# Gathering Sequences: BLAST

Searching............................done



Sequences producing significant alignments:

| | | Score (bits) | Value |
|---|---|---|---|
| sp|P80050| PRVA_MACFU (PVALB) PARVALBUMIN ALPHA. [PARVALBUMIN, MUS... | | 187 | 7e-48 |
| sp|P20472| PRVA_HUMAN (PVALB) PARVALBUMIN ALPHA. [Homo sapiens] | | 181 | 3e-46 |
| sp|P02625| PRVA_RAT (PVALB.) PARVALBUMIN ALPHA. [Rattus norvegicus] | | 175 | 2e-44 |
| sp|P80080| PRVA_GERSP PARVALBUMIN ALPHA. [Gerbillus sp.] | | 173 | 1e-43 |
| sp|P02624| PRVA_MOUSE (PVALB) PARVALBUMIN ALPHA. [Oryctolagus cuni... | | 172 | 2e-43 |
| sp|P32848| PRVA_MOUSE (PVALB.) PARVALBUMIN ALPHA. [Mus musculus] | | 168 | 3e-42 |
| sp|P80079| PRVA_FELCA (PVALB) PARVALBUMIN ALPHA. [Felis silvestris... | | 163 | 8e-41 |
| sp|P80026| PRVM_CHICK PARVALBUMIN, MUSCLE. [Gallus gallus] | | 154 | 5e-38 |
| sp|P51434| PRVA_CAVPO (PVALB) PARVALBUMIN (FRAGMENT). [Cavia ... | | 154 | 5e-38 |
| sp|P18087| PRVA_RANCA PARVALBUMIN ALPHA. (PA 4.97). [Rana catesbei... | | 123 | 1e-28 |
| sp|P02627| PRVA_RANES PARVALBUMIN ALPHA. [Rana esculenta] | | 121 | 4e-28 |
| sp|P02626| PRVA_AMPME PARVALBUMIN ALPHA. [Amphiuma means] | | 111 | 4e-25 |
| sp|P02614| PRVB_GRAGE PARVALBUMIN BETA. [Esox lucius] | | 111 | 5e-25 |
| sp|P30563| PRVA_TRISE PARVALBUMIN ALPHA. [Triakis semifasciata] | | 107 | 6e-24 |
| sp|P02594| MYC_PATSE MYOSIN REGULATORY LIGHT CHAIN 2 SKELETED YD... | | 32 | 0.038 |
| sp|P02623| MYB_CHMYI MYOSIN REGULATORY LIGHT CHAIN 2 SKELETED YD... | | 32 | 0.038 |
| sp|Q00143| IE KDY CYTCININ-BINDING PROTEIN [ECC SVBICE ... | | 31 | 0.050 |
| sp|P04104| CB23E.) CYTBIUDIU-3S DIDosobli'S MeTeNOG- ... | | 30 | 0.010 |
| sp|P52030| LCH3 RBZTH (LCH5) CYLWODDIU-HINDIHG bboLeIU'S ... | | 30 | 0.003 |
| sp|Q13249| MY4 ACACY PARVALBUMIN BETY [Bos bpeacicco] | | 40 | 0.003 |
| sp|Q03430| PRVA BYCLC PARVALBUMIN BETY [bos baeacicco] | | 31 | 4e-13 |
| sp|Q03421| PRVA BOVCO PARVALBUMIN BETY [Bos baeacicco] | | 71 | 4e-13 |
| sp|Q13212| PRVA WEBWE PARVALBUMIN BETY [MehIoccius baeacicco] | | 72 | 4e-14 |
| sp|Q01485| PRVA WEBWE PARVALBUMIN BETY [MehIsccius mexicaus] | | 90 | 1e-12 |
| sp|Q31485| PRVA WEBWE PARVALBUMIN BETY [Mehpaccipa meTisnaus] | | 83 | 1e-12 |
| sp|Q01493| PRVA FYLCH PARVALBUMIN BETY S (CYOME S4-I) [WYJOR Ayp... | | 81 | 7e-11 |
| sp|Q0263S| OMCO FYL (OCW OMCOWODLIU) (OW) [Denscica Cphaina] [kg... | | 88 | 4e-18 |
| sp|R21834| OMCO MOUSE PARVALBUMIN BETY (A) [Denscica Cphaina] | | 88 | 4e-18 |
| sp|O32208| OMCO MOUSE PARVALBUMIN BETY (OCW) [Wehiccius bpivisyis] | | 88 | 4e-18 |
| sp|P33530| OMCO HUMMA PARVALBUMIN BETY (OCW OMCOWODLIU) (OW.) | | 88 | 3e-18 |
| sp|Q00353| PRVA CAPCY PARVALBUMIN BETY (OW) [Cabills ca bicl] | | 80 | 1e-18 |
| sp|O32047| PRVA Q63TA PARVALBUMIN BETY (iv) [Cabills cebio] | | 80 | 1e-18 |
| sp|O6173S| PRVA YWBWE PARVALBUMIN BETY [Gens eschenia] | | 84 | 6e-20 |
| sp|O6161S| PRVA YWBWE PARVALBUMIN BETY [ambpinia means] | | 84 | 3e-20 |
| sp|P43302| PRVA CHICK PARVALBUMIN THAHIC CBA3 (PARVALBUMIN 3).[... | | 82 | 1e-20 |
| sp|O5069A| PRVA E3OfU PARVALBUMIN BETY.[Esox lucius] | | 81 | 1e-20 |
| sp|O5069A| PRVA CAVCY PARVALBUMIN BETY.[Cabills cebio] | | 35 | 2e-51 |

---

# Common Mistake: Sequences Too Closely Related

```
PRVA_MACFU    SMTDLLNAEDIKKAVGAFSAIDSFDHKKFFQMVGLKKKSADVKKVFHILDKDKSGFIEE
PRVA_HUMAN    SMTDLLNAEDIKKAVGAFSAIDSFDHKKFFQMVGLKKKSADVKKVFHMLDKDKSGFIEE
PRVA_GERSP    SMTDLLSAEDIKKAIGAFAAAADSFDHKKFFQMVGLKKKTPDDVKKVFHILDKDKSGFIEE
PRVA_MOUSE    SMTDVLSAEDIKKAIGAFAAAADSFDHKKFFQMVGLKKKNPDEVKKVFHILDKDKSGFIEE
PRVA_RAT      SMTDLLSAEDIKRAIGAFTAADSFDHKKFFQMVGLKKKSADDVKKVFHILDKDKSGFIEE
PRVA_RABIT    AMTELLNAEDIKKAIGAFAAAAESFDHKKFFQMVGLKKKSTEDVKKVFHILDKDKSGFIEE
              :*..:::*,.**********,::**:.* :**********************::******

PRVA_MACFU    DELGFILKGFSPDARDLSAKETKTLMAAGDKDGDGKIGVDEFSTLVAES
PRVA_HUMAN    DELGFILKGFSPDARDLSAKETKMLMAAGDKDGDGKIGVDEFSTLVAES
PRVA_GERSP    DELGFILKGFSSDARDLSAKETKTLLAAGDKDGDGKIGVEEFSTLVSES
PRVA_MOUSE    DELGSILKGFSSDARDLSAKETKTLLAAGDKDGDGKIGVEEFSTLVAES
PRVA_RAT      DELGSILKGFSSDARDLSAKETKTLLAAGDKDGDGKIGVEEFSTLVAES
PRVA_RABIT    EELGFILKGFSPDARDLSVKETKTLMAAGDKDGDGKIGADEFSTLVSES
              :*** .******.******:*.********:*.*********.******.**
```

- MULTIPLE SEQUENCE ALIGNMENTS THRIVE ON DIVERSITY...

- IDENTICAL SEQUENCES BRING NO INFORMATION FOR THE MULTIPLE SEQUENCE ALIGNMENT

-MULTIPLE SEQUENCE ALIGNMENT

# Selecting Diverse Sequences (Opus I)

Searching............................done

| | | Score | E |
| --- | --- | --- | --- |
| Sequences producing significant alignments: | | (bits) | Value |
| sp|P80050|PRVA_MACFU | PARVALBUMIN ALPHA (PARVALBUMIN, MUS... | 187 | 7e-48 |
| sp|P20472|PRVA_HUMAN | (PVALB) PARVALBUMIN ALPHA. [Homo sapiens] | 181 | 3e-46 |
| sp|P02625|PRVA_RAT | (PVALB..) PARVALBUMIN ALPHA. [Rattus norvegicus] | 175 | 2e-44 |
| sp|P80080|PRVA_GERSP | PARVALBUMIN ALPHA. [Gerbillus sp.] | 173 | 1e-43 |
| sp|P02624|PRVA_RABIT | (PVALB) PARVALBUMIN ALPHA. [Oryctolagus cuni... | 172 | 1e-43 |
| sp|P32848|PRVA_MOUSE | (PVALB..) PARVALBUMIN ALPHA. [Mus musculus] | 168 | 2e-42 |
| sp|P80029|PRVA_FELCA | (PVALB) PARVALBUMIN ALPHA. [Felis silvestris... | 163 | 8e-41 |
| sp|P80026|PRVM_CHICK | PARVALBUMIN, MUSCLE. [Gallus gallus] | 154 | 3e-38 |
| sp|P18087|PRVA_CAVPO | (PVALB) PARVALBUMIN. [Cavia porcellus... | 154 | 1e-31 |
| sp|P51434|PRVA_RANCA | PARVALBUMIN ALPHA (PA 4.97). [Rana catesbei... | 132 | 1e-28 |
| sp|P02627|PRVA_RANES | PARVALBUMIN ALPHA. [Rana esculenta] | 121 | 4e-25 |
| sp|P02626|PRVA_AMPME | PARVALBUMIN ALPHA. [Amphiuma means] | 111 | 4e-25 |
| sp|P02621|PRVA_RABIT | PARVALBUMIN ALPHA. [Esox lucius] | 111 | 5e-24 |
| sp|P30563|PRVB_GRAGE | PARVALBUMIN BETA. [Gadus geographica] | 107 | 6e-24 |
| sp|P02618|PRVA_TRISE | PARVALBUMIN ALPHA. [Triakis semifasciata] | 107 | 6e-24 |

## Respect Information!

(multiple sequence alignment block for PRVA_MACFU, PRVA_HUMAN, PRVA_GERSP, PRVA_MOUSE, PRVA_RAT, PRVA_RABIT, TPCC_MOUSE)

- This Alignment Is not Informative about the relation Betwwen TPCC MOUSE and the rest of the sequences.

- A better Spread of the Sequences is needed

TPCC MOUSE

PRVA GERSP
PRVA MACFU
PRVA RABIT
PRVA HUMAN
PRVA RAT
PRVA MOUSE

0.1

Searching..................done

Sequences producing significant alignments:

| | | Score (bits) | E Value |
|---|---|---|---|
| sp|P20472|PRVA_MACFU PARVALBUMIN ALPHA. [Macaca fascicu... | | 187 | 7e-48 |
| sp|P80050|PRVA PARVALBUMIN ALPHA. (PARVALBUMIN, MUS... | | 181 | 3e-46 |
| sp|P20472|PRVA HUMAN (PVALB) PARVALBUMIN ALPHA. [Homo sapiens] | | 181 | 7e-46 |
| sp|P02625|PRVA RAT (PVALB.) PARVALBUMIN ALPHA. [Rattus norvegicus] | | 175 | 2e-44 |
| sp|P80080|PRVA GERSP PARVALBUMIN ALPHA. [Gerbillus sp.] | | 173 | 1e-43 |
| sp|P02624|PRVA RABIT (PVALB) PARVALBUMIN ALPHA. [Oryctolagus cuni... | | 172 | 2e-43 |
| sp|P32848|PRVA MOUSE (PVALB.) PARVALBUMIN ALPHA. [Mus musculus] | | 168 | 3e-42 |
| sp|P80079|PRVA FELCA (PVALB) PARVALBUMIN ALPHA. [Felis silvestris... | | 163 | 8e-41 |
| sp|P80026|PRVM CHICK PARVALBUMIN, MUSCLE. [Gallus gallus] | | 154 | 5e-38 |
| sp|P51434|PRVA CAVPO (PVALB) PARVALBUMIN. [Cavia... | | 123 | 1e-31 |
| sp|P18087|PRVA RANCA PARVALBUMIN ALPHA (FRAGMENT). [Cavia... | | 123 | 1e-28 |
| sp|P02627|PRVA RANCA PARVALBUMIN ALPHA. [PA 4.97). [Rana catesbei... | | 121 | 4e-28 |
| sp|P02626|PRVA AMPME PARVALBUMIN ALPHA. [Amphiuma means] | | 111 | 4e-25 |
| sp|P02620|PRVA ESOLU PARVALBUMIN ALPHA. [Esox lucius] | | 111 | 5e-25 |
| sp|P02614|PRVB GRAGE PARVALBUMIN BETA. [Graptemys geographica] | | 107 | 6e-24 |
| sp|P30563|PRVA TRISE PARVALBUMIN ALPHA. [Triakis semifasciata] | | 107 | 6e-24 |
| sp|P21076|PRV1 SALSA PARVALBUMIN BETA. [Salmo salar] | | 38 | 0.055 |
| sp|P00100|xxx CYP4A1 e KDV CYTIDINE-BINDING PROTEIN, PUTATIVE [Esc... | | 35 | 0.010 |
| sp|b4A1044|CB3S DROME (CB23E-) CYTIDIDIN-35 [Drosophi's melano... | | 33 | 0.003 |
| sp|B2030|LCHS 4EAVLH (LCHS)CYTMODIFIN-BETAILE PROTEIN'S, LONCH-... | | 30 | 0.005 |
| sp|b32d3|BRAM 2CACV CATMODIIMI-[PARVICMENI].LCAIIokymira caric... | | 18 | 4e-15 |
| sp|b50ea9|PRVB CALCY PARVALBUMIN BETA.[Gado sphalrica] | | 17 | 9e-14 |
| sp|b52o34|PRVB CYPCA PARVALBUMIN BETA.[VILEGEIN CYD.C.]·[CYD.C... | | 13 | 2e-14 |
| sp|b50e5S|PRVB OMCO CYVO.OMCOWODIIM (OW) (PARVALBUMIN BETA)·[... | | 80 | 1e-12 |
| sp|b50e53|PRVB MEBME PARVALBUMIN BETA.[Merlandira meriandra] | | 83 | 1e-12 |
| sp|b50e3S|PRVB CYPCA CYPCA.[CYD.C.]·[CYD.C... | | 12 | 5e-14 |
| sp|b50e50|PRVB MEBME PARVALBUMIN BETA.[Merlandira meriandra] | | 83 | 1e-12 |
| sp|b0148S|PRVA 4V1AV PARVALBUMIN BETA.[Sono escleris] | | 86 | 1e-11 |
| sp|b0149S|PRVB 4V1AV PARVALBUMIN BETA.[CYD.C.4.1)... | | 36 | 5e-11 |
| sp|b50es3|PRVB FV1CH (PVALB)·[Parchmets cyetonuse] | | 87 | 1e-11 |
| sp|b50e31|PRVB OMCO KY1.(OW) OMCOWODIIM (OW) (PARVALBUMIN BETA)·[Sr... | | 88 | 4e-18 |
| sp|b0030S|PRVB 4V1AV PARVALBUMIN BETA.[er... | | 88 | 4e-18 |
| sp|b50e04|PRVB FE1CH PARVALBUMIN BETA.[Inchacra carbera] | | 18 | 4e-18 |
| sp|b92003|PRVB WEEBI PARVALBUMIN BETA.[CYD.C.Pilineeria] | | 85 | 4e-18 |
| sp|b21834|PRVB THAMIC (THAMIC CBJ3.(PARVALBUMIN 3)·[... | | 82 | 4e-18 |
| sp|b50ea9|PRVB CHICK PARVALBUMIN BETA.[ohd... | | 92 | 4e-18 |
| sp|b21834|PRVB MOUSE (OW) OMCOWODIIM (OW) (PARVALBUMIN BETA)·[... | | 80 | 3e-18 |
| sp|b50e09|PRVB XENIV PARVALBUMIN BETA.[Xenopra laevis] | | 80 | 3e-18 |
| sp|O32208|OMCO CYVAO (OCW) OMCOWODIIM (OW) (PARVALBUMIN BETA)·(... | | 80 | 3e-18 |
| sp|b5a30|OMCO HOWTN (OCW) OMCOWODIIM (OW) (PARVALBUMIN BETA)·(... | | 80 | 3e-18 |
| sp|b0a54A|PRVA CALCY 4EAVLH (V) (CABIIIng carbo)... | | 80 | 2e-18 |
| sp|b02d4l|PRVB Obai4V PARVALBUMIN III·[-iobsan... | | 83 | 1e-18 |
| sp|b50e09|PRVB 4V1AV PARVALBUMIN BETA S.[CYbME S4.1)... | | 84 | 3e-50 |
| sp|b00ea7|PRVB 4v1AV PARVALBUMIN BETA.[Sono escleris] | | 84 | 3e-50 |
| sp|b50ea4|PRVB YMBME PARVALBUMIN BETA.[ambruina means] | | 84 | 3e-50 |
| sp|b433O2|PRVA CHICK PARVALBUMIN BETA.[CYD.C]·[... | | 85 | 3e-50 |
| sp|b50e0e|PRVA E20FU PARVALBUMIN BETA.[Esox lucius] | | 32 | 1e-50 |
| sp|b50e18|PRVB CALCY PARVALBUMIN BETA.[CABIIIng carbo] | | 85 | 3e-51 |

---

PRVB_CYPCA    -AFAGVLNDADIAAALEACKAADSFNHKAFFAKVGLITSKS...ADDVKKAFAIIDQDKSGFIE
PRVB_BOACO    -AFAGILSDADIAAGLQSCQAADSFSCKTFFAKSGLHSKSKDQLITKVFGVIDRDKSGYIE
PRV1_SALSA    MACAHICKEADIKTALEACKAADTFSFKTFFHTIGFASKSADDVKKAFKVIDQDASGFIE
PRVB_LATCH    -AVAKLLAAADVTAALEGCKADDSFNHKVFFQKTGLAKKSNEELEAIFKIILDQDKSGFIE
PRVB_RANES    -SITDIVSEKDIDAALESVKRAGSFNYKIFFQKVGLAGKSADDAKKVFEILDRDKSGFIE
PRVA_MACFU    -SMTDLLNAEDIKKAVGAFSAIDSFDHKKFFQMVGLKKKSADDVKKVFHIILDKDKSGFIE
PRVA_ESOLU    --AKDILKADDIKKALDAVKAEGSFNHKKFFALVGLKAMSANDVKKVFKAIDADASGFIE
              :          .  *  :       *   **     * .*  .   *  .  *.  **.**

PRVB_CYPCA    EDELKLFLQNFKADARALTDGETKTFLKAGDSDGDGKIGVDEFTALVKA-
PRVB_BOACO    EDELKFLQNFDGKARDLTDKETAEFLKEGDTDGDGKIGVEEFVVLVTKG
PRV1_SALSA    VEEIKLFLQNFCPKARELTDAETKAFLKAGDADGDGMIGIDEFAVIVKQ-
PRVB_LATCH    DEELELFLQNFSAGARTLTKTETEFLKAGDSDGDGKIGVDEFQKLVKA-
PRVB_RANES    QDELGIFLQNFRASARVLSDAETSAFLKAGDSDGDGKIGVDEFQAIVKA-
PRVA_MACFU    EDELGFIILKGFSPDARDLSAKETKTIMAAGDKDGDGKIGIDEFSTLVAES
PRVA_ESOLU    EEELKFVLKSFAADGRDLTDAETKAFLKAADKDGDGKIGIDEFETLVHEA
              ::     *.  *      *   **  ****  **..** **.*

---

- Going Further: Remote Homologues.

- A REASONABLE Model Now Exists.

Phylogenetic tree:
PRV1 SALSA
PRVB BOACO
PRVA ESOLU
PRVB LATCH
PRVB CYPCA
PRVB RANES
PRVA MACFU

```
PRVA_MACFU    ----------------SMTDLLNA----EDIKKA
PRVA_ESOLU    ----------------AKDLLKA-----DDIKKA
PRVB_CYPCA    ----------------AFAGVLND----ADIAAA
PRVB_BOACO    ----------------AFAGILSD----ADIAAG
PRV1_SALSA    ----------------AFAGILSD----ADIAAG
PRVB_LATCH    ----------------MACAHLCKE---ADVTAA
PRVB_RANES    ----------------AVAKILAA----ADIKTA
TPCS_RABIT    -TDQQAEARSYLSEEMIAEFKAAFDMFDADGG--GDISVKELGTVMRMLGQTPTKEELDAI
TPCS_PIG      -TDQQAEARSYLSEEMIAEFKAAFDMFDADGG--GDISVKELGTVMRMLGQTPTKEELDAI
TPCC_MOUSE    MDDIYKAAVEQLTEEQKNEFKAAFDIFVLGAEDGCISTKELGKVMRMLGQNPTPEELQEM
                                                                     :       ::

PRVA_MACFU    VGAFSAIDS--FDHKKFFQMVG------------LKKKSADDVKKVFHILDKDKSGFIEEDELGFI
PRVA_ESOLU    LDAVKAEGS--FNHKKFFALVG------------LKAMSANDVKKVFKAIDADASGFIEEEKFV
PRVB_CYPCA    LEACKAADS--FNHKAFFAKVG------------LTSKSADDVKKAFAIIDQDKSGFIEEDELKLF
PRVB_BOACO    LQSCQAADS--FSCKTFFAKSG------------LHSKSKDQLTKVFGVIDRDKSGYIEEDELKKF
PRV1_SALSA    LEACKAADT--FSFKTFFHTIG------------FASKSADDVKKAFKVIDQDASGFIEVEELKLF
PRVB_LATCH    LEGCKRADDS--FNHKVFFQKTG-----------LAKKSNEELEAIFKILDQDKSGFIEDEELELF
PRVB_RANES    LESVKRAAGS--FNYKIFFQKVG-----------LAGRSAADAKKVFEILDRDKSGFIEQDELGLF
TPCS_RABIT    IEEVDEDGSGTIDFEEFLVMMVRQMKEDAKGKSEEELAECFRIFDRNADGYIDAEELAEI
TPCS_PIG      IEEVDEDGSGTIDFEEFLVMMVRQMKEDAKGKSEEELAECFRIFDRNADGYIDAEELAEI
TPCC_MOUSE    IDEVDEDGSGTVDFDEFLVMMVRCMKDDSKGKSEEELSDLFRMFDKNADGYIDLDELKMM
              :     :          *            :   .:   :.  :  *:.:    .. *:.:.**

PRVA_MACFU    LKGFSPDARDLSAKETKTLMAAGDKDGDGKIGVDEFSTLVAES-
PRVA_ESOLU    LKSFAADGRDLTDAETKAFLKAADRDGDGEIGIDEFETLVHEA-
PRVB_CYPCA    LQNFKADARALTDGETKTFLKAGDSDGDGKIGIDEFETALVKA-
PRVB_BOACO    LQNFEDGKARDLTDKETAEFLKEGDTDGDGMIGIDEFAVLVKG-
PRV1_SALSA    LQNFCPKARELTDAETKAFLKAGDADGDGMIGIDEAVLVKQ--
PRVB_LATCH    LQNFSAGARTLTKTETETFLKAGDSDGDGKIGVDEFQKLVKA--
PRVB_RANES    LQNFRASARVLSDAETSAFLKAGDSDGDGKIGVEEFQALVKA--
TPCS_RABIT    FR----ASGEHVTDEEIESLMKDGDKNNDGRIDFDEFLKMMEGVQ--
TPCS_PIG      FR----ASGEHVTDEEIESLMKDGDKNNDGRIDFDEFLKMMEGVQ--
TPCC_MOUSE    LQ---ATGETITEDDIEELMKDGDKNNDGRIDYDEFLEFMKGVE
              :          :.            **          :       **
```

General information about the entry

| | |
|---|---|
| Entry name | TPC_PATYE |
| Primary accession number | P35622 |

Comments

• **FUNCTION**: TROPONIN IS THE CENTRAL REGULATORY PROTEIN OF STRIATED MUSCLE CONTRACTION. TN CONSISTS OF THREE COMPONENTS: TN-I WHICH IS THE INHIBITOR OF ACTOMYOSIN ATPASE, TN-T WHICH CONTAINS THE BINDING SITE FOR TROPOMYOSIN AND TN-C. THE BINDING OF CALCIUM TO TN-C ABOLISHES THE INHIBITORY ACTION OF TN ON ACTIN FILAMENTS.
• **MISCELLANEOUS**: THIS PROTEIN BINDS ONE CALCIUM ION PER MOLECULE.
• **SIMILARITY**: TO OTHER EF-HAND CALCIUM BINDING PROTEINS.

Features

| | | | |
|---|---|---|---|
| MOD_RES | 1 | 1 | ACETYLATION. |
| DOMAIN | 22 | 33 | ANCESTRAL CALCIUM SITE 1. |
| DOMAIN | 58 | 69 | ANCESTRAL CALCIUM SITE 2. |
| DOMAIN | 95 | 106 | ANCESTRAL CALCIUM SITE 3. |
| CA_BIND | 131 | 142 | SITE 4. |

```
PRVA_MACFU    VGAFSAIDS--FDHKKFFQMVG------------LKKKSADDVKKVFHILDKDKSGFIEEDELGFI
PRVB_BOACO    LQSCQAADS--FSCKTFFAKSG------------LHSKSKDQLTKVFGVIDRDKSGYIEEDELKKF
PRV1_SALSA    LEACKAADT--FSFKTFFHTIG------------FASKSADDVKKAFKVIDQDASGFIEVEELKLF
TPCS_RABIT    IEEVDEDGSGTIDFEEFLVMMVRQMKEDAKGKSEEELAECFRIFDRNADGYIDAEELAEI
TPCS_PIG      IEEVDEDGSGTIDFEEFLVMMVRQMKEDAKGKSEEELAECFRIFDRNADGYIDAEELAEI
TPCC_MOUSE    IDEVDEDGSGTVDFDEFLVMMVRCMKDDSKGKSEEELSDLFRMFDKNADGYIDLDELKMM
TPC_PATYE     SDEMDEEATGRLNCDAWIQLFER----KLKEDLDEREIKEAFRVLDKEKKGVIKVTLRWI
              :     :          *                 :. :  :     .:. . *: :. :
```

```
PRVA_MACFU    LKGFSPDARDLSAKETKTLMAAGDKDGDGKIGVDEFSTLVAES--
PRVB_BOACO    LQNFEDGKARDLTDKETAEFLKEGDTDGDGKIGVEEFVVLVTKG--
PRV1_SALSA    LQNFCPKARELTDAETKAFLKAGDADGDGMIGIDEFAVLVKQ--
TPCS_RABIT    FR----ASGEHVTDEEIESLMKDGDKNNDGRIDFDEFLKMMEGVQ--
TPCS_PIG      FR----ASGEHVTDEEIESLMKDGDKNNDGRIDFDEFLKMMEGVQ--
TPCC_MOUSE    LQ---ATGETITEDDIEELMKDGDKNNDGRIDYDEFLEFMKGVE-
TPC_PATYE     LS---SIGDELTEEEIENMIAEFDTDGSGTVDYEEFKCLMMSSDA
              :          :.            **          :  .:  :
```

# WHAT MAKES A GOOD ALIGNMENT...

- THE MORE DIVERGEANT THE SEQUENCES, THE BETTER

- THE FEWER INDELS, THE BETTER

- NICE UNGAPPED BLOCKS SEPARATED WITH INDELS

- DIFFERENT CLASSES OF RESIDUES WITHIN A BLOCK:
  - Completely Conserved
  - Conserved For Size *and* Hydropathy
  - Conserved For Size *or* Hydropathy

- THE ULTIMATE EVALUATION IS A MATTER OF PERSONNAL JUDGEMENT AND KNOWLEDGE.

---

# DO NOT OVERTUNE!!!

```
chite   ----ADKPKRPLSAYMLWLNSARESIKRENPDFK-VTEVAKKGGELWRGLKD
wheat   --DPNKPKRAPSAFFVFMGEFREEFKQKNPKNKSVAAVGKAAGERWKSLSE
trybr   KKDSNAPKRAMTSFMFFSSDFRS----KHSDLS-IVEMSKAAGAAWKELGP
mouse   -----KPKRPRSAYNIYVSESFQ----EAKDDS-AQGKLKLVNEAWKNLSP
             ***   .:::  ..::  ..      .    *:.  *

chite   AATAKQNYIRALQEYERNGG-
wheat   ANKLKGEYNKAIAAYNKGESA
trybr   AEKDKERYKREM---------
mouse   AKDDRIRYDNEMKSWEEQMAE
        *  :.*.:    :
```

DO NOT PLAY WITH PARAMETERS IF YOU KNOW THE ALIGNMENT YOU WANT: MAKE IT YOURSELF!

```
chite   ----ADKPKRPL-SAYMLWLNSARESIKRENPDFK-VTEVAKKGGELWRGLKD
wheat   --DPNKPKRAP-SAFFVFMGEFREEFKQKNPKNKSVAAVGKAAGERWKSLSE
trybr   KKDSNAPKRAMTSFMFFSSDFRS----KHSDLS-IVEMSKAAGAAWKELGP
mouse   -----KPKRPR-SAYNIYVSESFRS---KHSDLS-IVEMSKAAGAAWKELGP
              ***.  .:.  .:..:          .     *

chite   AATAKQNYIRALQEYERNGG-
wheat   ANKLKGEYNKAIAAYNKGESA
trybr   AEKDKERYKREM---------
mouse   AKDDRIRYDNEMKSWEEQMAE
        *  :.*.:    :
```
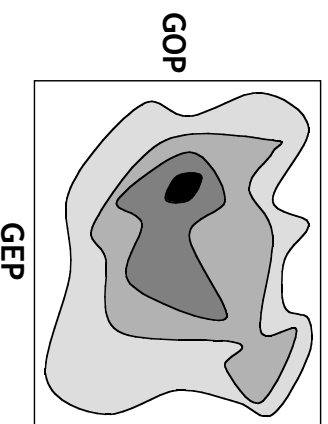
# TUNING or NOT TUNING?

-PARAMETERS TO TUNE USUALLY INCLUDE:
  • GOP/ GEP
  • MATRIX
  • SENSITIVITY Vs SPEED

GOP

GEP

| Substitution Matrices (Etzold and al. 1993) | |
|---|---|
| Gonnet | 61.7 |
| Blosum50 | 59.7 |
| Pam250 | 59.2 |

-MOST METHODS ARE TUNED FOR WORKING WELL ON AVERAGE

-PARAMETERS BEHAVIOUR DO NOT NECESSARILY FOLLOW THE THEORY (i.e. Substitution Matrices).

-A GOOD ALIGNMENT IS USUALLY ROBUST(i.e. Changes little).

-TUNE IF YOU WANT TO CONVINCE YOURSELF.

---

# KEEP A BIOLOGICAL PERSPECTIVE

```
chite   ----ADKPKRPLSAYMIWLNSARESIKRENPDFK-VTEVAKKGGELWRGLKD
wheat   ---DPNKPKRAPSAFFVFMGEFREEFKQKNPKNKSVAAVGKAAGERWKSLSE
trybr   KKDSNAPKRAMTSFMFFSSDFRS----KHSDLS-IVEMSKAAGAAWKELGP
mouse   -----KPKRPRSAYNIYVSESFQ----EAKDDS-AQGKLKLVNEAAWKNLSP
             ***  .::. ::  .      .   .  *  . *

chite   AATAKQNYIRALQEYERNGG-
wheat   ANKLKGEYNKAIAAYNKGESA
trybr   AEKDKERYKREM--------
mouse   AKDDRIRYDNEMKSWEEQMAE
        *   :  :.*  .
```
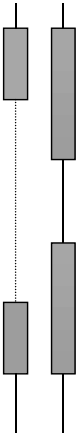
DIFFERENT PARAMETERS

```
chite   AD--K----PKR-PLYMIWLNS-ARESIKRENPDFK-VT-EVAKKGGELWRGL-
wheat   -DPNK----PKRAP-FFVFMGE-FREEFKQKNPKNKSVA-AVGKAAGERWKSLS
trybr   -K--KDSNAPKR-AMT-MFFSSDFR-S-KH-S-DLS-IV-EMSKAAGAAWKELG
mouse   ----K----PKR-PRYNIYVSESFQEA-K--D-D-S-AQGKL-KLVNEAAWKNLS
             *    ***   .:: :::..     : .  *       : *       *: *

chite   KSEWEAKAATAKQNY-I--RALQE-YERNG-G-
wheat   KAPYVAKANKLKGEY-N--KAIAA-YNK-GESA
trybr   RKVYEEMAEKDKERY----K--RE-M------
mouse   KQAYIQLAKDDRIRYDNEMKSWEEQMAE-----
        :  .    :  *    .
```
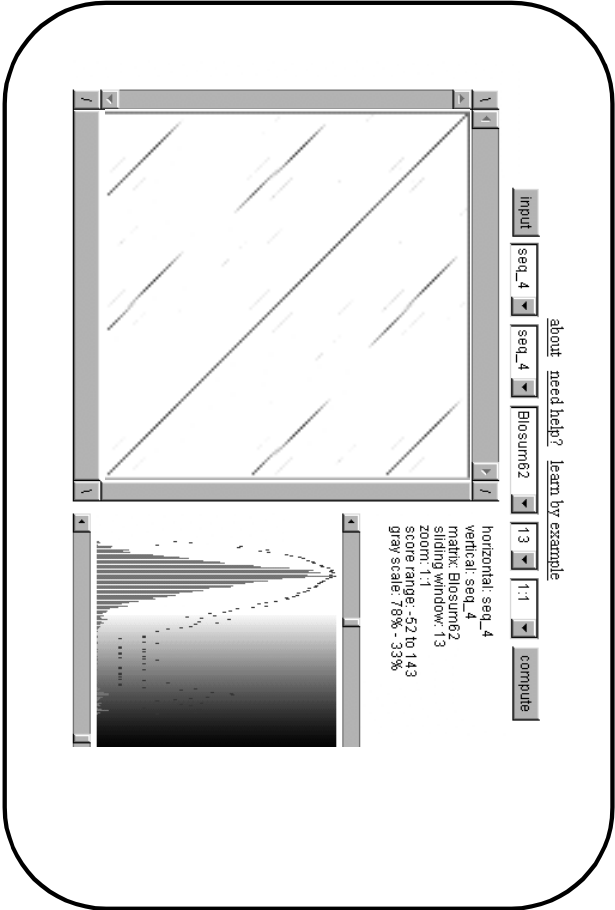
# REPEATS

THERE IS A PROBLEM WHEN TWO SEQUENCES DO NOT CONTAIN
THE SAME NUMBER OF REPEATS

IT IS THEN BETTER TO MANUALLY EXTRACT THE REPEATS AND
TO ALIGN THEM. INDIVIDUAL REPEATS CAN BE RECOGNIZED
USING DOTTER



---

# Choosing The Right Method

| PROBLEM | PROGRAM | METHOD |
|---|---|---|
|  | ClustalW |  |
|  | ClustalW |  |
|  | MSA |  |
|  | DIALIGN II |  |
|  | DIALIGN II |  |

## Playing With Blocks: Mbh1

```
primer          Basic - Helix I        Loop        Helix II
------
Mbh1      SNILERNKARDLALAIRDSER          3     AQVEIITDGEEPAEMIQVLGPK
c-Myc     HNVLERQRRNELKRSFFALRD         11     PKVVILKKATAYILSIQADEHK
TFE-3     HNLIERRRFFNINDRIKELGT         12     NKGTILKASVDYIRKLQKEQQR
E12       NNARERLRVRDINEAFKELGR         12     TKLLILHQAVSVILNIEQQVRE
MyoD      ATMRERRLSKVNEAFETLKR          10     PKVEILRNAIRYIEGLQALLRD
```

"maybe"  
"maybe"  
"no"  
"no"  

4-12  
$5.4e^{-1}$  
1

**Figure 2.** Evaluation of block alignments involving Mbh1 using the MACAW program (Schuler et al. 1991) as described by Prendergast and Ziff (1991).

Blocks are reproduced as in their Figure 1 except that three conflicting residues in the second block for c-Myc are taken from the reference that they cite. MACAW scores a block as "yes," "no," or "maybe" depending on how the block compares to a random protein model. Brackets show the interval and MACAW evaluation for the first block, where successive columns of primer-related sequences are trimmed off. For the second block, all MACAW evaluations were "no", so that the actual P-score (an exponential number) is reported. Prendergast and Ziff reported that all negative controls gave P=1. For all loop sizes from 4 to 12, a score of P=1 was obtained.

**A**

Motif 1

| | | | | |
|---|---|---|---|---|
| ECOLYS | 194 | IRQEMVNRGFMEVET | P | MM |
| YSTLYS | 257 | IRRFLDQRKFIEVET | P | MM |
| ECOASP | 148 | VRRFMDDHGFLDIET | P | ML |
| RATASP | 206 | FRETLINKGFVEIQT | P | KI |
| YSTASP | 258 | FREYLATKKFTEVHT | P | KL |
| ECOHIS | 29 | LKNVLGSYGYSEIRL | P | IV |
| HAMHIS | 80 | ICCCFKRHGAEVIDT | P | VE |
| YSTHIS | 79 | LSGLFKKHGGVTIDT | P | VF |
| ECOPHE | 132 | EI.EDDYHNFDALNI | P | GH |
| MITPHE | 91 | EPVVTTMENFDSLGF | P | KD |
| YSTPHE | 253 | QYVETGFWNFDALYV | P | QQ |
| ECOASN | 148 | LHREENEQGFFWVST | P | LI |
| YSTTHR | 367 | LRTEYRKRGYEEVIT | P | NM |
| **ECOTHR** | **281** | **VRSKLKEYQYQEVKG** | **P** | **FM** |
| **ECOPRO** | **57** | **VREEMNNAGAIEVSM** | **P** | **VV** |
| **ECOSER** | **181** | **LDLHTEQHGYSENYV** | **P** | **YL** |
| YSTSER | 196 | GLQFLAAKGYIPLQA | P | VM |

**B**

Scrambled set

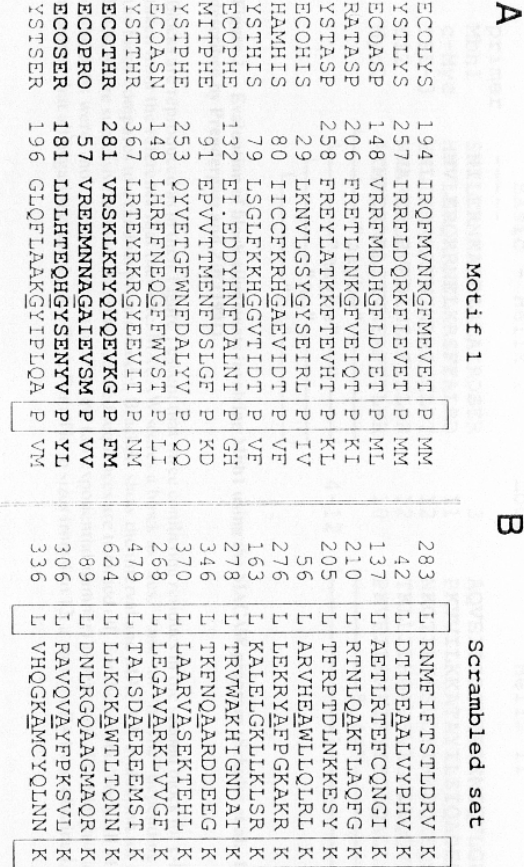| | | | |
|---|---|---|---|
| 283 | L | RNMEIFTSTLDRV | K |
| 42 | L | DTIDEAALVYPHV | K |
| 137 | L | AETLRTEFCQNGI | K |
| 210 | L | RTNLQAKFLAQFG | K |
| 205 | L | TFRPTDLNKKESY | K |
| 56 | L | ARVHEAWLLQLRL | K |
| 276 | L | LEKRYAFPGKAKR | K |
| 163 | L | KALELGKLLKLSR | K |
| 278 | L | TRVWAKHIGNDAI | K |
| 346 | L | TKENQAARDDEEG | K |
| 370 | L | LAARVASEKTEHL | K |
| 268 | L | LEGAVARKLVVGF | K |
| 479 | L | TAISDAEREEMST | K |
| 624 | L | LLKCKAWTLTQNN | K |
| 89 | L | DNLRGQAAGMAQR | K |
| 306 | L | RAVQVAYFPKSVL | K |
| 336 | L | VHQGKAMCYQLNN | K |

**Figure 1.** Comparison of a block derived from aminoacyl tRNA synthetase sequences to a block constructed from scrambled sequences.

(A) Motif 1 from Figure 3 of Eriani et al. (1990), using their criteria for strictly and strongly conserved regions. The sequences in bold type show the initial alignment from their Fig. 2. (B) A block from scrambled sequences with average searching strength found by the MOTIF program (Smith et al. 1990) in the 17 randomized synthetase sequences. The position of the first residue is indicated for each segment. The single-letter code for the amino acid residues is as follows: A: Ala; C: Cys; D: Asp; E: Glu; F: Phe; G: Gly; H: His; I: Ile; K: Lys; L: Leu; M: Met; N: Asn; P: Pro; Q: Gln; R: Arg; S: Ser; T: Thr; V: Val; W: Trp; and Y: Tyr.

---

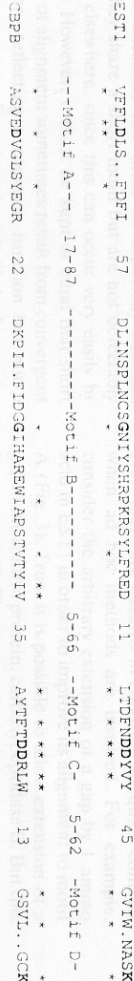| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| EST1 | VFFDLS..FDFI | 57 | DLINSPLNCSGNIYSHRPKRSYLFRED | 11 | LCDFNDDYVY | 45 | GVIW NASK |
| | --Motif A--- | 17-87 | -------------Motif B------- | 5-66 | --Motif C- | 5-62 | --Motif D- |
| CBPB | ASVEDVGLSYEGR | 22 | DKPII.FIDGGIHAREWIAPSTVTYIV | 35 | AYTFTDDRLW | 13 | GSVL..GCK |

**Figure 3.** Alignment of EST1 and Carboxypeptidase B (CBPB) with four blocks derived from selected RNA-dependent RNA polymerases using the procedure of Lundblad and Blackburn (1990).

The numbers indicate distance in amino acid residues between the blocks for the known polymerases at the top, and for EST1 and CBPB beneath. The EST1 alignment is taken from their figure. Asterisks indicate conserved residues for both the reverse transcriptases and the RNA polymerases. To align CBPB using the same rules, Motif C was anchored at the [F Y]xDD and the other 3 invariant residues (in bold) were found by scanning the flanking regions while maintaining interblock distances consistent with those seen in known polymerases. A single gap extension was allowed, consistent with their arbitrary extension of the 1-amino acid gap to 2-amino acids in Motif A. Hydrophobic residues are MILVFWAP.

# Conclusion

**The Best Alignment Method:**
- Your Brain
- The Right Data

**The Best Evaluation:**
- Your Eyes
- Experimental Information (SwissProt)

**What Can I Conclude:**
- Homology=> Information Extrapolation

**How Can I go Further?:**
- PrositePatterns.
- PrositeProfiles.