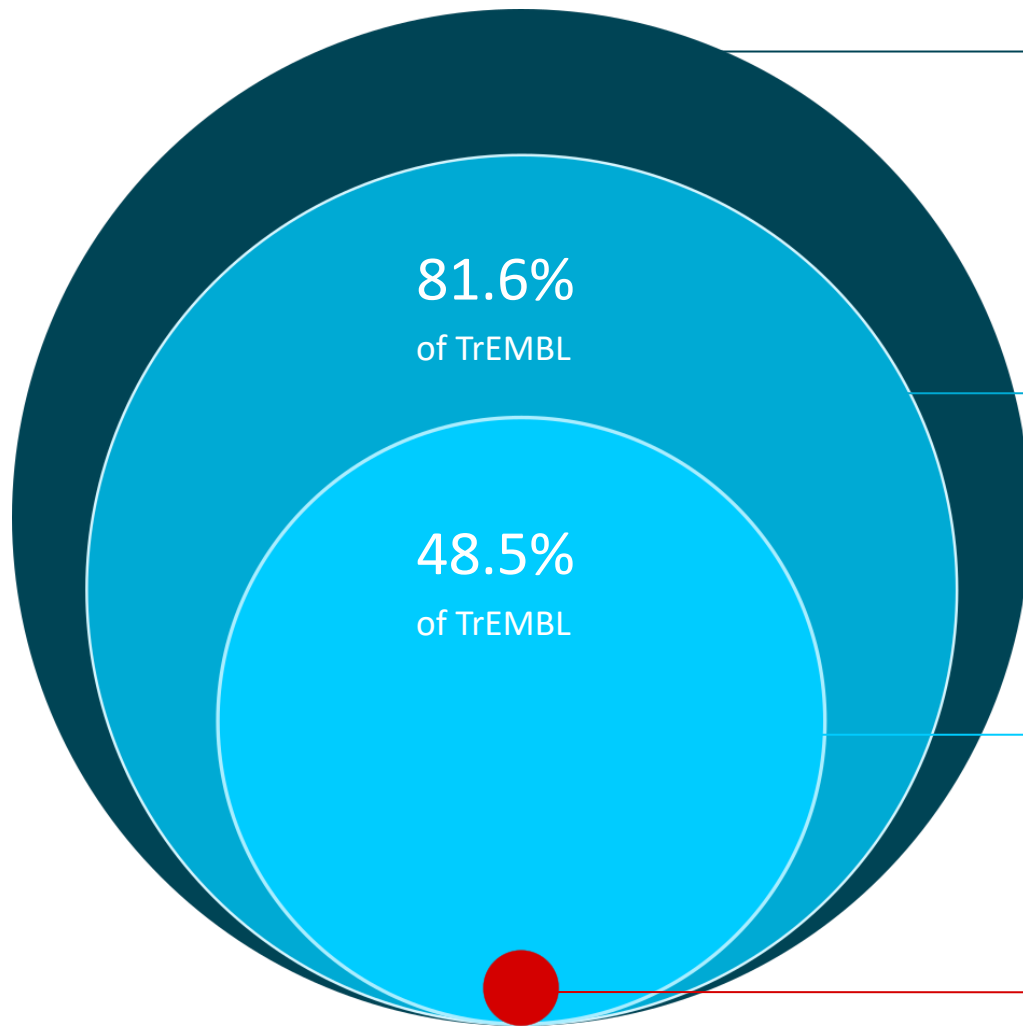


Rules are created by a modified **Decision Tree** algorithm.

## Additional information



**Sequence-specific information** from external providers such as InterPro or via SAM.



**214,406,868 TrEMBL**

Non-validated names from submitter;  
no further annotations

**174,854,319**

Sequence feature annotation  
signal peptides, transmembrane  
peptides, domains

**104,086,334 (UniRule + ARBA)**

Functional annotation  
protein names, enzymatic activity,  
interactors, subcellular location, GO terms  
etc.

**564,638 Swiss-Prot**

Release 2021\_02



# How to access automatic annotation data?

## Find your protein

UniProtKB + Advanced | List Search

UniProtKB  
UniProtKB & AlphaFold predictions  
UniRef  
UniParc  
Proteomes  
Taxonomy  
Keywords  
Literature citations  
Human diseases  
Cross-referenced databases  
Subcellular locations  
**UniRule**  
ARBA

0067, organism\_id:9606

high-quality, comprehensive and freely accessible resource of protein sequence and functional information. [Cite UniProt](#)

### Species Proteomes

Protein sets for species with sequenced genomes from across the tree of life

### Protein Clusters UniRef

Clusters of protein sequences at 100%, 90% & 50% identity

### Sequence Archive UniParc

Non-redundant archive of publicly available protein sequences seen across different databases

### AlphaFold structures

Search with all the power of the UniProt search engine for proteins with an AlphaFold prediction provided by DeepMind

### ProtNLM Predictions

Explore all the entries annotated with Google's ProtNLM predictions

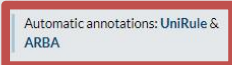


### UniProt COVID-19 portal

UniProt portal for the latest SARS-CoV-2 coronavirus protein entries and receptors, updated independent of the general UniProt release cycle

### Supporting Data

Human diseases | Cross-referenced databases | Subcellular locations | **Automatic annotations: UniRule & ARBA** | Literature Citations

Reviewed (Swiss-Prot) 568,363 | Unreviewed (TrEMBL) 229,928,140



# Using advanced search in UniRule

Advanced Search

Searching in  
UniRule

Protein Name [DE] Flavin prenyltransferase Remove

Gene Name [GN] ydj1 Remove

Protein Name [DE] 05067, cdc7 human Remove

Add Field

Search for field

- Enzyme classification [EC]
- Cofactors
- Catalytic activity
- Activity regulation
- Pathway
- Subcellular location
- Expression
- Family and Domains
- Gene Ontology [GO]
- Keyword [KW]

Cancel Search

Release 2022\_04 | Statistics Help

## Find your protein

Advanced | List | Search

comprehensive and freely accessible resource of protein sequence and functional information. [Cite UniProt](#)

### Species Proteomes

Protein sets for species with sequenced proteomes from across the tree of life

### Protein Clusters UniRef

Clusters of protein sequences at 100%, 90% & 50% identity

### Sequence Archive UniParc

Non-redundant archive of publicly available protein sequences seen across different databases

Feedback Help

# Using advanced search in UniRule

UniProt BLAST Align Peptide search ID mapping SPARQL UniRule (protein\_name:"Flavin prenyltransferase") Advanced | List Search

Superkingdom  
Bacteria (3)  
Archaea (1)  
Eukaryota (1)

## UniRule 3 results

Download View: Cards Table Customize columns Share

UniRule ID	Statistics	Taxonomic scope	Annotation covered	Predicted protein name	Template entries
<input type="checkbox"/> UR000195519	17,794 UniProtKB entries 30 reviewed UniProtKB entries 17,764 unreviewed UniProtKB entries	Bacteria Archaea	protein name catalytic activity function similarity keyword 2 more annotations	Flavin prenyltransferase UbiX, 2.5.1.129	Q9HX08 P0AG03
<input type="checkbox"/> UR000375818	642 UniProtKB entries 1 reviewed UniProtKB entry 641 unreviewed UniProtKB entries	Fungi (Excluding) Bacteria	protein name catalytic activity function subcellular location similarity 4 more annotations	Flavin prenyltransferase PAD1, mitochondrial, 2.5.1.129	P33751
<input type="checkbox"/> UR000195377	1,369 UniProtKB entries 6 reviewed UniProtKB entries 1,363 unreviewed UniProtKB entries	Bacteria	protein name catalytic activity function similarity subunit 2 more annotations	Probable UbiX-like flavin prenyltransferase, 2.5.1.129, Phenolic acid decarboxylase subunit B, PAD	P69772 P94404

Feedback Help

# Evaluating information in UniProtKB entries using its source

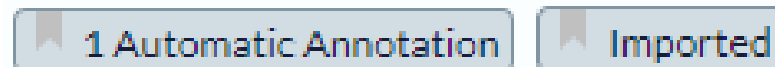
---

All information in an entry is attributed to its original source.

- manual assertions:



- automatic assertions:



# Filtering UniProtKB results for assertion method

**Advanced Search** ✕

Searching in UniProtKB

Function [CC] enzyme

Taxonomy [OC] human

Keyword [KW] chromosomal

All a4\_human, P05067, cdc7 human

**Evidence**

- Any assertion method
- Any manual assertion
- Any automatic assertion
- Any experimental assertion

**Manual assertions**

- Experimental
- Non-traceable author statement
- Curator inference
- Sequence similarity
- Sequence model
- Combinatorial
- Imported information




**Automatic assertions**

- Sequence model
- Combinatorial
- Imported information
- Sequence motif match (InterPro)

**Add Field**

Type \* in the search box to search for all values for the selected field.

**Search**

Release 2022\_04 | Statistics    Help

## our protein

Advanced | List **Search**

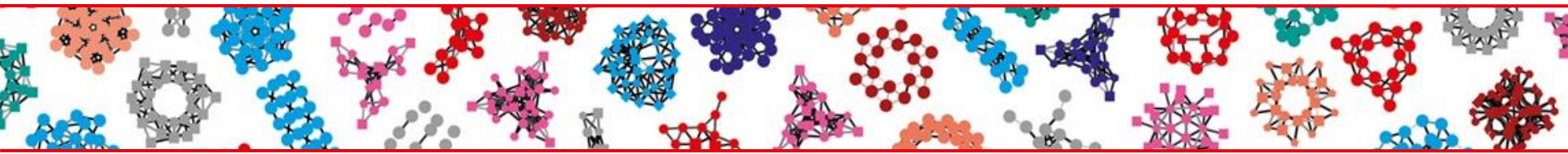
accessible resource of protein sequence and functional information. [Cite UniProt](#)

**Protein Clusters**  
UniRef

Clusters of protein sequences at 100%, 90% & 50% identity

**Sequence Archive**  
UniParc

Non-redundant archive of publicly available protein sequences seen across different databases



- Why do we need predictive annotation tools?
- Protein signatures for homology detection – A short primer
- Annotation rules for functional annotation
- HAMAP and PROSITE - automatic annotation in UniProtKB
- **HAMAP and PROSITE - services for external users**
- Practical exercises (afternoon)



## Database of protein domains, families and functional sites

PROSITE consists of documentation entries describing protein domains, families and functional sites as well as associated patterns and profiles to identify them [\[More... / References / Commercial users\]](#).

PROSITE is complemented by [ProRule](#), a collection of rules based on profiles and patterns, which increases the discriminatory power of profiles and patterns by providing additional information about functionally and/or structurally critical amino acids [\[More...\]](#).

**Release 20.113 of 26-Mar-2015 contains 1718 documentation entries, 1308 patterns, 1112 profiles and 1112 ProRule.**

### Search

e.g. PDOC00022, PS50089, SH3, zinc finger

### Browse

- [by documentation entry](#)
- [by ProRule description](#)
- [by taxonomic scope](#)
- [by number of positive hits](#)

### Quick Scan mode of ScanProsite

Quickly find matches of your protein sequences to PROSITE signatures (max. 10 sequences). [\[?\] Examples](#)

Enter UniProtKB accessions or identifiers or PDB identifiers or sequences in FASTA format

[Exclude motifs with a high probability of occurrence from the scan](#)

For more scanning options go to [ScanProsite](#)

### Other tools

- [PRATT](#) - allows to interactively generate conserved patterns from a series of unaligned proteins.
- [MyDomains - Image Creator](#) - allows to generate custom domain figures.





This form allows you to scan proteins for matches against the [PROSITE collection of motifs](#) as well as against your own patterns.

- Option 1 - Submit PROTEIN sequences to scan them against the PROSITE collection of motifs.**
- Option 2 - Submit MOTIFS to scan them against a PROTEIN sequence database.
- Option 3 - Submit PROTEIN sequences and MOTIFS to scan them against each other.

STEP 1 - Submit PROTEIN sequences [\[help\]](#)

- Submit PROTEIN sequences (max. 10) [Examples](#)
- Submit a PROTEIN database (max. 16MB) for repeated scans (The data will be stored on our server for 1 month).

Supported input:

- UniProtKB accessions e.g. [P98073](#) or identifiers e.g. [ENTK\\_HUMAN](#)
- PDB identifiers e.g. [4DGJ](#)
- Sequences in [FASTA format](#)

STEP 2 - Select options [\[help\]](#)

- Exclude motifs with a high probability of occurrence from the scan
- Exclude profiles from the scan
- Run the scan at high sensitivity (show weak matches for profiles)

STEP 3 - Select output options and submit your job

Output format:

Retrieve complete sequences:  If you choose this option, not all output formats are available.

---

Receive your results by email





# ScanProsite Results Viewer

Output format: Graphical view - this view shows ScanProsite results together with ProRule-based predicted intra-domain features [\[help\]](#).

## Hits for all PROSITE (release 20.121) motifs on sequence P98073 [UniProtKB/Swiss-Prot (release 2015\_12 of 09-Dec-15: 550116 entries)]:

found: 13 hits in 1 sequence

P98073 ENTK\_HUMAN (1019 aa)

RecName: Full=Enteropeptidase; EC=3.4.21.9; AltName: Full=Enterokinase; AltName: Full=Serine protease 7; AltName: Full=Transmembrane protease serine 15; Contains: RecName: Full=Enteropeptidase non-catalytic light chain; Full=Enteropeptidase catalytic light chain; Flags: Precursor; . *Homo sapiens (Human)*

```
MGSKRGISSRHHSLSSYEIMFAALFAILVVLCAGLIAVVSCLTIKESQRGAALGQSSHEARATFKITS
GVTYNPNLQDKLSDVDFKVLAFDLQQMIDEIFLSSNLKNEYKNSRVLQFENGSIIVVFDLFFAQWVS
DENVKEELIQGLEANKSSQLVTFHIDLNSVDILDKLTITSHLATPGNVSIECLPGSSPCTDALTCI
KADLFCDGVEVNCVPGSDSDNKMCAIVCDGRFLLTGSSGSFQATHYKPKSETSVVQCWIIRVNOGLS
IKLSFDDFNYYTDLIDYIEYGVGSSKILRASIWETNPGTIRIFSNQVITATFLIESDESDYVGFNAT
YTAFNSSSELNNYEKINCNEFDGFCFWVQDLNDDNEWERIQGSTIFSPFTGPNFDHTFGNASGFYIST
PTGPGGRQERVGLLSPLDPTLEPACLSFWYHMYGENVHKLSINISNDQNMKTVFQKEGNYGDNW
NYGQVTLNETVFKVAFNAFKNKILSDIALDDISLTYGICNGSLYPEPTLVPTPPPELPTDCGGPF
ELWEPNITFSSINFPNSYPNLAFVCVWILNAQKGNIQLFHQFDLENINDVVEIRDGEADSLLLA
VYTGPGFVKDVFSTINRMTVLLITNDVLRGGFKANFTTGYHLGIPPECKADHFQCKNGECVPLVN
LCDGHLHCEDGSDAECVRFNNGTNNNGLVFRFRIQSIWHTACAENWTQISNDVQCQLLGLGSGNS
SKPIFPDGGPFVKLNTAPDGLIILTPSQCLQDSLRLQCNHKSCKGLLAAQDITPKIVGGSNAK
EGAWFVWVGLYGGRLCCGASLVSSDWLVSAAHCVYGRNLEPSKWTAILGLHMKSNLTSFQIVPRL
IDEIVINPHYNRRRKNNDIAMMHLEFKVNYTDYIQPICLPEENQVFPGRNCSIAGWGTVVYQGIT
ANILQEAADVPLLSNERCQQMPEYNI TENMICAGYEEGGIDSCQDGGSGPLMCQENNRWFLAGVTS
FGYKCALPNRPGVYARVRSFTEWISFLH
```

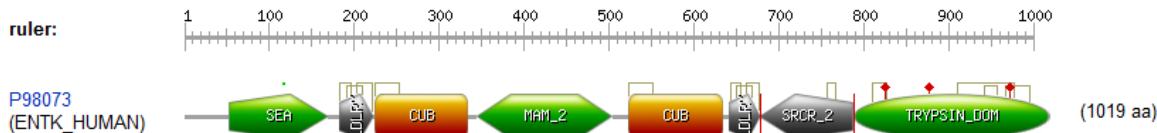
### Legend:



Please note that the graphical representations of domains displayed hereafter are for illustrative purposes only, and that their colors and shapes are not intended to indicate homology or shared function. For more information about how these graphical representations are constructed, go to <http://prosite.expasy.org/mydomains/>.

## hits by profiles: [8 hits (by 6 distinct profiles) on 1 sequence]

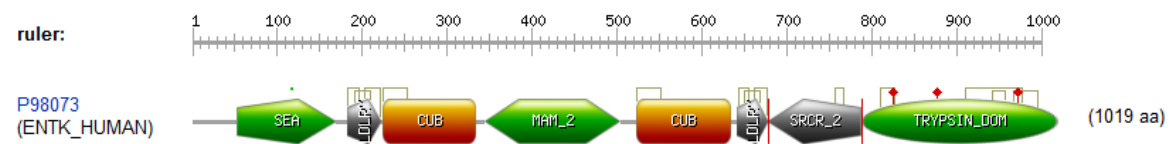
Upper case represents match positions, lower case insert positions, and the '-' symbol represents deletions relative to the matching profile.



RecName: Full=Enteropeptidase; EC=3.4.21.9; AltName: Full=Enterokinase; AltName: Full=Serine protease 7; AltName: Full=Transmembrane protease serine 15; Contains: RecName: Full=Enteropeptidase non-catalytic light chain; Full=Enteropeptidase catalytic light chain; Flags: Precursor; . *Homo sapiens (Human)*

**hits by profiles: [8 hits (by 6 distinct profiles) on 1 sequence]**

Upper case represents match positions, lower case insert positions, and the '.' symbol represents deletions relative to the matching profile.



RecName: Full=Enteropeptidase; EC=3.4.21.9; AltName: Full=Enterokinase; AltName: Full=Serine protease 7; AltName: Full=Transmembrane protease serine 15; Contains: RecName: Full=Enteropeptidase non-catalytic light chain; Full=Enteropeptidase catalytic light chain; Flags: Precursor; . *Homo sapiens (Human)*

**PS50024 SEA SEA domain profile :**

**54 - 169: score = 32.979**

QSHEARATFKITaGVTYNPNLQDKLSVDFKVLAFDLQQMIDEIFLSSNLKNEYKNSRVLQ  
 FEN--GSIIIVVFDLFFAQWVSD--ENVKEELIQGLEANKssQLVTFHIDLNSVDILDKLT

**Predicted features:**

DOMAIN	54	169	SEA	[condition: none]
SITE	117	118	Cleavage; by autolysis	[condition: none]

**PS50068 LDLRA\_2 LDL-receptor class A (LDLRA) domain profile :**

**183 - 222: score = 10.75**

E L P G S S P C T D a L T I K A D L F C D G E V N C P D G S D e d n R M C A

**Predicted features:**

DOMAIN	183	222	LDL-receptor class A	[condition: none]
DISULFID	184	197		[condition: C-x*-C]
DISULFID	191	210		[condition: C-x*-C]
DISULFID	204	221		[condition: C-x*-C]

**642 - 678: score = 13.3**

P C K A D H F Q C K N G E C V P L V N L C D G H L H C E D G S D E A D C V

**Predicted features:**

DOMAIN	642	678	LDL-receptor class A	[condition: none]
DISULFID	643	655		[condition: C-x*-C]
DISULFID	650	668		[condition: C-x*-C]
DISULFID	662	677		[condition: C-x*-C]

**PS01180 CUB CUB domain profile :**

**225 - 334: score = 13.293**

C D G R F L L I g s S G S F Q A T H Y P K --- P S E t a V V C Q W I I R V N Q G L -- S I K L S F D D - F N T Y ---  
 -- Y T D I L D I Y E G v g s S K I L R A S I W E T N P G T I R I F S N Q V T A T F L I E S D E S D y V G F N A T Y T A  
 F

**Predicted features:**

DISULFID	225	253		[condition: C-x*-C]
----------	-----	-----	--	---------------------



Search HAMAP [Search]

## High-quality Automated and Manual Annotation of Proteins

HAMAP is a system for the classification and annotation of protein sequences. It consists of a collection of manually curated family profiles for protein classification, and associated, manually created [annotation rules](#) that specify annotations that apply to family members. HAMAP is used to annotate protein records in UniProtKB via [UniProt's automatic annotation pipeline](#). We also provide an interface to scan user sequences against HAMAP family profiles [\[More...\]](#).

### Search

Field: All | Term: [ ] | Search | Reset

e.g. [recA](#), [MF\\_00191](#), [Iron](#)  
The wildcard \* is supported.

Last updated  
29-Apr-2015 for UniProt release 2015\_05

Number of family profiles: 2'015  
Number of annotation rules: 2'035

### HAMAP-Scan

Quickly find matches of your protein sequences to HAMAP family profiles (max. 1'000 sequences). [Examples](#)

Enter UniProtKB accessions or identifiers or sequences in FASTA format

Scan | Reset

### More options

If you want to submit more than 1'000 sequences, or if you want to have your sequences annotated in UniProtKB format by using HAMAP annotation rules, use the

[Advanced HAMAP-Scan submission form](#)

Questions? Comments? Please send us your [feedback](#).



## Browse annotation rules

Home | [Browse](#) | HAMAP-Scan | Proteomes | Documents | Downloads



The HAMAP annotation rules are created by expert curators. They are written in the [UniRule format](#) and contain data which is used to annotate bacterial, archaeal and eukaryotic proteins in the [UniProtKB format](#) [\[More...\]](#).

[List all rules](#) | **or Browse by:** [enzyme class](#) | [feature key](#) | [GO term](#) | [keyword](#) | [pathway](#) | [taxonomic scope](#)

Name	Gene name	Rule	Profile	Matches
16SrRNA_methyltr_A	ksmA / rsmA	MF_00607	MF_00607	8336
16SrRNA_methyltr_B	rsmB / sun	MF_01856	MF_01856	822
16SrRNA_methyltr_C	rsmC	MF_01862	MF_01862	1756
16SrRNA_methyltr_F	rsmF	MF_01579	MF_01579	962
16SrRNA_methyltr_G	rsmG	MF_00074	MF_00074	7913
16SrRNA_methyltr_H	rsmH	MF_01007	MF_01007	8454
16SrRNA_methyltr_I	rsmI	MF_01877	MF_01877	7638
16SrRNA_methyltr_J	rsmJ	MF_01523	MF_01523	2102
23SrRNA_methyltr_B	rlmB	MF_01887	MF_01887	2319
23SrRNA_methyltr_Cfr	cfr	MF_01873	MF_01873	78
23SrRNA_methyltr_F	rlmF	MF_01848	MF_01848	1479
23SrRNA_methyltr_G	rlmG	MF_01859	MF_01859	1342
23SrRNA_methyltr_H	rlmH	MF_00658	MF_00658	6693
23SrRNA_methyltr_I	rlmI	MF_01857	MF_01857	862
23SrRNA_methyltr_J	rlmJ	MF_00934	MF_00934	2872



## HAMAP-Scan

Classify (and annotate) your own protein sequences by using the collection of HAMAP family profiles and annotation rules [\[help\]](#).

=> [Retrieve results of a previous scan.](#)

### STEP 1 - Enter PROTEIN sequences [?]

- Enter PROTEIN sequences (max. 1'000) [Examples](#)
- or
- Upload a file (up to 100'000 sequences)

Enter UniProtKB accessions or identifiers or sequences in FASTA format

#### Supported input:

- UniProtKB accessions e.g. Q8ZHG0 or identifiers e.g. AGUA\_YERPE
- Amino acid sequences in [FASTA format](#)

### STEP 2 - Choose 'Scan' or 'Scan & Annotate' [?]

- 'Scan' - Find matches for your sequence(s) to HAMAP family profiles. You will get a list of all (trusted and weak) matches of your sequences along with their match score.
- or
- 'Scan & Annotate' - Have trusted matches of your sequence(s) to HAMAP family profiles annotated in UniProtKB format by the associated HAMAP annotation rules. To get the best annotation, all your sequences must originate from the same organism and you must indicate the corresponding taxonomic identifier (TaxID).

TaxID: \*  Enter the [TaxID](#) that represents the organism from which the sequences you're submitting originate from.

### STEP 3 - Submit your job [?]

Email: \*  You will be notified by email once the job you submitted has completed (even if no match has been found). Your results will be available on our server for 1 month.

Job title:  If available, the job title will be included in the subject of the email.

Password:  If you enter a password, then the same password will be requested before you can download your results.

\*Required field

# HAMAP-Scan: 'Scan' output



## HAMAP-Scan results

[Home](#) | [Search](#) | [HAMAP-Scan](#) | [Proteomes](#) | [Documents](#) | [Downloads](#)

Showing results for 10 sequences

Download: [\[ Tab-delimited \]](#) [\[ Excel \]](#) | Filter:  Trusted match  Weak match (below cutoff)  No match [Reload](#)

Your sequence [?]	Profile AC [?]	Profile name [?]	Trusted cutoff [?]	Match score [?]	Match region [?]	Match quality [?]
sp B1LQ15 SYA_ECOSM (876 aa)	<a href="#">MF_00036_B</a>	Ala_tRNA_synth_B	78.0	210.208	1-875	<b>Trusted</b>
sp B1LQ15 SYA_ECOSM (876 aa)	<a href="#">MF_03134</a>	Ala_tRNA_synth_plantC	29.724	20.217	2-876	Weak
sp B1LQ15 SYA_ECOSM (876 aa)	<a href="#">MF_00036_A</a>	Ala_tRNA_synth_A	18.492	10.689	1-868	Weak
sp B1LI51 ACP_ECOSM (78 aa)	<a href="#">MF_01217</a>	Acyl_carrier	13.362	22.327	3-77	<b>Trusted</b>
sp B1LGN8 AROE_ECOSM (272 aa)	<a href="#">MF_00222</a>	Shikimate_DH_AroE	26.377	33.106	1-269	<b>Trusted</b>
sp B1LGN8 AROE_ECOSM (272 aa)	<a href="#">MF_01578</a>	Shikimate_DH_YdiB	67.475	14.249	1-270	Weak
sp B1LI10 MNMA_ECOSM (368 aa)	<a href="#">MF_00144</a>	tRNA_thiouridyl_MnmA	31.753	40.978	6-361	<b>Trusted</b>
sp B1LI10 MNMA_ECOSM (368 aa)	<a href="#">MF_01633</a>	QueC	38.469	8.948	6-196	Weak
sp B1LF99 LSRB_ECOSM (340 aa)	-	-	-	-	-	No match
tr B1LFT9 B1LFT9_ECOSM (338 aa)	<a href="#">MF_00492</a>	Transaldolase_1	45.973	52.598	23-338	<b>Trusted</b>
tr B1LFT9 B1LFT9_ECOSM (338 aa)	<a href="#">MF_00496</a>	F6P_aldolase	48.944	10.611	33-299	Weak
tr B1LFT9 B1LFT9_ECOSM (338 aa)	<a href="#">MF_00494</a>	Transaldolase_3b	42.738	15.62	33-295	Weak
tr B1LDG3 B1LDG3_ECOSM (197 aa)	<a href="#">MF_01405</a>	Non_canon_purine_NTPase	36.059	51.874	2-196	<b>Trusted</b>
tr B1LDG3 B1LDG3_ECOSM (197 aa)	<a href="#">MF_03148</a>	HAM1_NTPase	49.594	21.502	2-196	Weak
tr B1LD69 B1LD69_ECOSM (375 aa)	<a href="#">MF_01899</a>	RNase_D	35.982	46.667	5-371	<b>Trusted</b>
tr B1LI94 B1LI94_ECOSM (367 aa)	<a href="#">MF_02121</a>	ASADH	31.464	34.304	1-367	<b>Trusted</b>
tr B1LKL8 B1LKL8_ECOSM (67 aa)	<a href="#">MF_00903</a>	TatE	24.392	27.09	1-67	<b>Trusted</b>
tr B1LKL8 B1LKL8_ECOSM (67 aa)	<a href="#">MF_00236</a>	TatA_E	14.87	21.147	1-62	<b>Trusted</b>

# HAMAP-Scan: 'Scan & Annotate'



**Request code:** LJJ

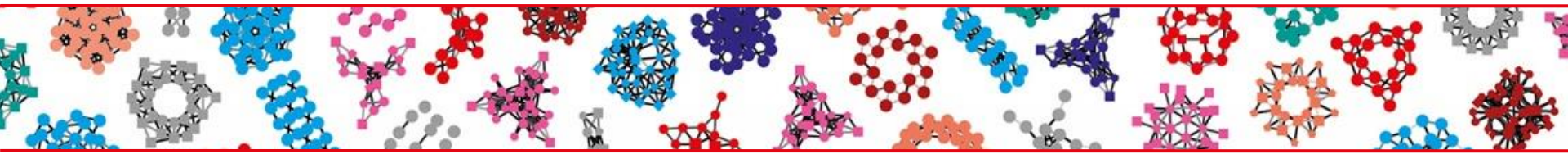
**Status:** Completed

**Input file:** 4305 entries

**Result files:** Annotated matches  
"Not annotated entries which app

HA

```
ID CYSC_ECOLX Unreviewed; 201 AA.
AC 3LJY468;
DE RecName: Full=Adenylyl-sulfate kinase;
DE EC=2.7.1.25;
DE AltName: Full=APS kinase;
DE AltName: Full=ATP adenosine-5'-phosphosulfate 3'-phosphotransferase;
DE AltName: Full=Adenosine-5'-phosphosulfate kinase;
GN Name=cysC;
OS Escherichia coli.
OC Bacteria; Proteobacteria; Gammaproteobacteria; Enterobacteriales;
OC Enterobacteriaceae; Escherichia.
OX NCBI_TaxID=562;
CC -!- FUNCTION: Catalyzes the synthesis of activated sulfate.
CC -!- CATALYTIC ACTIVITY: ATP + adenylyl sulfate = ADP + 3'-
CC phosphoadenylyl sulfate.
CC -!- PATHWAY: Sulfur metabolism; hydrogen sulfide biosynthesis; sulfite
CC from sulfate: step 2/3.
CC -!- SIMILARITY: Belongs to the APS kinase family.
DR UniPathway; UPA00140; UER00205.
DR GO; GO:0004020; F:adenylylsulfate kinase activity; IEA:HAMAP.
DR GO; GO:0005524; F:ATP binding; IEA:HAMAP.
DR GO; GO:0000103; P:sulfate assimilation; IEA:HAMAP.
DR HAMAP; MF_00065; Adenylyl_sulf_kinase; 1.
KW ATP-binding; Kinase; Nucleotide-binding; Phosphoprotein; Transferase.
FT NP_BIND 35 42 ATP.
FT ACT_SITE 109 109 Phosphoserine intermediate.
**
** ##### INTERNAL SECTION #####
**HA FAM; Method MF_00065; CYSC; Trusted match; 38.735 (+7.3).
**HA SAM; Annotated by praise 1.9.4; MF_00065.30; MF_00065; 06-JAN-2014 11:07:03.
SQ SEQUENCE 201 AA; 22321 MW; 11E15BB8F9D2FD4B CRC64;
MALHDENVVW HSHPVTVQQR ELHHGHRGVV LWFTGLSGSG KSTVAGALEE ALHKLGVSTY
LLDGDNVVRHG LCSDLGFSDA DRKENIRRVG EVANLMVEAG LVVLTAFISP HRAERQMVRE
RVGEGRFIEV FVDTPLAICE ARDPKGLYKK ARAGELRNFT GIDSVYEAPE SAEIHLNGEQ
LVTNLVQQLL DLLRQNDIIR S
//
```



- Why do we need predictive annotation tools?
- Protein signatures for homology detection – A short primer
- Annotation rules for functional annotation
- HAMAP and PROSITE - automatic annotation in UniProtKB
- HAMAP and PROSITE - services for external users
- **Practical exercises**