**Supplemental Figure 1. Fold Recognition Performance of Ada-BLAST with Different Settings Given Fold-specific Libraries.** (a) Relative frequency of # of PSSMs generated from each TZ-SABmark queries. The average number of PSSMs generated from the queries is 3.6, with the most being 31 PSSMs which were generated from *d1tfra1,* a *SAM-domain like* fold. If a TZ-SABmark query was not in the version of PDB which we used to construct fold-specific libraries, PSSMs were not generated from the query. There are 60 such TZ-SABmark queries. (b) Comparison of ROC curves of Ada-BLAST at different coverage thresholds when e-value threshold is fixed at $10^{10}$. (c) Comparison of ROC curves of Ada-BLAST of e-value 0.01 & no coverage, e-value 0.01 & 80% coverage, e-value $10^{10}$ & no coverage, and e-value $10^{10}$ & 80% coverage.

**Supplemental Figure 2. Characterization of alignments used by Ada-BLAST at e-value 0.01 and $10^{10}$ thresholds.** Pairwise identity of the alignments collected by Ada-BLAST using rps-BLAST with e-value 0.01, no coverage and e-value $10^{10}$, 80% coverage. The analysis is done separately for the alignments between queries and the PSSMs of true fold-specific library and those between the queries and the PSSMs of false fold-specific library.

**Supplemental Figure 3. Comparison of Ada-BLAST dendrograms of e-value 0.01 and $10^{10}$ thresholds.** A portion of the e-value 0.01, no coverage dendrogram containing the queries mis-clustered (*top left box)* and a portion of the e-value $10^{10}$, 80% coverage dendrogram containing the same queries. The queries in blue boxes (queries in red in the dendrogram of e-value 0.01, no coverage in Fig. 2b) were improperly clustered in the result of e-value 0.01, no coverage, but properly clustered in the result of e-value $10^{10}$, 80% coverage.

**Supplemental Figure 4-5. Hierarchical clustering of transmembrane containing proteins with additional Ada-BLAST settings.** 74 sequences representing multiple classes of transmembrane containing proteins were hierarchically clustered and visualized by Cluster and Treeview [21]. The horizontal lines represent the correlation
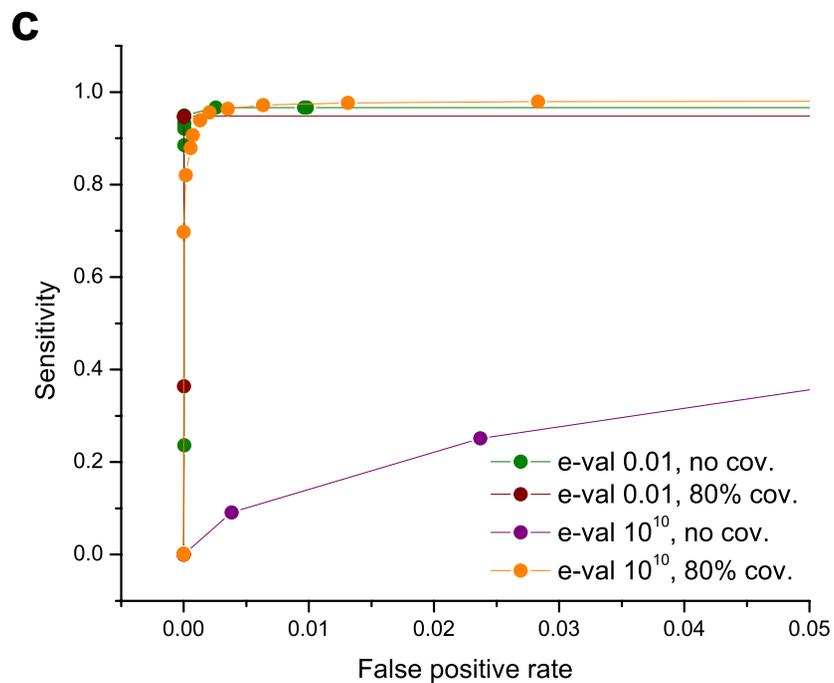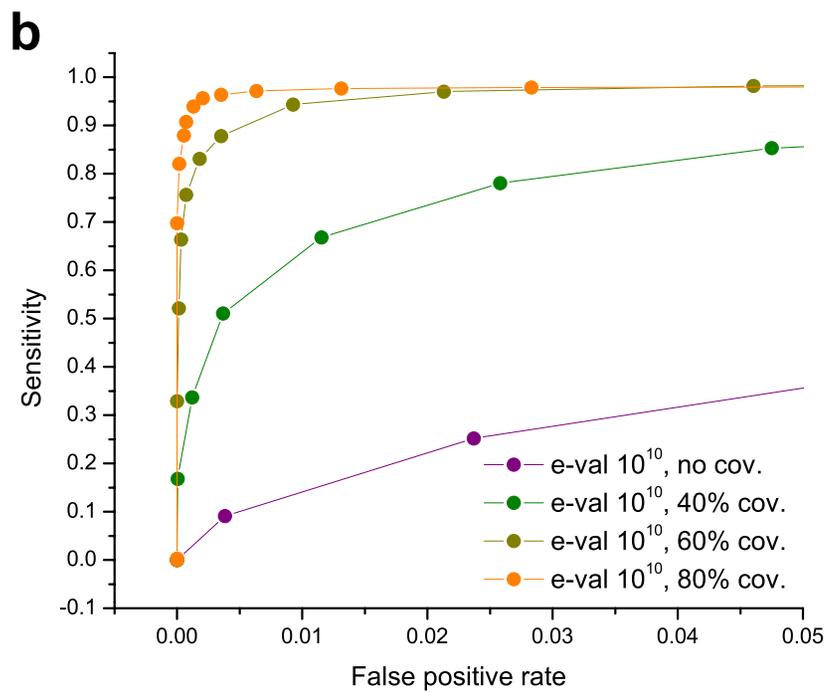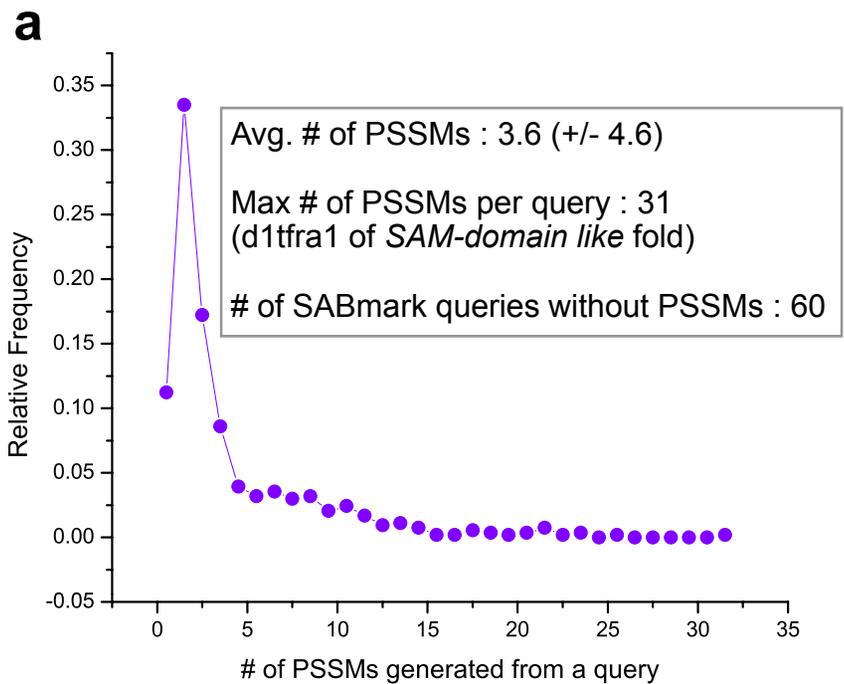
scores derived from the analysis. Alignments for the ILB DB PSSMs were derived from embedded alignments (Fig. S4) and e=0.01 (Fig. S5) threshold.
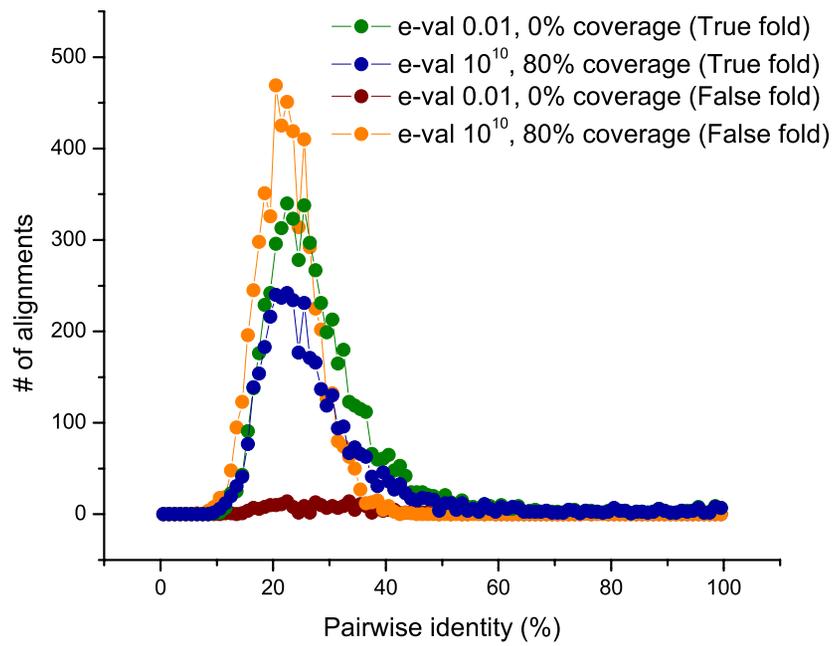
**Supplemental Figure 6. rps-BLAST Analysis of 1F88.** The primary amino acid sequence for 1F88 (gi|197107531) was screened by rps-BLAST at e-value=1. This figure shows the web-based output for domain identification.

**Supplemental Figure 7. The Characterization of Structural Elements in 1F88.** This graph shows the performance of Hidden Markov Models (TMHMM) vs embedded Ada-BLAST in determining the membrane spanning domains in Bovine Rhodopsin as determined by X-ray Crystallography (green= Beta pleated sheets, red=helices, loops not shown). This protein was analyzed with an expanded set of PSSMs representing a large variety of transmembrane domains (~30K PSSMs). It is resonable to consider that the amino-acids within transmembrane spanning helicies will be more conserved than the intervening loop residues. The support of this hypothesis is presented herein. The structural features are annotated with droplines. The transmembrane probability determined by TMHMM is shown in teal. The other graph depicts a curve-fitted positional score for embedded Ada-BLAST (see Figure 4 for raw and smoothed data). The positional score was quantified in the following manner. For each positive PSSM, the alignment boundaries are determined by the overlapping alignments obtained from Ada-BLAST. These regions were extracted and realigned by the Smith-Waterman algorithm with a BLOSUM62 substitution matrix. Using the alignments, each residue was scored with substitution scores of BLOSUM62 if the residue is identically or positively (non-identical but conserved) aligned. This process was repeated for all positive PSSMs and the results were summed for each amino acid in the protein. The positional results were normalized to zero by subtracting the average positional score across the protein length from each point, and each amino acid position was then subjected to smoothing (Fast Fourier-transform point=8) and discontinous baselining using Origin Lab 7.5© . Baseline correction was performed by baselining to every local minimum across the entire curve.
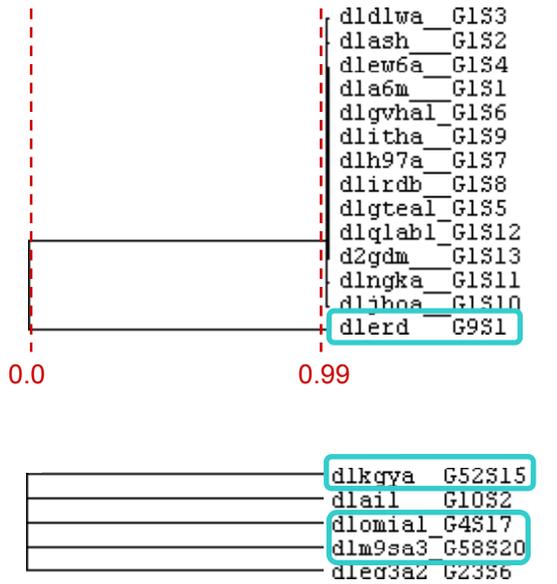
**Supplemental Figure 8. Ada-BLAST comparison with TMHMM and rps-BLAST for TRPC3.** The top graphic depicts the domain organization predicted by rps-BLAST. Below this is a graph containing TM probabilities predicted by hidden-markov model TMHMM (magenta, right y-axis) compared with postitional Ada-BLAST data from embedded alignments with >30K TM PSSMs.

**Supplemental Figure 9. Ada-BLAST comparison with TMHMM and rps-BLAST for TRPV5.** The top graphic depicts the domain organization predicted by rps-BLAST. Note that rps-BLAST does not show a domain for the known channel region. Below this is a graph containing TM probabilities predicted by hidden-markov model TMHMM (magenta, right y-axis) compared with postitional Ada-BLAST data from embedded alignments with >30K TM PSSMs.

**a**

Avg. # of PSSMs : 3.6 (+/- 4.6)

Max # of PSSMs per query : 31
(d1tfra1 of *SAM-domain like* fold)

# of SABmark queries without PSSMs : 60

**b**

- e-val $10^{10}$, no cov.
- e-val $10^{10}$, 40% cov.
- e-val $10^{10}$, 60% cov.
- e-val $10^{10}$, 80% cov.

**c**

- e-val 0.01, no cov.
- e-val 0.01, 80% cov.
- e-val $10^{10}$, no cov.
- e-val $10^{10}$, 80% cov.

**e-val 0.01, no cov.**

dldlwa__G1S3
dlash___G1S2
dlew6a__G1S4
dla6m___G1S1
dlgvha1_G1S6
dlitha__G1S9
dlh97a__G1S7
dlirdb__G1S8
dlgtea1_G1S5
dlqlab1_G1S12
d2gdm___G1S13
dlngka__G1S11
dljhoa__G1S10
dlerd___G9S1

0.0    0.99

dlkgya__G52S15
dlail___G10S2
dlomia1_G4S17
dlm9sa3_G58S20
dleg3a2_G23S6

**e-val $10^{10}$, 80% cov.**

dlerd___G9S1
dlhd6a__G9S2
d2er1___G9S3
1.0

dlgwma__G52S7
dlgmya__G52S4
dli5pa1_G52S9
dlbhga2_G52S1
dld7pm__G52S3
dlof4a__G52S18
dlgu3a__G52S5
dlgu1a__G52S6
dlk12a__G52S13
dlju3a1_G52S11
dljhja__G52S10
dlh6ya__G52S8
dlk3ia2_G52S14
dljz8a3_G52S12
dlciy_1_G52S2
dlot5a1_G52S19
dlxnaa__G52S20
dl17la__G52S16
dlf1ca1_G53S2
dlkgya__G52S15
dllnsa2_G52S17

0.53    0.75    0.94

dlg2ba__G58S7
dlgcqc__G58S8
dlm9sa3_G58S20
dlbf4a__G58S2
dlh3za__G58S9
dlfx7a3_G58S6
dlilja__G58S12
dlbb9___G58S1
dli07a__G58S11
dle0ba__G58S5
dljb0e__G58S14
dldzla__G58S4
dlvie___G58S23
dlpht___G58S22
dlixda__G58S13
dllpla__G58S18
dloila1_G58S21
dlm9sa2_G58S19
d2ahjb__G58S24
dlkhia1_G58S17

0.97        0.99

dlomia1_G4S17
dlrsha__G4S8
dlbjaa__G4S1
dlmd0a__G4S15
dlbm9a__G4S2
dlldja1_G4S13
dlka8a__G4S12
dlrepc2_G4S21
dlfc3a__G4S5
dljhga__G4S11
dlopc___G4S18
dld5va__G4S3
dlfp2a1_G4S7
dlfpld1_G4S6
dlf1za1_G4S4
dlig6a__G4S9
dligna1_G4S10
dlp4xa1_G4S19
dllj9a__G4S14
d2irfg__G4S25
dlsmta__G4S22
d2ez1___G4S23
dlp4xa2_G4S20
d2hts___G4S24
dlo7fa1_G4S16
dlh6ka1_G42S6
d2cb1a1_G23S16
dleg3a2_G23S6
dlg6ea__G51S4
dlail___G10S2

0.51    0.71    0.92

**b**

NP_003296.1 ⎱ TRPC
NP_065122.1 ⎰
NP_671737.1 TRPV4
NP_004612.2 TRPC6
NP_542435.2 ⎱
NP_057197.2 ⎱
NP_659505.1 ⎱ TRPV
NP_062815.2 ⎱
NP_061116.2 ⎰
NP_015628.2 TRPA1
NP_036603.1 TRPC5
NP_003295.1 TRPC1
NP_079247.5 TRPM3
NP_060132.3 TRPM6
NP_057263.1 TRPC4
NP_000288.1 TRPP2
NP_003298.1 TRPM2
NP_057196.2 ⎱ TRPP
NP_055201.2 ⎰
NP_076985.4 ⎱
NP_002411.3 ⎱ TRPM
NP_055370.1 ⎱
NP_060142.3 ⎰
NP_065394.1 ⎱
NP_694991.2 ⎱ TRPML
NP_060768.8 ⎰
NP_000714.3 Voltage-gated calcium channel
NP_000715.2 Voltage-gated calcium channel
NP_001005746.1 Voltage-gated calcium channel
NP_990562.1 Neurotrophic tyrosine kinase
NP_001690.2 AXL receptor tyrosine kinase
NP_000716.2 Voltage-gated calcium channel
NP_001013165.1 AXL receptor tyrosine kinase
NP_001101233.2 Leukocyte receptor tyrosine kinase
NP_996844.1 Leukocyte receptor tyrosine kinase
NP_996824.1 Leukocyte receptor tyrosine kinase
NP_062633.1 ⎱
NP_000211.1 ⎱ Potassium inwardly-recifying channel
NP_722450.1 ⎰
XP_001054644.1 Potassium channel
NP_852006.1 ⎱
NP_001157149.1 ⎱ Calcium-activated potassium channel
NP_055320.4 ⎰
NP_060106.2 TRPM4
XP_001651299.1 Tyrosine kinase receptor
NP_974783.1 ⎱ Cyclic-nucleotide-gated cation channel
XP_538462.2 ⎰
YP_039438.1 ⎱ Acetamide transporter
YP_002532980.1 ⎰
XP_592255.3 ⎱ Potassium channel
XP_001926284.1 ⎰
YP_003058108.1 Phosphate transporter
NP_175261.1 ⎱
NP_598537.2 ⎱
NP_543141.3 ⎱
NP_570822.1 ⎱
NP_665730.2 ⎱ G protein receptor
NP_071715.3 ⎱
NP_722561.1 ⎱
NP_001007200.2 ⎱
NP_001076117.1 ⎰
ZP_00518439.1 Calcium/proton exchanger
NP_988188.1 ⎱ Sodium/calcium exchanger
NP_150287.1 ⎰
XP_538675.2 ⎱ Sodium/potassium/calcium exchanger
XP_001926832.1 ⎰
XP_001950220.1 ⎱
XP_001931622.1 ⎱ Sodium/hydrogen exchanger
XP_545197.2 ⎱
XP_002291556.1 ⎰
NP_198067.1 Sodium ion transmembrane transporter
NP_001031990.1 ⎱ Voltage-gated chloride channel
NP_198905.1 ⎰
YP_003057227.1 Amino-acid transporter

**Ada-BLAST (rps-BLAST alignments)**
ILB DB (38,155 PSSMs)
74 queries
e-value threshold: 0.01
Score: avg. identity x max. coverage
Hierarchical clustering
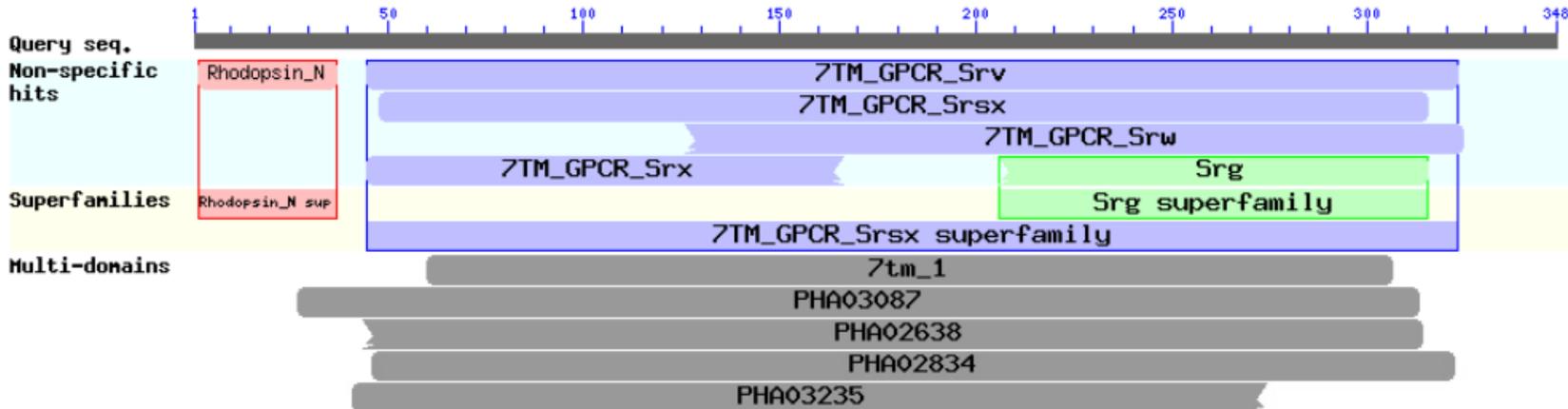Pearson's correlation

0.24        0.7   0.95

# 1F88

## Conserved domains on [lcl|seqsig_9090fcc463d253ce23ada88e98fdc7e7]

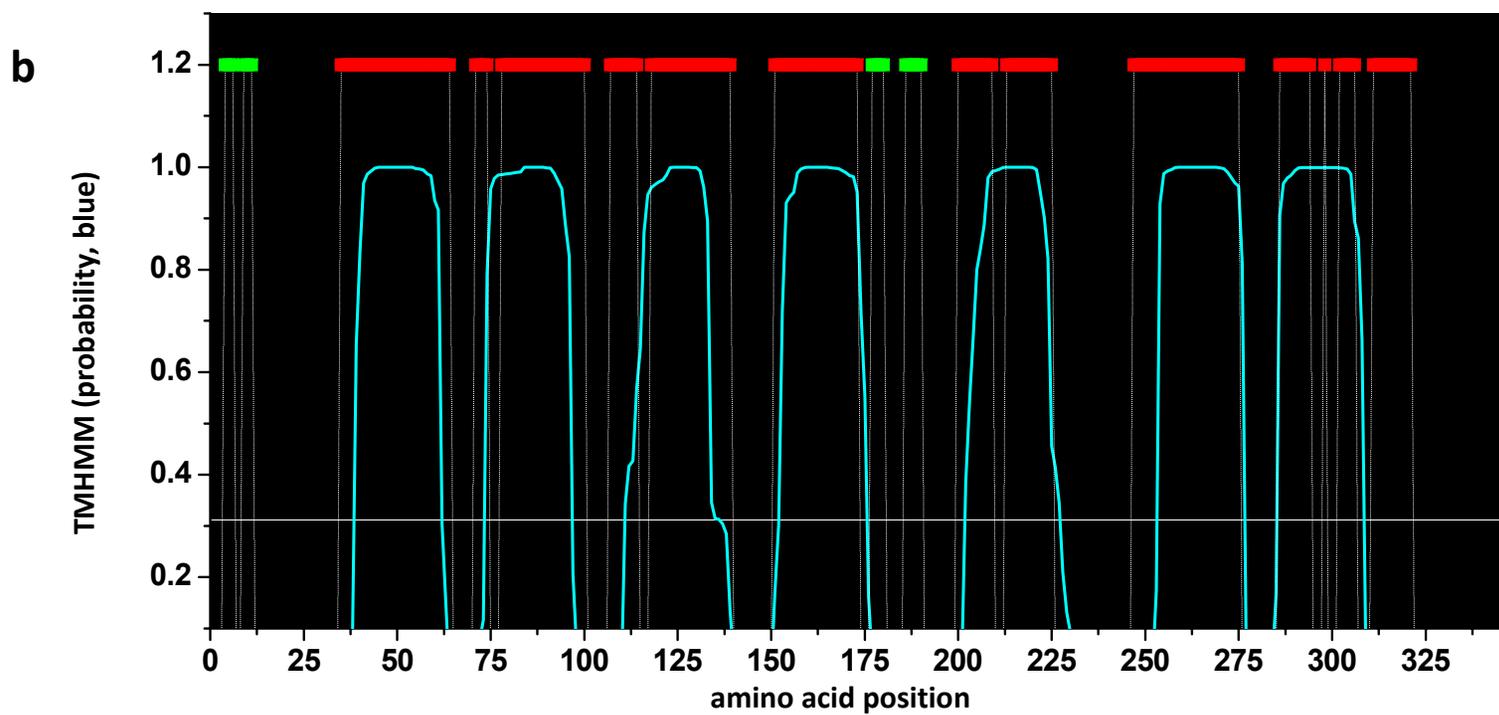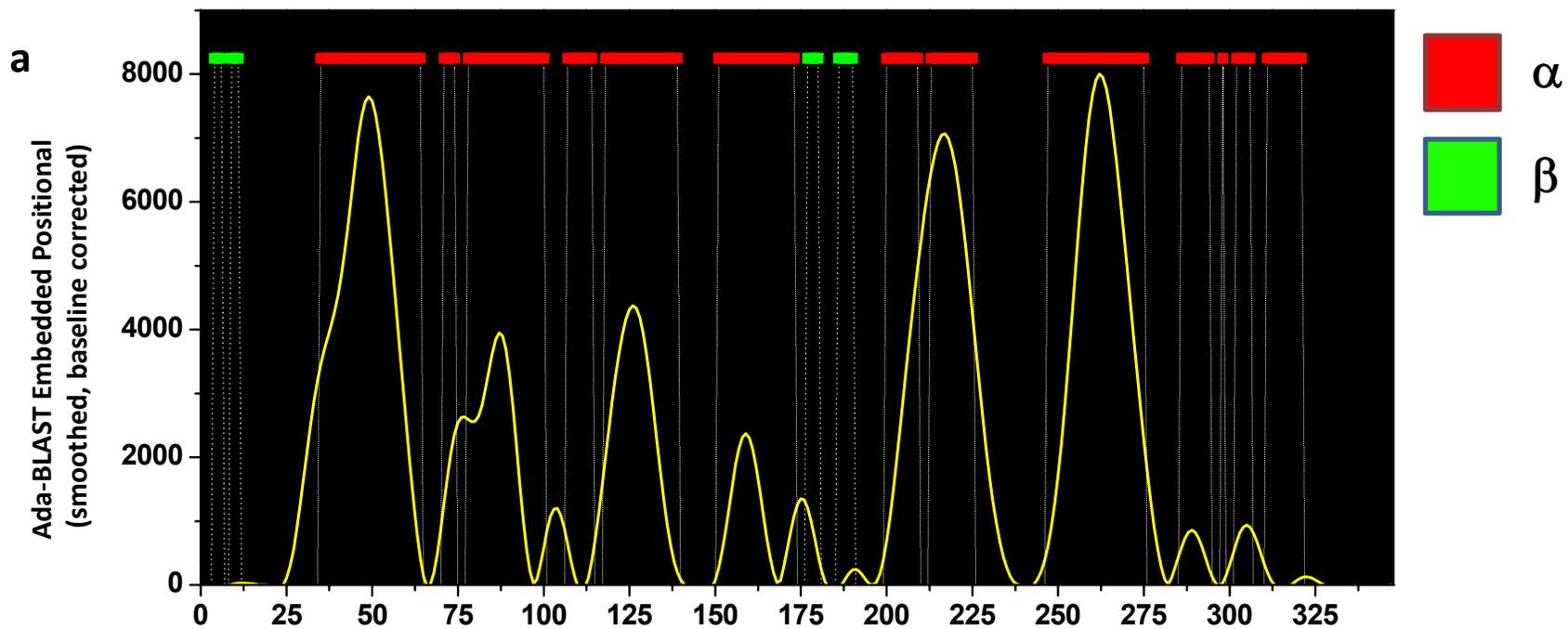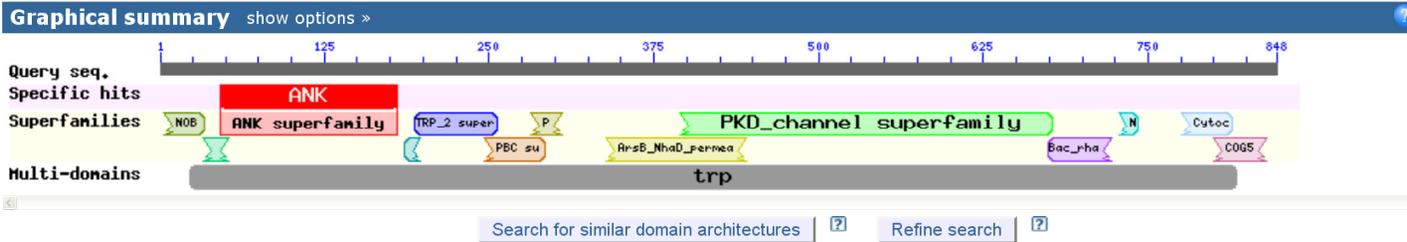Local query sequence

### Graphical summary   show options »



### List of domain hits

| Description | PssmId | Multi-dom | E-value |
|---|---|---|---|
| Rhodopsin_N[pfam10413], Rhodopsin is the archetypal G-protein-coupled receptor. Such receptors participate in... | 150994 | no | 2.65e-16 |
| 7TM_GPCR_Srv[pfam10323], Chemoreception is mediated in Caenorhabditis elegans by members of the seven-... | 150919 | no | 7.96e-06 |
| 7TM_GPCR_Srsx[pfam10320], Chemoreception is mediated in Caenorhabditis elegans by members of the seven-... | 150916 | no | 6.83e-03 |
| 7TM_GPCR_Srw[pfam10324], Chemoreception is mediated in Caenorhabditis elegans by members of the seven-... | 150920 | no | 0.03 |
| 7TM_GPCR_Srx[pfam10328], Chemoreception is mediated in Caenorhabditis elegans by members of the seven-... | 150924 | no | 0.15 |
| Srg[pfam02118], Srg family chemoreceptor. | 145331 | no | 0.89 |
| 7tm_1[pfam00001], This family contains, amongst other G-protein-coupled receptors (GCPRs), members of the... | 143794 | yes | 8.55e-42 |
| PHA03087[PHA03087], G protein-coupled chemokine receptor-like protein; Provisional | 165371 | yes | 2.10e-11 |
| PHA02638[PHA02638], CC chemokine receptor-like protein; Provisional | 165021 | yes | 2.67e-09 |
| PHA02834[PHA02834], chemokine receptor-like protein; Provisional | 165177 | yes | 3.02e-04 |
| PHA03235[PHA03235], DNA packaging protein UL33; Provisional | 165494 | yes | 2.58e-03 |

rps-BLAST
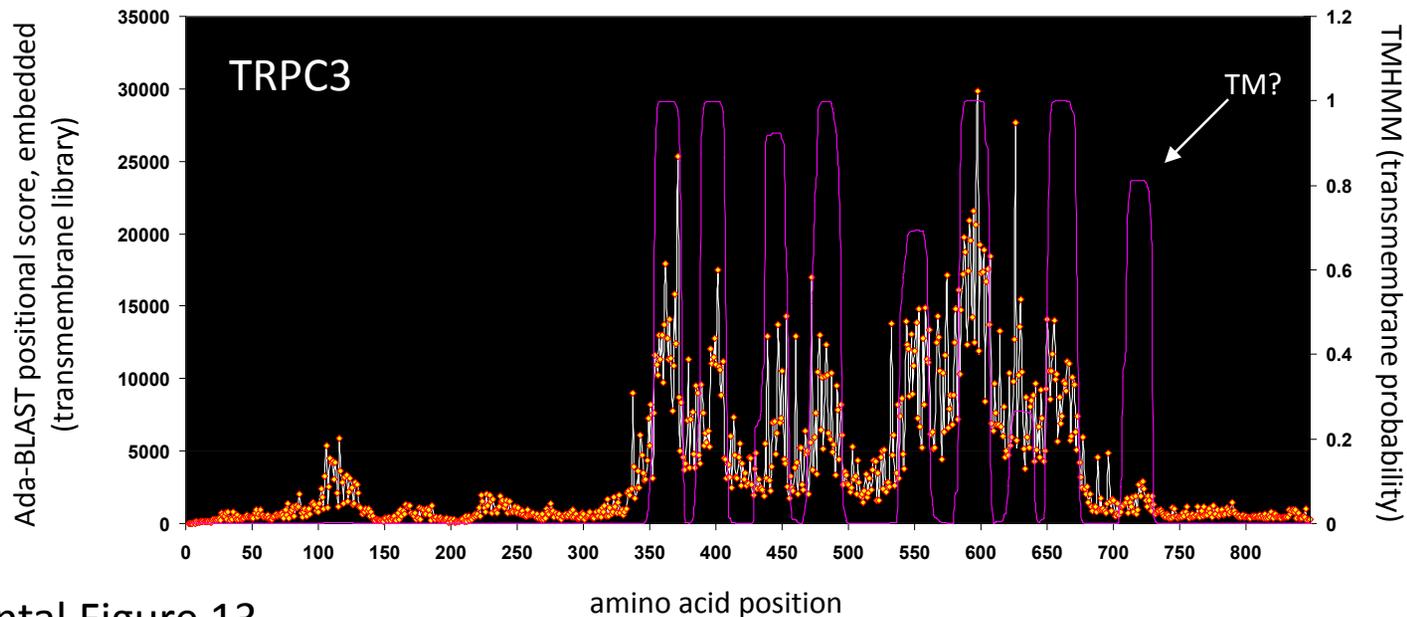
## Conserved domains on [gi|4507687|ref|NP_003296.1|]

View full result

short transient receptor potential channel 3 isoform b [Homo sapiens]

**Graphical summary** show options »



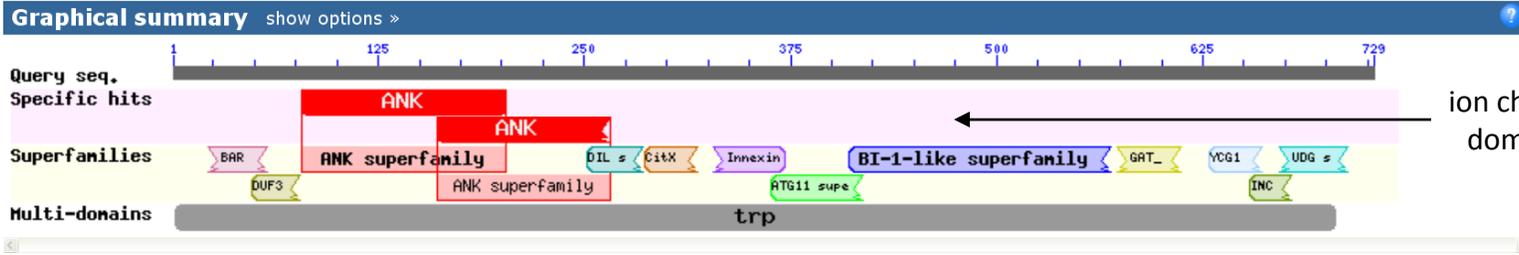Search for similar domain architectures | Refine search

**List of domain hits**

| Description | Pssmid | Multi-dom | E-value |
|---|---|---|---|
| ANK[cd00204], ankyrin repeats; ankyrin repeats mediate protein-protein interactions in very diverse... | 29261 | yes | 1.07e-10 |
| TRP_2 super family[cl07091], This domain is found in the transient receptor ion channel (Trp) family of proteins.... | 149414 | no | 3.93e-27 |
| PKD_channel super family[cl10691], This family contains the cation channel region of PKD1 and PKD2 proteins. | 116626 | no | 7.63e-14 |
| PBC super family[cl04264], The PBC domain is a member of the TALE (three-amino-acid loop extension) superclass of... | 146435 | no | 1.6 |
| ArsB_NhaD_permease super family[cl09110], Anion permease ArsB/NhaD. These permeases have been shown to translocate sodium,... | 175004 | no | 3.7 |
| Cytochrom_B562 super family[cl01546], This family contains the bacterial cytochrome b562. This forms a four-helix bundle that... | 174645 | no | 4.2 |
| NOB1_Zn_bind super family[cl10704], This domain corresponds to a zinc ribbon and is found on the RNA binding protein NOB1. | 149739 | no | 7.2 |
| COG5381 super family[cl02225], Uncharacterized protein conserved in bacteria [Function unknown] | 34944 | no | 11 |
| NR_DBD_like super family[cl02596], DNA-binding domain of nuclear receptors is composed of two C4-type zinc fingers. Each... | 155002 | no | 22 |
| Peptidase_M14_like super family[cl11393], The M14 family of metallocarboxypeptidases (MCPs) are zinc-binding carboxypeptidases (... | 175235 | no | 23 |
| Bac_rhamnosid super family[cl01801], This family consists of bacterial rhamnosidase A and B enzymes. L-Rhamnose is abundant... | 174672 | no | 27 |
| TT_ORF2 super family[cl03800], TT virus (TTV), isolated initially from a Japanese patient with hepatitis of unknown... | 145881 | no | 29 |
| Peptidase_S48 super family[cl11616], Peptidase family S48. | 159582 | no | 49 |
| trp[TIGR00870], after chronic exposure to capsaicin. (McCleskey and Gold, 1999). | 162078 | yes | 0e+00 |



Supplemental Figure 13

rps-BLAST

## Conserved domains on [gi|17505200|ref|NP_062815.2|]

transient receptor potential cation channel subfamily V member 5 [Homo sapiens]

**Graphical summary**  show options »

Query seq.
Specific hits
ANK
ANK

Superfamilies
BAR   ANK superfamily   ANK   DIL s  CitX   Innexin   BI-1-like superfamily   GAT_   YCG1   UDG s
DUF3   ANK superfamily   ATG11 supe   INC

Multi-domains
trp

ion channel domain?

Search for similar domain architectures     Refine search

**List of domain hits**

| Description | Pssmid | Multi-dom | E-value |
|---|---|---|---|
| ANK[cd00204], ankyrin repeats; ankyrin repeats mediate protein-protein interactions in very diverse... | 29261 | yes | 1.27e-15 |
| ANK[cd00204], ankyrin repeats; ankyrin repeats mediate protein-protein interactions in very diverse... | 29261 | yes | 1.54e-08 |
| BI-1-like super family[cl00473], BAX inhibitor (BI)-1 like protein family. Mammalian members of this family of small... | 174226 | no | 0.08 |
| ATG11 super family[cl11043], This is a family of proteins involved in telomere maintenance. In Schizosaccharomyces... | 150965 | no | 3.2 |
| CitX super family[cl01498], | 174633 | no | 4.9 |
| GAT_1 super family[cl00020], Type 1 glutamine amidotransferase (GATase1)-like domain. This group contains proteins... | 173987 | no | 6.4 |
| YCG1 super family[cl09228], Chromosome condensation complex Condensin, subunit G [Chromatin structure and dynamics /.. | 34815 | no | 7.5 |
| INCENP_ARK-bind super family[cl04337], This region of the inner centromere protein has been found to be necessary and... | 146526 | no | 12 |
| BAR super family[cl12013], BAR domains are dimerization, lipid binding and curvature sensing modules found in many... | 159673 | no | 14 |
| UDG super family[cl00483], | 174234 | no | 28 |
| DUF3603 super family[cl13637], This protein is found in bacteria and eukaryotes. Proteins in this family are about 250... | 152662 | no | 34 |
| Innexin super family[cl03000], This family includes the drosophila proteins Ogre and shaking-B, and the C. elegans... | 174746 | no | 45 |
| DIL super family[cl03379], The DIL domain has no known function. | 145158 | no | 62 |
| trp[TIGR00870], after chronic exposure to capsaicin. (McCleskey and Gold, 1999). | 162078 | yes | 0e+00 |

TRPV5

Ada-BLAST positional score, embedded (transmembrane library)

TMHMM (transmembrane probability)

amino acid position