# Protein structure prediction

Miguel Andrade
Faculty of Biology,
Institute of Organismic and Molecular Evolution
Johannes Gutenberg University
Mainz, Germany
andrade@uni-mainz.de
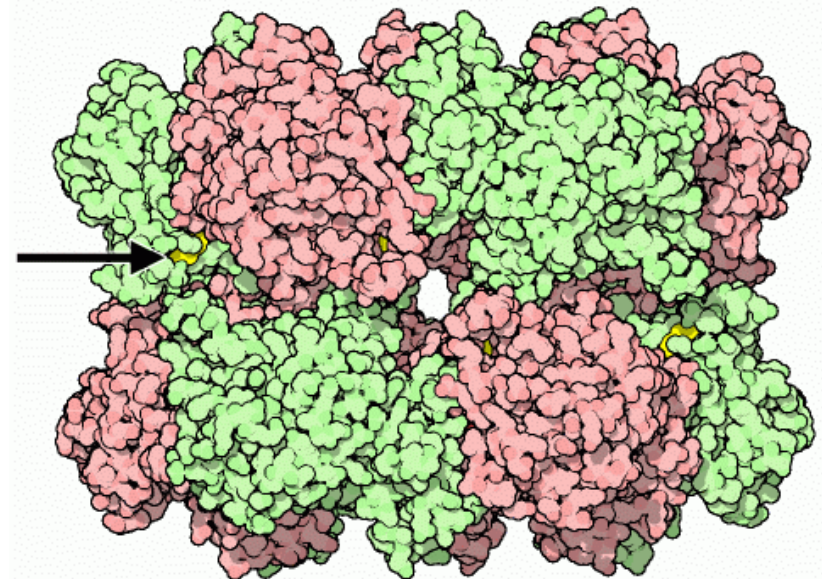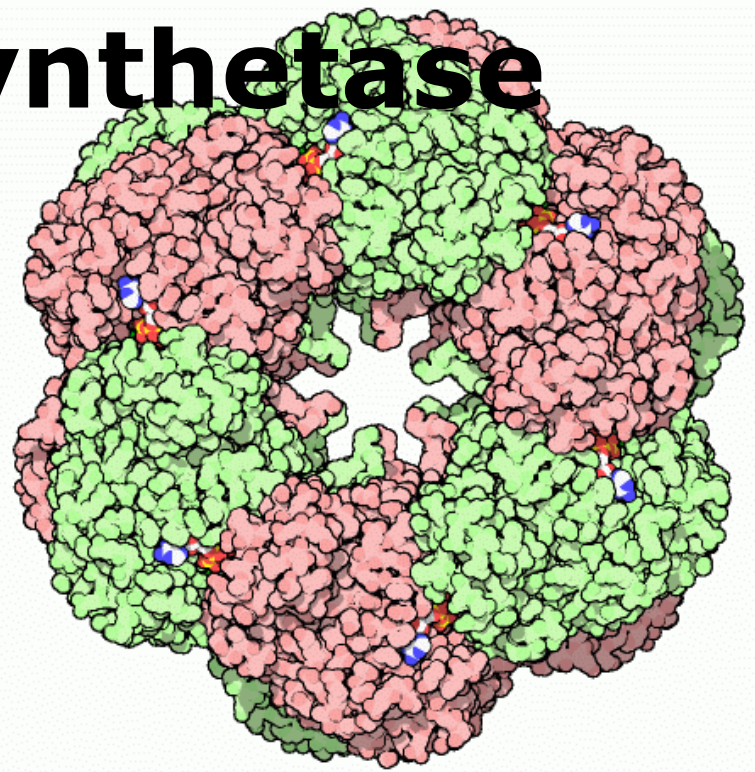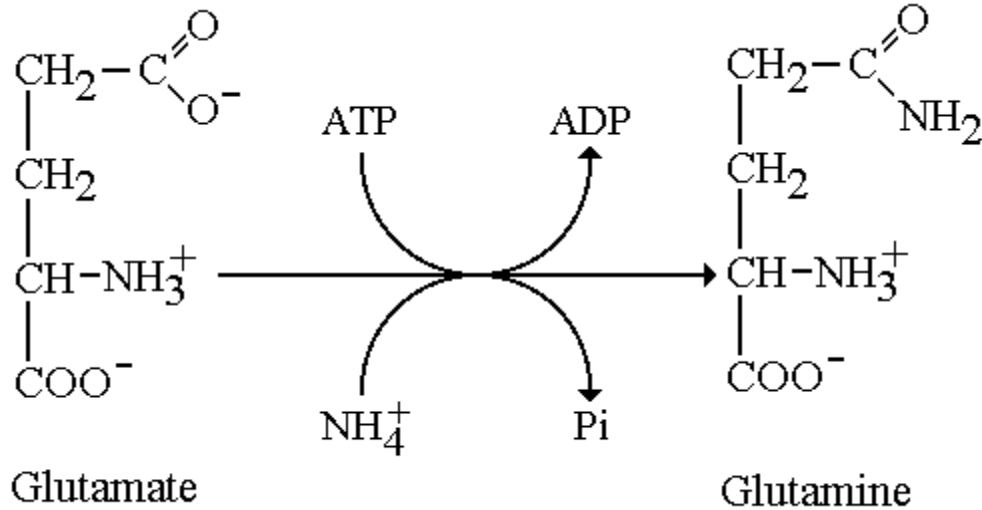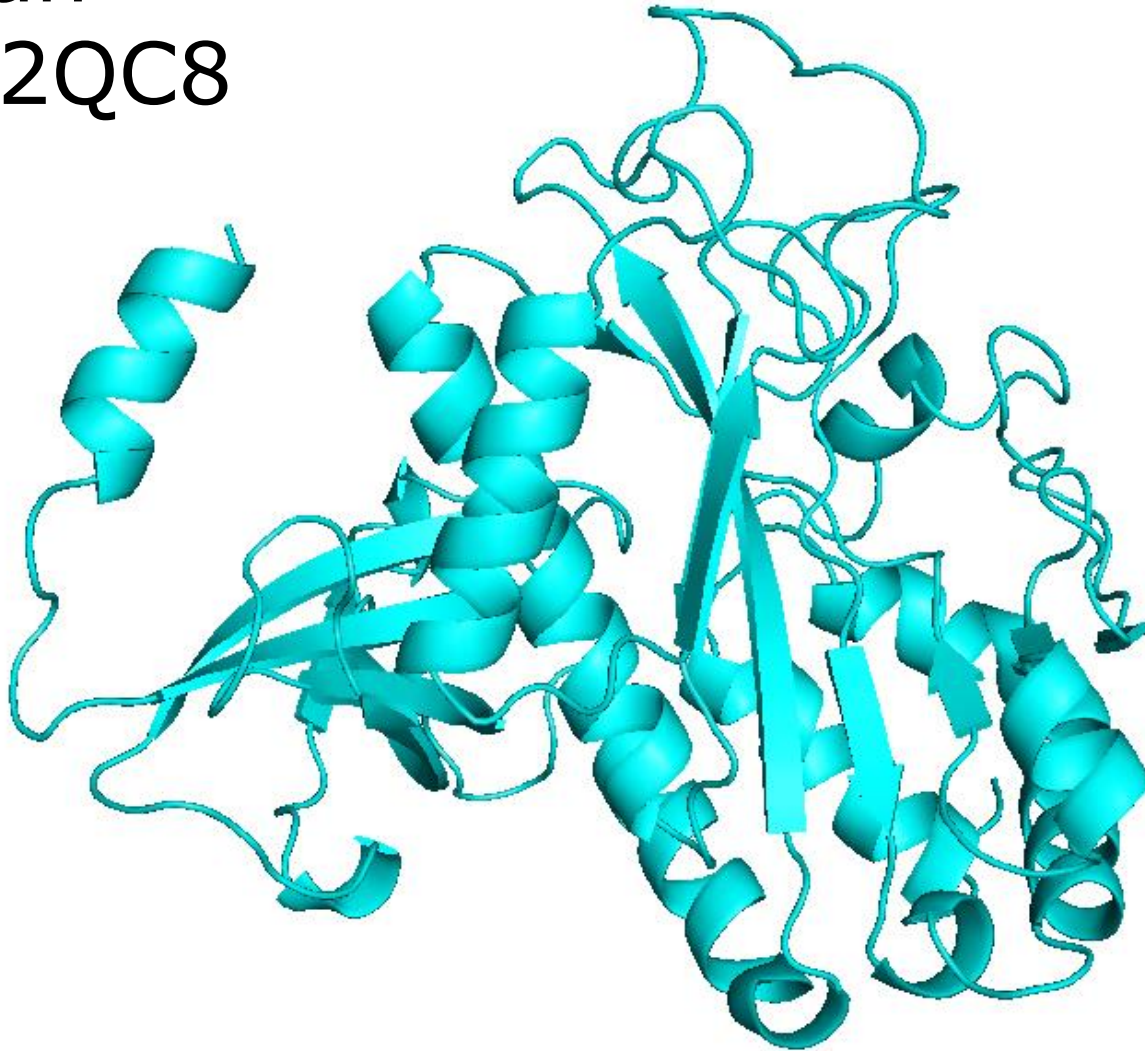
Mount Everest

Age: 60M years

# Glutamine synthetase
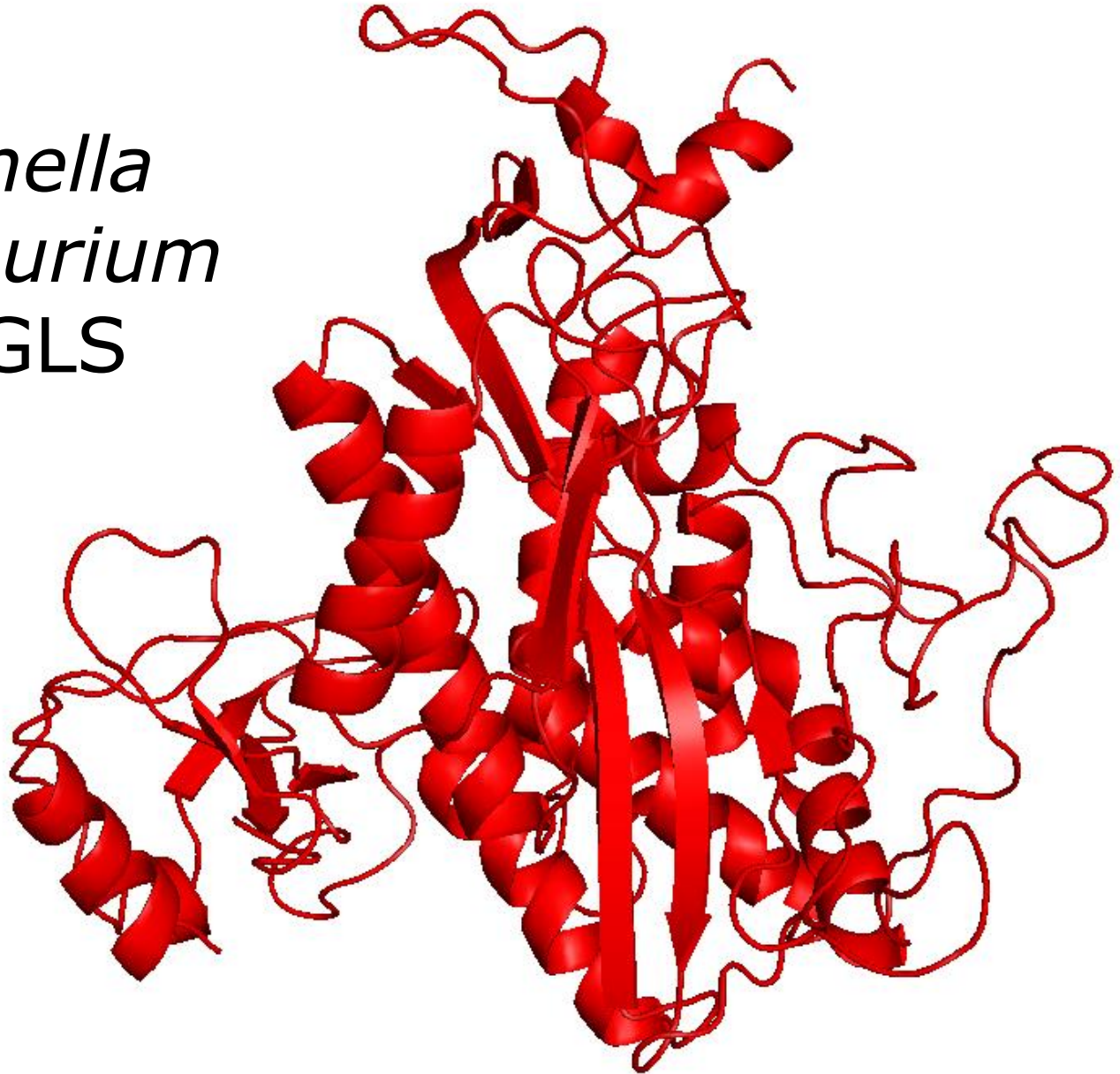


Age: +3500M years

# Glutamine synthetase

Human
PDB:2QC8

# Glutamine synthetase

*Salmonella typhimurium*
PDB:2GLS

# Glutamine synthetase

# Time line

Earth: 4.6 By

Origin of life: 3.9 By – 3.5 By

Last Common Ancestor:
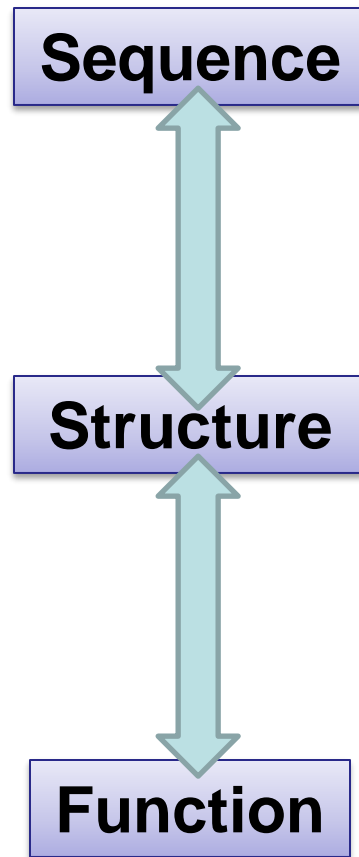3.5 – 3.8 By

Glansdorff & Labedan
(2008) *Biology Direct*

4.29 By

Sheridan *et al.* (2003)
*Geomicrobiology Journal*

# Sequence and function

Evolutionary constraints

**Sequence**

**Structure**

**Function**

MTQDELKKAVGWAALQYVQ
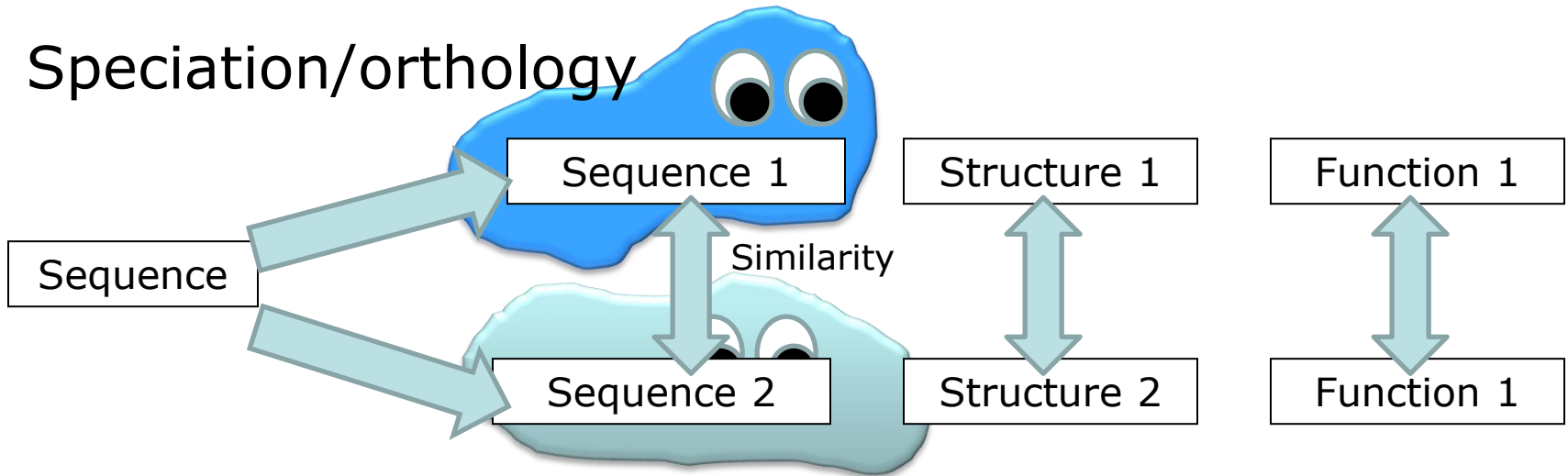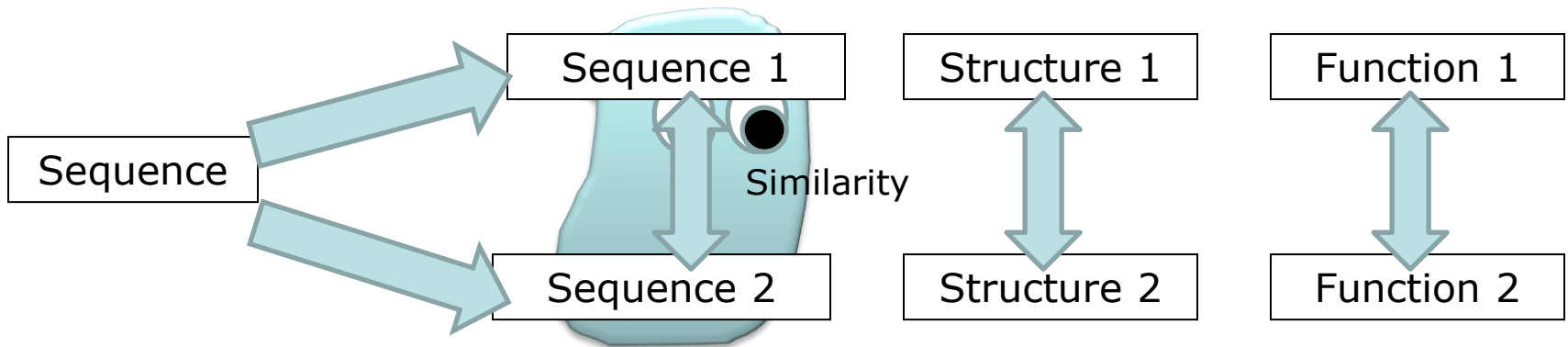PG                    DA
LG                    ST
EK

# Sequence and function

Evolutionary constraints

Speciation/orthology



Gene duplication/paralogy

# Sequence pairwise alignment

>gs_human gi|74271837|ref|NP_001028216.1| glutamine synthetase [Homo sapiens]
MTTSASSHLNKGIKQVYMSLPQGEKVQAMYIWIDGTGEGLRCKTRTLDSEPKCVEELPEWNFDGSSTLQS
EGSNSDMYLVPAAMFRDPFRKDPNKLVLCEVFKYNRRPAETNLRHTCKRIMDMVSNQHPWFGMEQEYTLM
GTDGHPFGWPSNGFPGPQGPYYCGVGADRAYGRDIVEAHYRACLYAGVKIAGTNAEVMPAQWEFQIGPCE
GISMGDHLWVARFILHRVCEDFGVIATFDPKPIPGNWNGAGCHTNFSTKAMREENGLKYIEEAIEKLSKR
HQYHIRAYDPKGGLDNARRLTGFHETSNINDFSAGVANRSASIRIPRTVGQEKKGYFEDRRPSANCDPFS
VTEALIRTCLLNETGDEPFQYKN


>gs_salmonella gi|16767272|ref|NP_462887.1| glutamine synthetase [Salmonella
enterica subsp. enterica serovar Typhimurium str. LT2]
MSAEHVLTMLNEHEVKFVDLRFTDTKGKEQHVTIPAHQVNAEFFEEGKMFDGSSIGGWKGINESDMVLMP
DASTAVIDPFFADSTLIIRCDILEPGTLQGYDRDPRSIAKRAEDYLRATGIADTVLFGPEPEFFLFDDIR
FGASISGSHVAIDDIEGAWNSSTKYEGGNKGHRPGVKGGYFPVPPVDSAQDIRSEMCLVMEQMGLVVEAH
HHEVATAGQNEVATRFNTMTKKADEIQIYKYVVHNVAHRFGKTATFMPKPMFGDNGSGMHCHMSLAKNGT
NLFSGDKYAGLSEQALYYIGGVIKHAKAINALANPTTNSYKRLVPGYEAPVMLAYSARNRSASIRIPVVA
SPKARRIEVRFPDPAANPYLCFAALLMAGLDGIKNKIHPGEAMDKNLYDLPPEEAKEIPQVAGSLEEALN
ALDLDREFLKAGGVFTDEAIDAYIALRREEDDRVRMTPHPVEFELYYSV

# Sequence pairwise alignment

## BLAST (Altschul et al, 1990)

```
>lcl|39919 unnamed protein product
Length=469

 Score = 70.5 bits (171),  Expect = 1e-17, Method: Compositional matrix adjust.
 Identities = 102/363 (28%), Positives = 138/363 (38%), Gaps = 96/363 (26%)


Query  62   FDGSSTLQSEGSN-SDMYLVPAA--MFRDPFRKDPNKLVLCEVFK------YNRRP----  108
            FDGSS    +G N SDM L+P A    DPF  D   ++ C++ +      Y+R P
Sbjct  50   FDGSSIGGWKGINESDMVLMPDASTAVIDPFFADSTLIIRCDILEPGTLQGYDRDPRSIA  109


Query  109  --AETNLRHTCKRIMDMVSNQHPWFGMEQEYTLMGTDGHPFGWPSNGF------------  154
              AE  LR T   I D V      FG E E+ L   D   FG   +G
Sbjct  110  KRAEDYLRATG--IADTV-----LFGPEPEFFLF--DDIRFGASISGSHVAIDDIEGAWN  160


Query  155  -------------PGPQGPYYCGVGADRAYGRDI--------------VEAHYRACLYAG  187
                         PG +G Y+     D A +DI              VEAH+      AG
Sbjct  161  SSTKYEGGNKGHRPGVKGGYFPVPPVDSA--QDIRSEMCLVMEQMGLVVEAHHHEVATAG  218


Query  188  VKIAGTNAEVMPAQWEFQIGPCEGISMGDHLWVARFILHRVCEDFGVIATFDPKPIPG-N  246
                 T   M  +                D + + ++++H V   FG  ATF PKP+ G N
Sbjct  219  QNEVATRFNTMTKK-------------ADEIQIYKYVVHNVAHRFGKTATFMPKPMFGDN  265


Query  247  WNGAGCHTNFSTKAMREENGLKYIEEAIEKLSKRHQYHIRAYDPKGGLDNA---------  297
             +G  CH + +        +G KY        LS++  Y+I            NA
Sbjct  266  GSGMHCHMSLAKNGTNLFSGDKY-----AGLSEQALYYIGGVIKHAKAINALANPTTNSY  320


Query  298  RRLTGFHETSNINDFSAGVANRSASIRIPRTVGQEKKGYFEDRRPSANCDPFSVTEALIR  357
            +RL  +E   + +SA   NRSASIRIP V   K    E R P    +P+    AL+
Sbjct  321  KRLVPGYEAPVMLAYSA--RNRSASIRIP-VVASPKARRIEVRFPDPAANPYLCFAALLM  377


Query  358  TCL  360
            L
Sbjct  378  AGL  380
```

# Multiple sequence alignment

```
>gs_human gi|74271837|ref|NP_001028216.1| glutamine synthetase [Homo sapiens]
MTTSASSHLNKGIKQVYMSLPQGEKVQAMYIWIDGTGEGLRCKTRTLDSEPKCVEELPEWNFDGSSTLQS
EGSNSDMYLVPAAMFRDPFRKDPNKLVLCEVFKYNRRPAETNLRHTCKRIMDMVSNQHPWFGMEQEYTLM
GTDGHPFGWPSNGFPGPQGPYYCGVGADRAYGRDIVEAHYRACLYAGVKIAGTNAEVMPAQWEFQIGPCE
GISMGDHLWVARFILHRVCEDFGVIATFDPKPIPGNWNGAGCHTNFSTKAMREENGLKYIEEAIEKLSKR
HQYHIRAYDPKGGLDNARRLTGFHETSNINDFSAGVANRSASIRIPRTVGQEKKGYFEDRRPSANCDPFS
VTEALIRTCLLNETGDEPFQYKN
>gs_vulca gi|307594850|ref|YP_003901167.1| glutamine synthetase [Vulcanisaeta
distributa DSM 14429]
MPTRNLEIEPADLWRILKASGIKYVKFIIVDINGAPRSEIVPIDMAKDLFIDGMPFDASSIPSYSTVNKS
DFVAYVDPRAVYVEYWQDGKVADVFTMVSDIADKPSPLDPRRVLNDALEQARSKGYEFLMGVEVEFFVIK
EDGGKPVFADPGIYFDGWNVTVQSQFMKELITAIADAGINYTKTHHEVAPSQYEVNIGATDPLRLADQIV
YFKIMAKDIARKYGLVATFMPKPFWGVNGSGAHTHISVWKDGKNLFQSSTGKITEECGYAISAILSNARA
LSSFVAPLVNSYKRLVPHYEAPTRIVWGYANRSAMIRIPQYKMRINRIEYRHPDPSMNPYLAFTAIIKTM
IRGLEEKKEPPPPTEEVAYELANALETPATLEDTLKELSKSFLATELPSELVNAYIKIKQNEWEDYLTNV
GPWEKTWNIITQWEYNKYLVTA
>gs_salmonella gi|16767272|ref|NP_462887.1| glutamine synthetase [Salmonella
enterica subsp. enterica serovar Typhimurium str. LT2]
MSAEHVLTMLNEHEVKFVDLRFTDTKGKEQHVTIPAHQVNAEFFEEGKMFDGSSIGGWKGINESDMVLMP
DASTAVIDPFFADSTLIIRCDILEPGTLQGYDRDPRSIAKRAEDYLRATGIADTVLFGPEPEFFLFDDIR
FGASISGSHVAIDDIEGAWNSSTKYEGGNKGHRPGVKGGYFPVPPVDSAQDIRSEMCLVMEQMGLVVEAH
HHEVATAGQNEVATRFNTMTKKADEIQIYKYVVHNVAHRFGKTATFMPKPMFGDNGSGMHCHMSLAKNGT
NLFSGDKYAGLSEQALYYIGGVIKHAKAINALANPTTNSYKRLVPGYEAPVMLAYSARNRSASIRIPVVA
SPKARRIEVRFPDPAANPYLCFAALLMAGLDGIKNKIHPGEAMDKNLYDLPPEEAKEIPQVAGSLEEALN
ALDLDREFLKAGGVFTDEAIDAYIALRREEDDRVRMTPHPVEFELYYSV
>gs_yeast gi|330443748|ref|NP_015360.2| Gln1p [Saccharomyces cerevisiae S288c]
MAEASIEKTQILQKYLELDQRGRIIAEYVWIDGTGNLRSKGRTLKKRITSIDQLPEWNFDGSSTNQAPGH
DSDIYLKPVAYYPDPFRRGDNIVVLAACYNNDGTPNKFNHRHEAAKLFAAHKDEEIWFGLEQEYTLFDMY
DDVYGWPKGGYPAPQGPYYCGVGAGKVYARDMIEAHYRACLYAGLEISGINAEVMPSQWEFQVGPCTGID
MGDQLWMARYFLHRVAEEFGIKISFHPKPLKGDWNGAGCHTNVSTKEMRQPGGMKYIEQAIEKLSKRHAE
HIKLYGSDNDMRLTGRHETASMTAFSSGVANRGSSIRIPRSVAKEGYGYFEDRRPASNIDPYLVTGIMCE
TVCGAIDNADMTKEFERESS
```
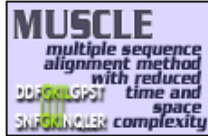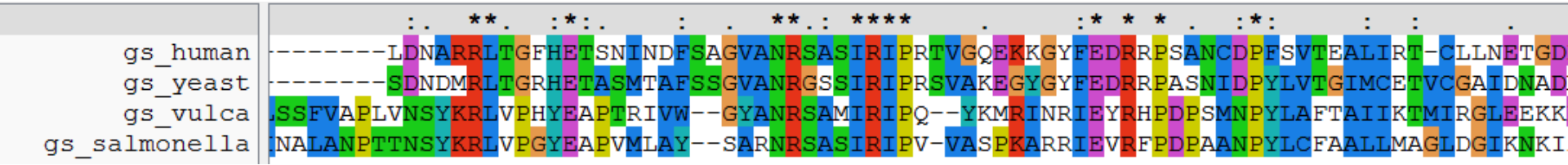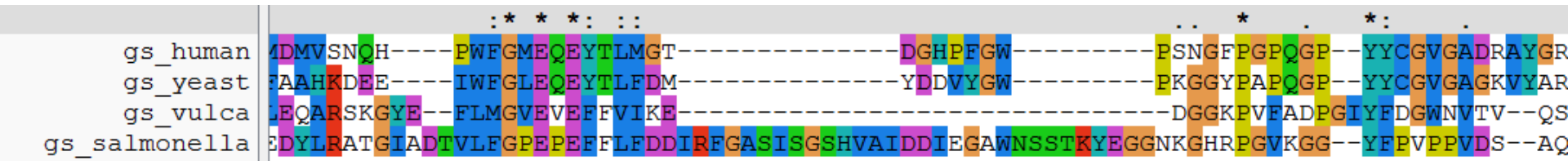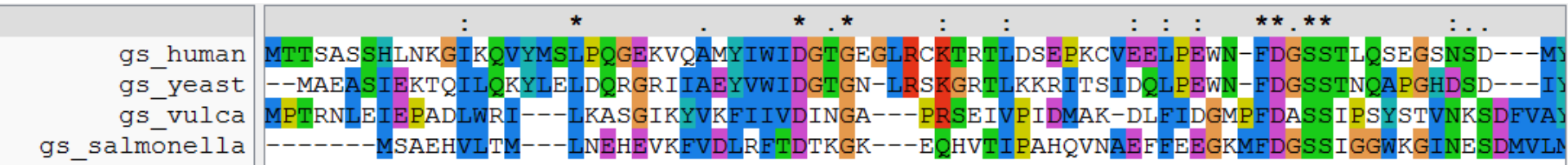
http://www.ebi.ac.uk/Tools/msa/muscle/

```
>gs_human gi|74271837|ref|NP_001028216.1| glutamine synthetase [Homo sapiens]
MTTSASSHLNKGIKQVYMSLPQGEKVQAMYIWIDGTGEGLRCKTRTLDSEPKCVEELPEW
N-FDGSSTLQSEGSNSD---MYLVPAAMFRDPFRKDPNKLVLCEVFKYNRRPA-ETNLRH
TCKRIMDMVSNQH----PWFGMEQEYTLMGT--------------DGHPFGW--------
-PSNGFPGPQGP--YYCGVGADRAYGRDIVEAHYRACLYAGVKIAGTNAEVMPA-QWEFQ
IGPCEGISMGDHLWVARFILHRVCEDFGVIATFDPKPIPGNWNGAGCHTNFSTKAMREEN
GLKYIEEAIEKLSKRHQYHIRAYDPKGG---------LDNARRLTGFHETSNINDFSAGV
ANRSASIRIPRTVGQEKKGYFEDRRPSANCDPFSVTEALIRT-CLLNETGDEP-------
------------------------------------------------------------
--------------FQYKN-----

>gs_yeast gi|330443748|ref|NP_015360.2| Gln1p [Saccharomyces cerevisiae S288c]
--MAEASIEKTQILQKYLELDQRGRIIAEYVWIDGTGN-LRSKGRTLKKRITSIDQLPEW
N-FDGSSTNQAPGHDSD---IYLKPVAYYPDPFRRGDNIVVLAACYNNDGTPN-KFNHRH
EAAKLFAAHKDEE----IWFGLEQEYTLFDM--------------YDDVYGW--------
-PKGGYPAPQGP--YYCGVGAGKVYARDMIEAHYRACLYAGLEISGINAEVMPS-QWEFQ
VGPCTGIDMGDQLWMARYFLHRVAEEFGIKISFHPKPLKGDWNGAGCHTNVSTKEMRQPG
GMKYIEQAIEKLSKRHAEHIKLYG------------SDNDMRLTGRHETASMTAFSSGV
ANRGSSIRIPRSVAKEGYGYFEDRRPASNIDPYLVTGIMCETVCGAIDNADMT-------
------------------------------------------------------------
--------------KEFERESS—

>gs_vulca gi|307594850|ref|YP_003901167.1| glutamine synthetase [Vulcanisaeta distributa DSM 14429]
MPTRNLEIEPADLWRI---LKASGIKYVKFIIVDINGA---PRSEIVPIDMAK-DLFIDG
MPFDASSIPSYSTVNKSDFVAYVDPRAVYVEYWQDGKVADVFTMVSDIADKPS-PLDPRR
VLNDALEQARSKGYE--FLMGVEVEFFVIKE--------------------------
--DGGKPVFADPGIYFDGWNVTV--QSQFMKELITAIADAGINYTKTHHEVAPS-QYEVN
IGATDPLRLADQIVYFKIMAKDIARKYGLVATFMPKPFWGV-NGSGAHTHIS---VWKDG
KNLF-QSSTGKITEECGYAISAILSNARALSSFVAPLVNSYKRLVPHYEAPTRIVW--GY
ANRSAMIRIPQ--YKMRINRIEYRHPDPSMNPYLAFTAIIKTMIRGLEEKKEPPPPTEEV
AYELA--NALETP---ATLEDTLK--ELSKSFLATE--LPSELVNAYIKIKQNEWEDYLT
NVGPWEKTWNIITQWEYNKYLVTA

>gs_salmonella gi|16767272|ref|NP_462887.1| glutamine synthetase [Salmonella enterica
-------MSAEHVLTM---LNEHEVKFVDLRFTDTKGK---EQHVTIPAHQVNAEFFEEG
KMFDGSSIGGWKGINESDMVLMPDASTAVIDPFFADSTLIIRCDILEPGTLQGYDRDPRS
IAKRAEDYLRATGIADTVLFGPEPEFFLFDDIRFGASISGSHVAIDDIEGAWNSSTKYEG
GNKGHRPGVKGG--YFPVPPVDS--AQDIRSEMCLVMEQMGLVVEAHHHEVATAGQNEVA
TRFNTMTKKADEIQIYKYVVHNVAHRFGKTATFMPKPMFGD-NGSGMHCHMS---LAKNG
TNLFSGDKYAGLSEQALYYIGGVIKHAKAINALANPTTNSYKRLVPGYEAPVMLAY--SA
RNRSASIRIPV-VASPKARRIEVRFPDPAANPYLCFAALLMAGLDGIKNKIHPGEAMDKN
LYDLPPEEAKEIPQVAGSLEEALNALDLDREFLKAGGVFTDEAIDAYIALRREEDDRVRM
TPHP----------VEFELYYSV-
```

gs_human     MTTSASSHLNKGIKQVYMSLPQGEKVQAMYIWIDGTGEGLRCKTRTLDSEPKCVEELPEWN-FDGSSTLQSEGSNSD---MY
gs_yeast     --MAEASIEKTQILQKYLELDQRGRIIAEYVWIDGTGN-LRSKGRTLKKRITSIDQLPEWN-FDGSSTNQAPGHDSD---IY
gs_vulca     MPTRNLEIEPADLWRI---LKASGIKYVKFIIVDINGA---PRSEIVPIDMAK-DLFIDGMPFDASSIPSYSTVNKSDFVAY
gs_salmonella -------MSAEHVLTM---LNEHEVKFVDLRFTDTKGK--EQHVTIPAHQVNAEFFEEGKMFDGSSIGGWKGINESDMVLM

gs_human     MDMVSNQH----PWFGMEQEYTLMGT--------------DGHPFGW---------PSNGFPGPQGP--YYCGVGADRAYGR
gs_yeast     FAAHKDEE----IWFGLEQEYTLFDM--------------YDDVYGW---------PKGGYPAPQGP--YYCGVGAGKVYAR
gs_vulca     LEQARSKGYE--FLMGVEVEFFVIKE----------------DGGKPVFADPGIYFDGWNVTV--QS
gs_salmonella EDYLRATGIADTVLFGPEPEFFLFDDIRFGASISGSHVAIDDIEGAWNSSTKYEGGNKGHRPGVKGG--YFPVPPVDS--AQ

gs_human     -------LDNARRLTGFHETSNINDFSAGVANRSASIRIPRTVGQEKKGYFEDRRPSANCDPFSVTEALIRT-CLLNETGD
gs_yeast     -------SDNDMRLTGRHETASMTAFSSGVANRGSSIRIPRSVAKEGYGYFEDRRPASNIDPYLVTGIMCETVCGAIDNAD
gs_vulca     LSSFVAPLVNSYKRLVPHYEAPTRIVW---GYANRSAMIRIPQ--YKMRINRIEYRHPDPSMNPYLAFTAIIKTMIRGLEEKK
gs_salmonella NALANPTTNSYKRLVPGYEAPVMLAY--SARNRSASIRIPV-VASPKARRIEVRFPDPAANPYLCFAALLMAGLDGIKNKI

ClustalW, JalView

# Determination of protein structure

X-ray crystallography (144K in PDB)
- need crystals

Nuclear Magnetic Resonance (NMR) (12K)
- proteins in solution
- lower size limit (600 aa)

Electron microscopy (7K)
- Low resolution (>5A)

# Determination of protein structure


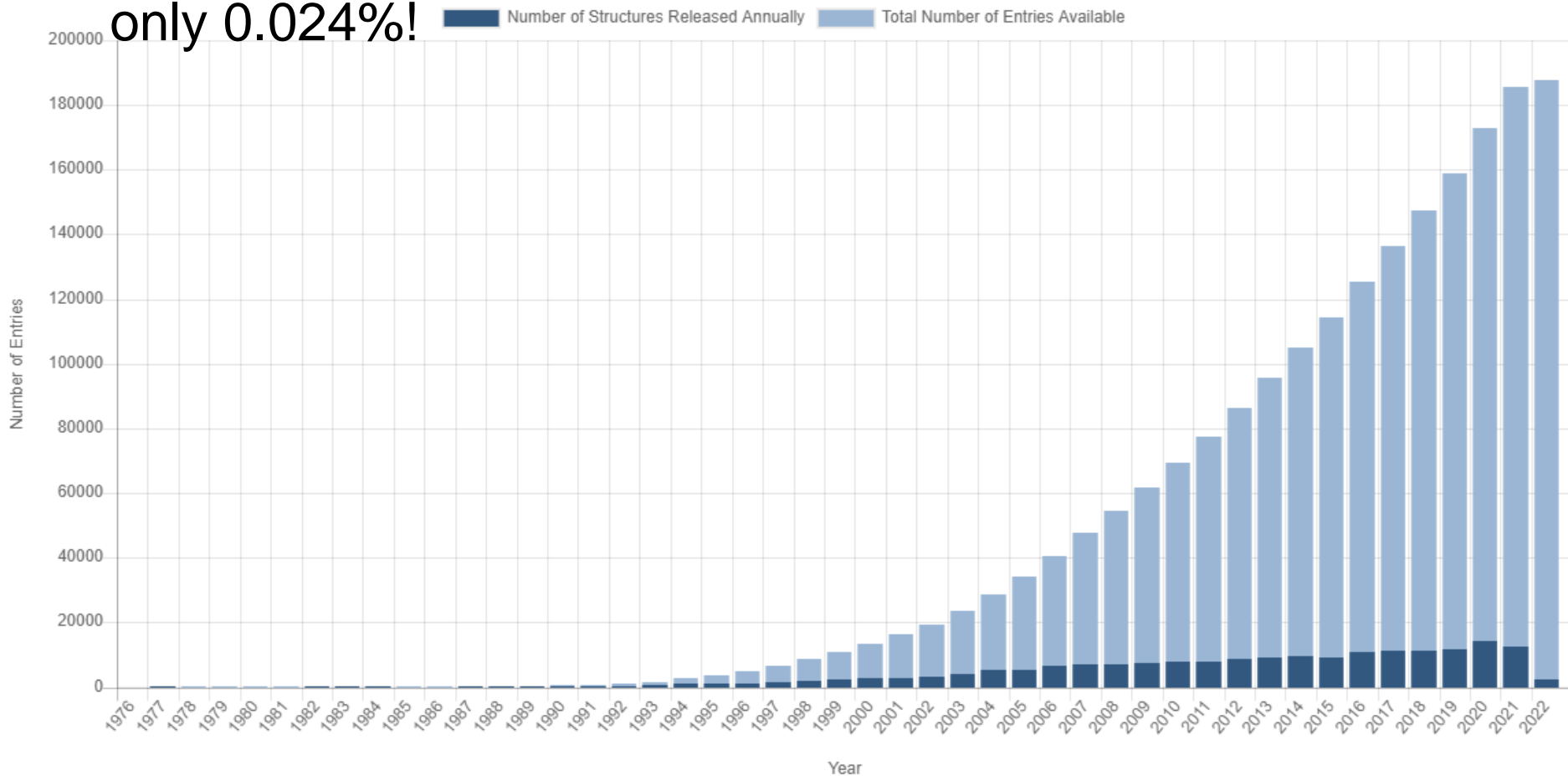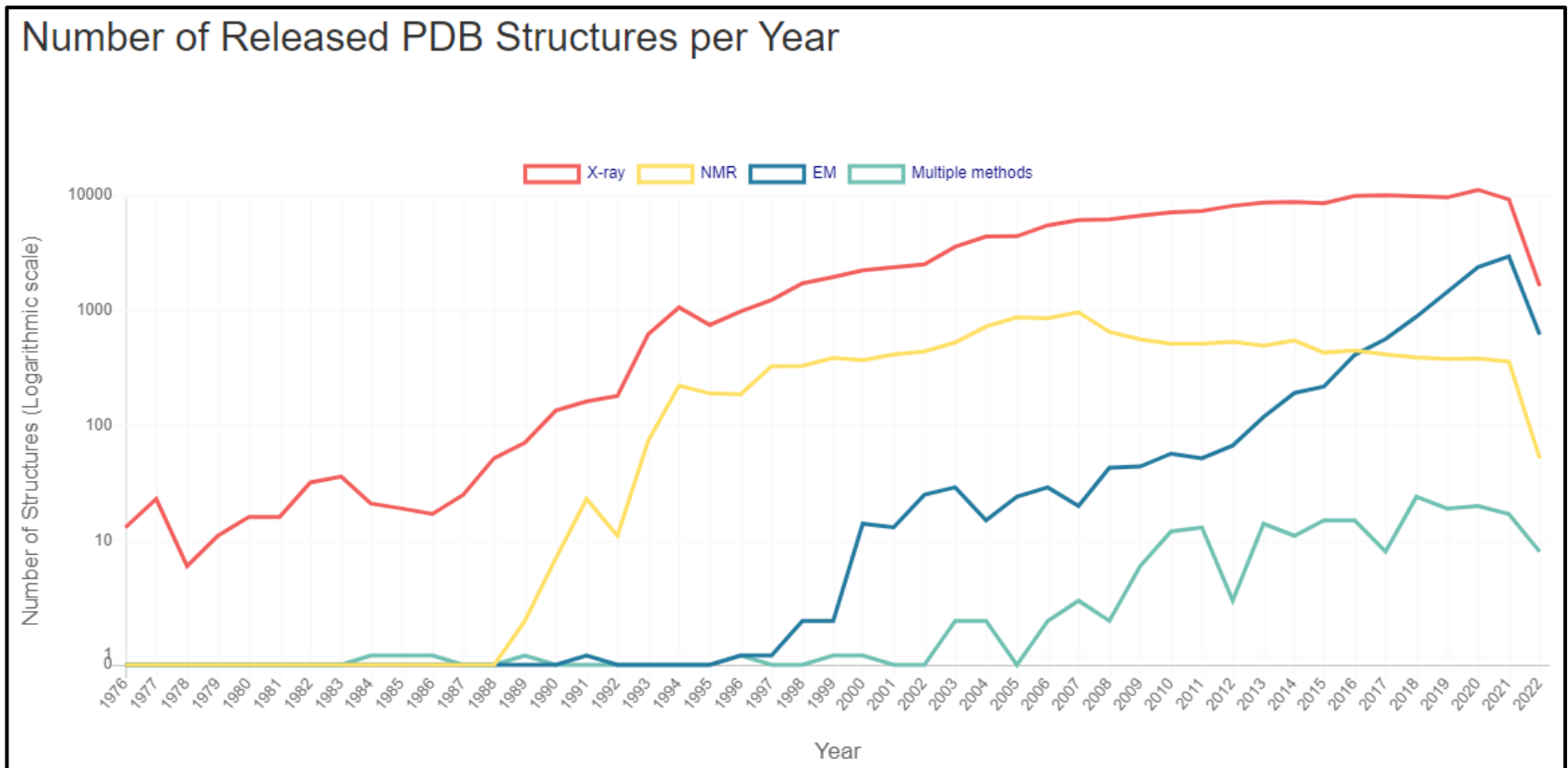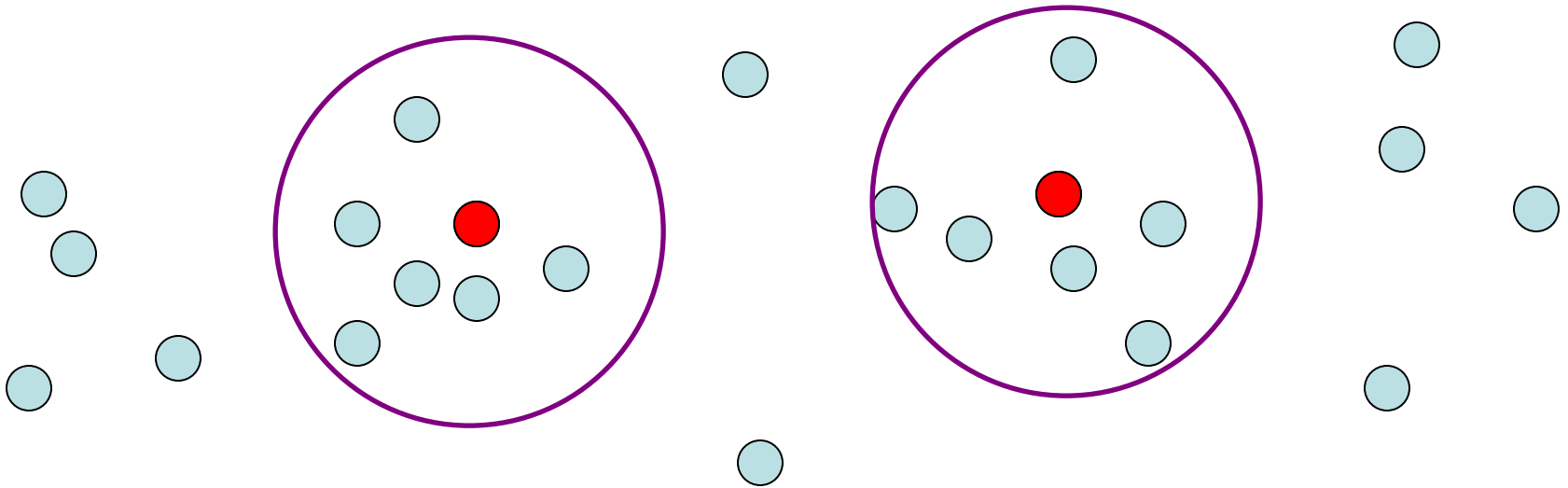
resolution 2.4 A

# Determination of protein structure



resolution 2.4 A

# Structural genomics

Currently: 187K protein 3D structures
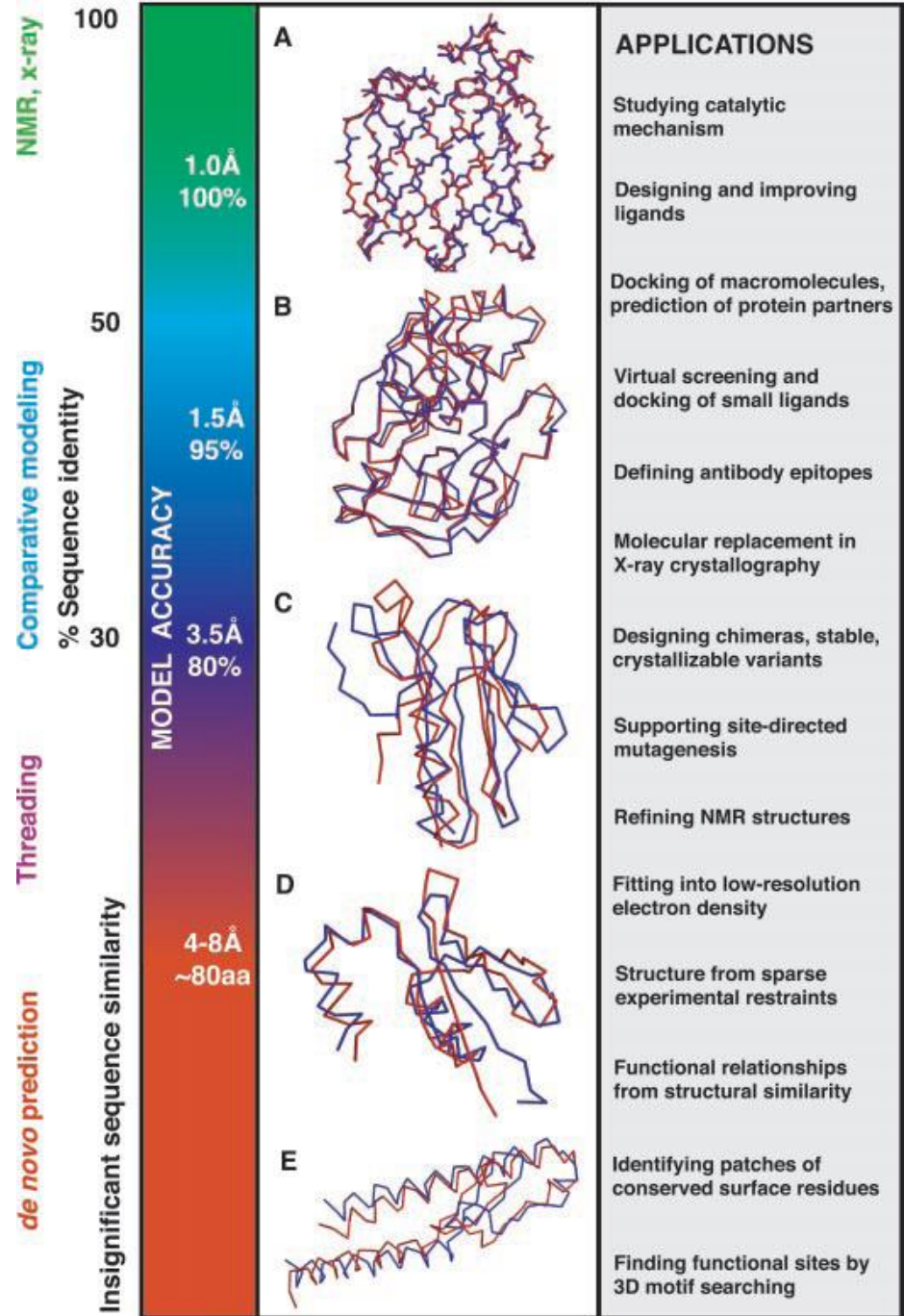from around 55K sequences in UniProt (how do I know?)
226M sequences in UniProt

only 0.024%!

# Structural genomics

Currently: 187K protein 3D structures
from around 55K sequences in UniProt (how do I know?)
226M sequences in UniProt

    only 0.024%!

# Structural genomics

Currently: 187K protein 3D structures
from around 55K sequences in UniProt (how do I know?)
226M sequences in UniProt

only 0.024%!



50% sequences covered (25% in 1995)

# Relation between sequence identity and accuracy/applications

Predicted structure (red) and real (blue)

From:
Baker and Sali (2001)
*Science*

# Homology modelling
## Applications: target design

Query sequence

| | G | K | |
|---|---|---|---|

similar to

| | L | G | |
|---|---|---|---|



known 3D

model 3D by homology

# Homology modelling
## Applications: fit to low res 3D

Query sequence 1

Query sequence 2

low resolution 3D
(electron microscopy)

# Homology modelling
# Phyre

Mike Sternberg  http://www.sbg.bio.ic.ac.uk/phyre2/

Kelley et al (2000) *J Mol Biol*
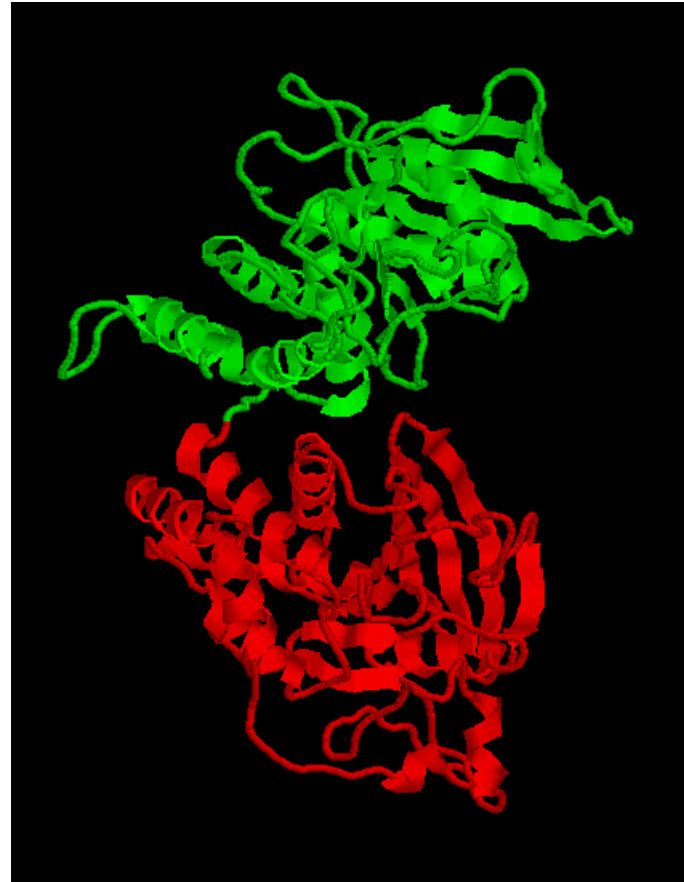Kelley et al (2015)
*Nature Protocols*



Processing time can be hours

# Domains

Protein domains are structural units (average 160 aa) that share:

Function
Folding
Evolution

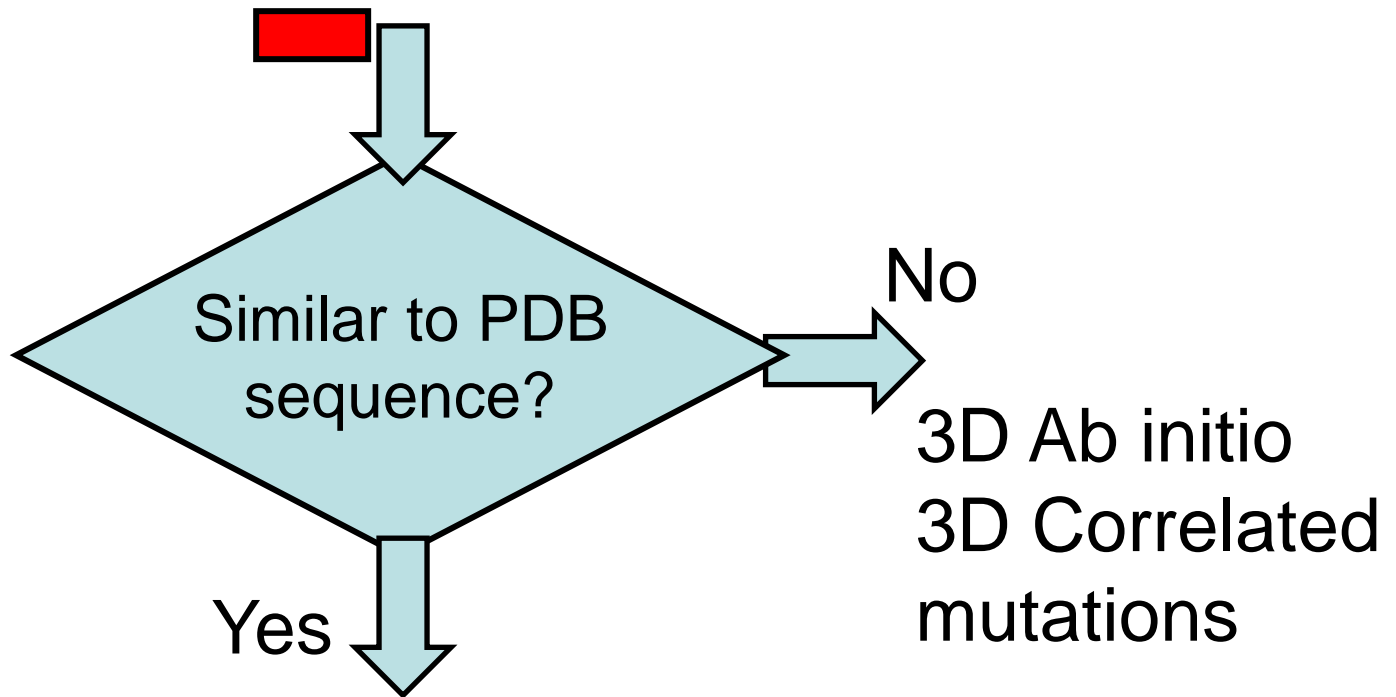Proteins normally are multidomain (average 300 aa)

# Domains

Protein domains are structural units (average 160 aa) that share:

Function
Folding
Evolution

Proteins normally are multidomain
(average 300 aa)

# Domains

Query Sequence

Predict domains

Cut

Similar to PDB sequence?

No

3D Ab initio
3D Correlated mutations

Yes

Protein structure modeling by homology

# Protein structure prediction
# Ab initio

Explore conformational space

Limit the number of atoms

Break the problem into fragments of sequence

Optimize hydrophobic residue burial and pairing of beta-strands

Limited success…

# Protein structure prediction

## Correlated mutations



Your query sequence

contact in 3D

correlated

https://commons.wikimedia.org/wiki/File:Correlated_mutation.png

# Protein structure prediction
# Combined

Homology to solved structures

Correlated sequence variation in homologs

Generation of a structure following physical constraints

# Protein structure prediction

**AlphaFold:** DeepMind, Google



Input sequence

Demis Hassabis



Jumper et al (2021) Nature

# Protein structure prediction

**AlphaFold:** DeepMind, Google

# Protein structure prediction

**AlphaFold:** DeepMind, Google

**Precomputed models**:
UniProt
https://alphafold.ebi.ac.uk/
(limited to model organisms)

**Colab notebook** (simplified version / limited server / takes hours)
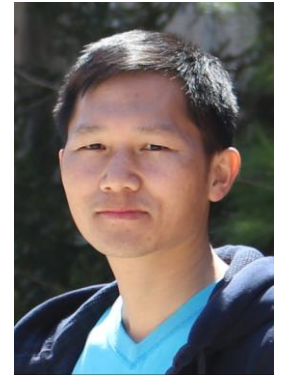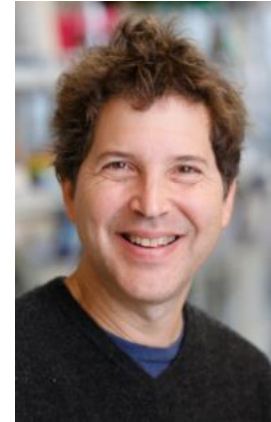
**Source code** (Needs 3 Tb disk space)

# Protein structure prediction

**trRosetta:** David Baker & Jianyi Yang

Needs large multiple sequence alignments to predict contacts

Predictions available for all PFAM domains
Example:
https://www.ebi.ac.uk/interpro/entry/pfam/PF07887/rosettafold/

# Run online at https://yanglab.nankai.edu.cn/trRosetta/

Du et al (2021) *Nature Protocols*

# Protein structure prediction

**C-I-Tasser:** Yang Zhang



Run online at
https://zhanggroup.org/C-I-TASSER/

Zheng et al (2021) *Cell Reports Methods*