

## Supplementary information

# Structural homology guided alignment of cysteine rich proteins

---

Thomas M A Shafee<sup>1</sup>, Andrew J Robinson<sup>2,3</sup>, Nicole van der Weerden<sup>1</sup>, Marilyn A Anderson<sup>1</sup>

<sup>1</sup> Department of Biochemistry, La Trobe Institute of Molecular Sciences, La Trobe University, Melbourne, Australia, 3086

<sup>2</sup> College of Science, Health and Engineering, La Trobe University, Melbourne, Australia, 3086

<sup>3</sup> Life Sciences Computation Centre, Victorian Life Sciences Computation Initiative, Melbourne, Australia, 3053

Email: T.Shafee@LaTrobe.edu.au

## SUPPLEMENTARY FIGURES

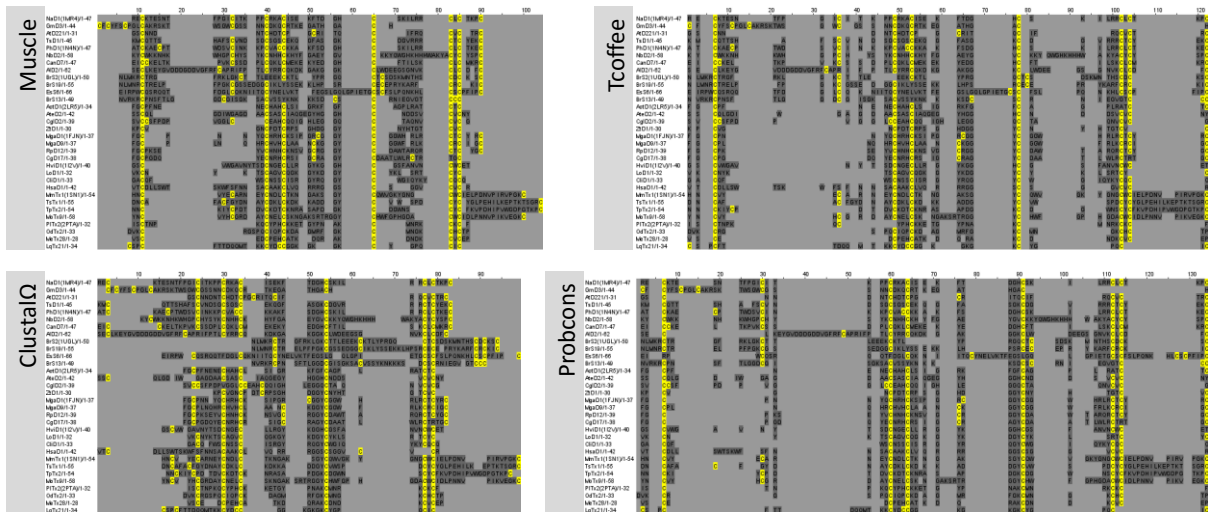
Figure S1   Misalignment by standard algorithms	2
Figure S2   Misalignment quantification	3
Figure S3   Alignments of sub-groups to closest sequence of known structure	4
Figure S4   Barcoded sequences	5
Figure S5   Misalignment by standard algorithms (larger data set)	6
Figure S6   Misalignment quantification (larger data set)	7
Figure S7   Summary of biophysical attributes by <i>looppproperties.xls</i>	8

## SUPPLEMENTARY DATA FOLDERS:

1. Scripts, barcodes and readme
2. Failed initial alignments
3. DALI structural alignment
4. Aligned sub-groups
5. Barcoded sub-groups
6. Final combined alignment

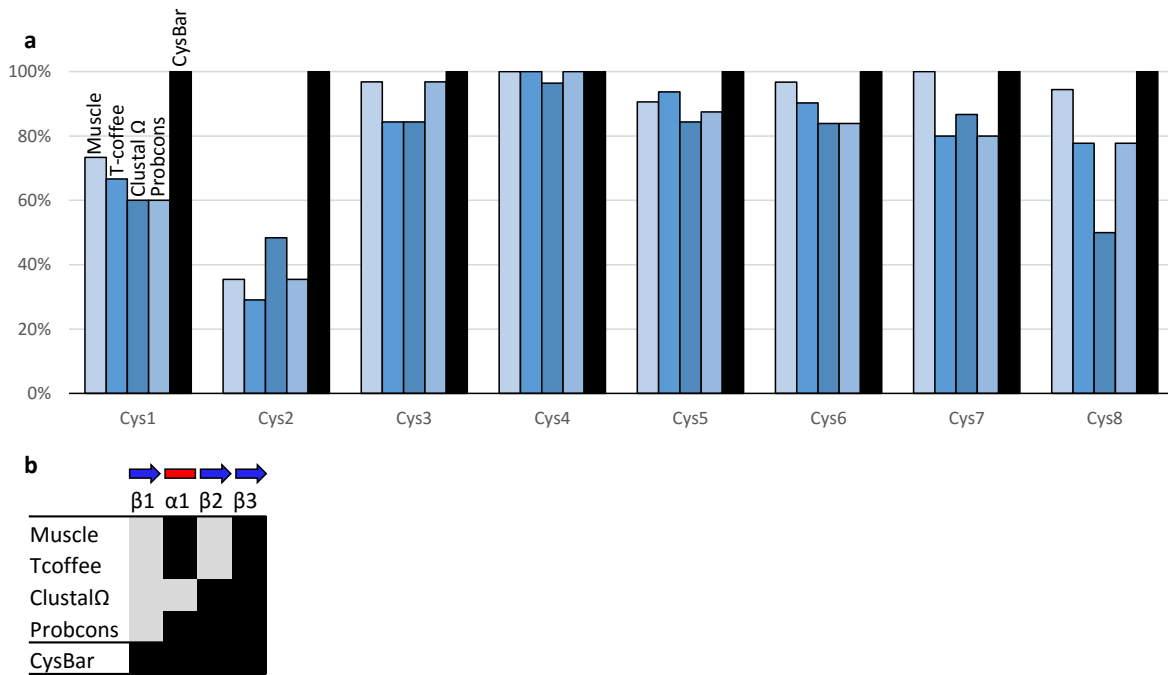
Webtool implementation available at [CysBar.science.latrobe.edu.au](http://CysBar.science.latrobe.edu.au)

Python scripts, readme, and example datasets available at [Github.com/TS404/CysBar](https://github.com/TS404/CysBar)



**Figure S1 | Misalignment by standard algorithms**

Alignments of structurally related defensins from plants, fungi and invertebrates using Muscle, Tcoffe, Clustal Q or Probcons. Default settings on these algorithms generate conflicting, irreproducible alignments in which structurally homologous cysteines fail to align. Sequences coloured by JalView with cysteines in yellow, any other residue in grey, gaps in light grey. *Genbank accession numbers: 159162710, 571550504, 74820403, 332196243, 38492523, 6552502, 297824339, 158853052, 7209504, 557088590, 270381566, 392935432, 156179583, 88178907, 398388411, 193806528, 560135893, 555699991, 193806210, 159162452, 56462336, 49182286, 386642833, 51317001, 378748964, 399762321, 346655549, 6573542, 1173404, 487523606, 41017872.*



**Figure S2 | Misalignment quantification**

(a) Percentage of correct alignment for the 8 cysteines with known structural homology for each of Muscle, Tcoffee, Clustal Ω, Probcons and CysBar alignments. (b) For each alignment, secondary structural elements with one or more columns displaying >50% insertion or deletions in each alignment are indicated in light grey. (For alignments, see fig S1).

```

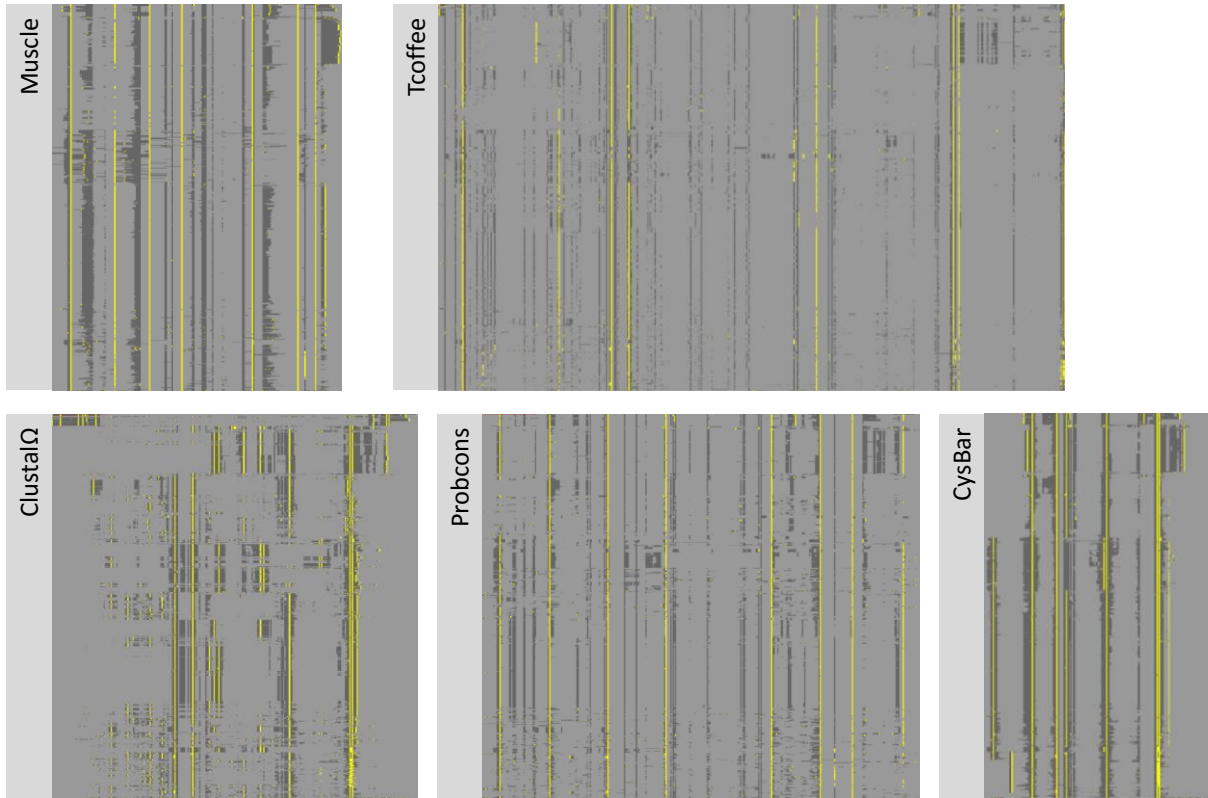
▶ NaD1(1MR4)/1-47      RECKTESNTFFPGICITKPPCRKAC ISEKFTDGHCSKILRRCLCTKPC
GmD3/1-44             CF CYFSC PGLCAKRSKTWSGWCGSSNNCDKQCRTEGATHGACH
AD221/1-31           GSCNNDNTCHDTC PGCRITGC IFRGCVCTR C
TsD1/1-45            KMCQTTSHAFS CVNDSGCSGS EKQGFASGK DGVRRRCICYEK
▶ PhD1(1N4N)/1-47     ATCKAE CPTWDSVCTN KKPVA CCKKAK FSDGHCSKILRR C LCTKEC
NbD2/1-58            KYCMMK NHKVMGPCHY SYKCNHHCKHYFGAEYGVCKKYQWGHKHHHMAKYACVYSPC
CanD7/1-47           EICCKE LTKPVKCSS DP LQK LQMEKEYEDGHCF TILSK C LCMKRC
AD2/1-62             SEC LKEYGVDDDDGDDVGFRC APRIFPTLCYRRC QKDKGAKGGK C LWDEEGSGNVK C L C DFC
▶ BrS2(1UGL)/1-50    NLMKRC TRGFRKLGKCTTLEEEKCKTLYPR GQ CTCSDSKMNTHSCDCKS C
BrS19/1-55           NLMNRC TRELFPFGKCGSSEDGGCIKLYSSEKKLHPSR CECEPRYKARF CRCKI C
EsS8/1-66            EIRPWCQSRQOTFDGLCDKNIITQCYNELVKTFEGLGQLGPIETGCS CFSLPQNKHLCS C PFI C
BrS13/1-49           NVRKRC PNSFTLGGDCGISGKSACVSSYKNKKK SD CSCRNIEGVGT C C C
▶ AotD1(2LR5)/1-34   FGC PFNENECHAHC LSIGR KFGFCAGPLRATCTC
AeD2/1-42            SSC QLGDWAGDAACSASC LAQGEY HGGHCND S VCVNY
CglD2/1-39           SVC C SFPD PVGGL CCEAHC QQIGHLEGGQCTAQ N VCVG
ZiD1/1-30            KPC VGN CPDTC RPSGH DGGYCNH TGTVC
▶ HviD1(1I2V)/1-40  BSCWGAVNYT SDCNGEC LLRGYKGGHC GSFANVNCWET
LoD1/1-32            VKC NYK TSCAGVCDGKGYKGYCYKLS RTCYC
CtiD1/1-33           GACDF WS CNSSCISRGYRGGYCWGIQYKYCQDQ
HsaD1/1-42           VTC DLLSWTSKWF SFNNSACAACLVRRRRGGSCSG GVCVCR
▶ MgaD1(1FJN)/1-37  FGCP NNYQCHRHCKSI PGR C GGYCGGMH RLRCTCYRC
MgaD8/1-37           FGCP LNQHRCHVHCLAA N CKGGYCGWF RLCRCIC
RpD12/1-39           FGCP KSEYVCHNHCKNSVG CRGGYCDAWTARQRCICYG
CgD17/1-38           FGCP GDQYECNRHCRS I G CRAGY CDAATLWLRCTR TGC
▶ MmTx1(1SN1)/1-54  HNCVYECR RNEYCNDLCTKNGAKS GYCD WVGKYGNGCWC IELPDNV PI RVP GK C
TsTx1/1-55           DNCAFACFGYDNAYCDKLC KDKKADD GYCV W SPDCYCYGLPEHILKEPT KTSGR C
TpTx2/1-54           NNCKIYCP DTDVCKDCKNRASAP DGKCDGW NSCYCFKVPDHI PWWGDPGTKPC
MoTx8/1-58           YNCVYHCG RDAYCNELC SKNGAKSRTGGYCH WFGPHGDACWCIDL PNNV PI KVEGK C
▶ PITx2(2PTA)/1-32  ISCTNPKQ CYPHCKKETGYPNAKCMNRKCKCF
OdTx2/1-33           DVKCRGSPQ CIQPC KDAGMRF GKCMNGKCHCTP
MeTx28/1-28          VSCED CPEHC ATKDQRAKCDNDKCVCEP
LqTx21/1-34         CSPCF TTDQQMTKKCYDCC GKGKGYGYPQCIC

```

**Figure S3 | Alignments of sub-groups to closest sequence of known structure**

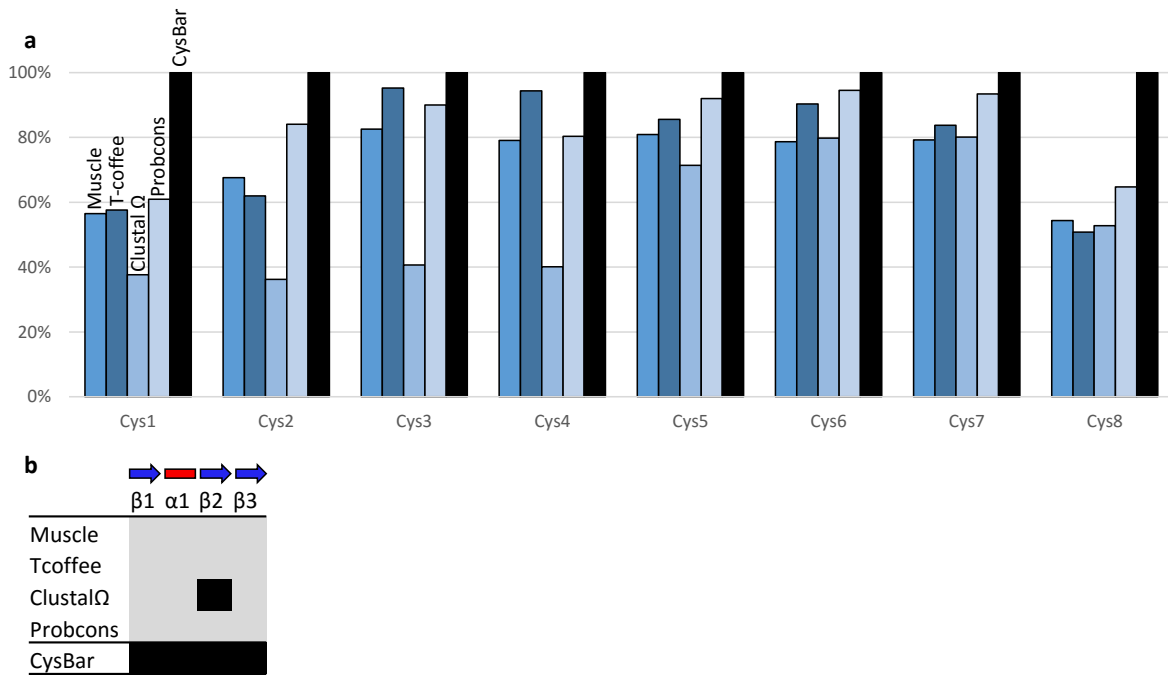
Alignments of sequences of unknown structure to their closest sequence of known structure (indicated by a black triangle) allows homologous cysteines to be found within sub-groups. Aligned with Muscle using default settings. Sequences coloured by JalView with cysteines in yellow, any other residue in grey, and gaps in light grey. *PDB accession numbers: 1MR4, 1N4N, 1UGL, 2LR5, 1I2V, 1FJN, 1SN1, 2PTA.*





**Figure S5 | Misalignment by standard algorithms (larger data set)**

For larger alignment of 965 sequences, using default settings on Muscle, Tcoffee, Clustal  $\Omega$  or Probcons generates conflicting, irreproducible alignments in which structurally homologous cysteines fail to align. The frequent misalignment of cysteines causes erroneous insertion and deletion predictions, leading to alignments with a large number of columns compared to the final CysBar alignment of the same sequences. Sequences coloured with cysteines in yellow, any other residue in grey, gaps in light grey.



**Figure S6 | Misalignment quantification (larger data set)**

(a) Percentage of correct alignment for the 8 cysteines with known structural homology for each of Muscle, Tcoffee, ClustalΩ, Probcons and CysBar alignments. (b) For each alignment, secondary structural elements with one or more columns displaying >50% insertion or deletions in each alignment are indicated in light grey. (For alignments, see fig S5).

