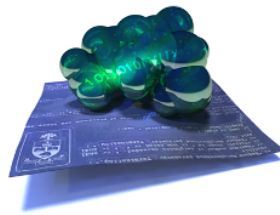


A
BIOINFORMATICS
COURSE

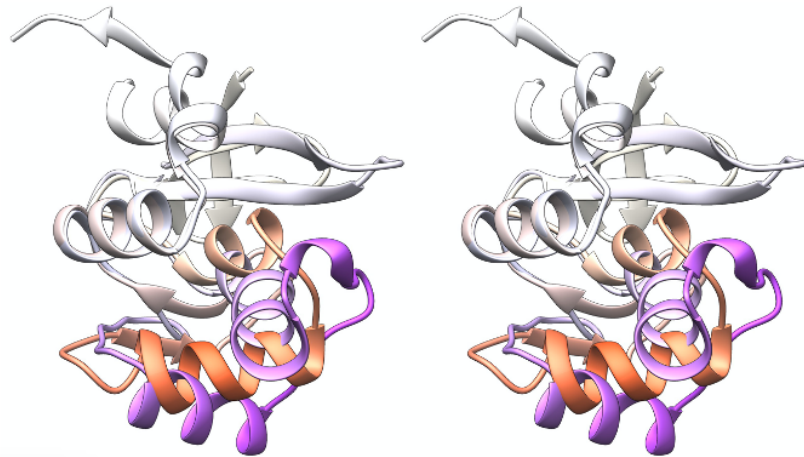
HOMOLOGY MODELLING



BORIS STEIPE

DEPARTMENT OF BIOCHEMISTRY – DEPARTMENT OF MOLECULAR GENETICS
UNIVERSITY OF TORONTO

STRUCTURAL SIMILARITY



```

1 -----QIYSARYSGVDVY 13
      |.|
1 ALBEEVRYLGPHEIPLTLTRDSETCGHFLKHFPLIQQY----- 40
14 EFTIHSSTGSIMKRKDDAVNATHILKAANFAKAKRTRILEKRVLKETHEKV 63
   |.|.|.|.|.|.|.|.|.|.|.|.|.|.|.|.|.|.|.|.|.|.|.|.|.|.|.|.|.|
41 ---HDTGNINETNFDSPFT-----DEERNKL 63
64 QGGFGKYQGT-----NVPLNIKQLAEKFSVDQKPLDF----- 99
   ..:|...| |:|...|...:|...:|...:|...:|...:|...:|
64 LAHYGIAVNTDRGELNILEKCLQLLNMLNLFLGLFQDAFEFKEPETDQD 113
100 ----- 99
114 EEDPSSHSLPEN 125
  
```

1BM8 - APSES domain
(Yeast *Mbp1*)

2XFV - APSES domain
(Yeast *Swi6*)

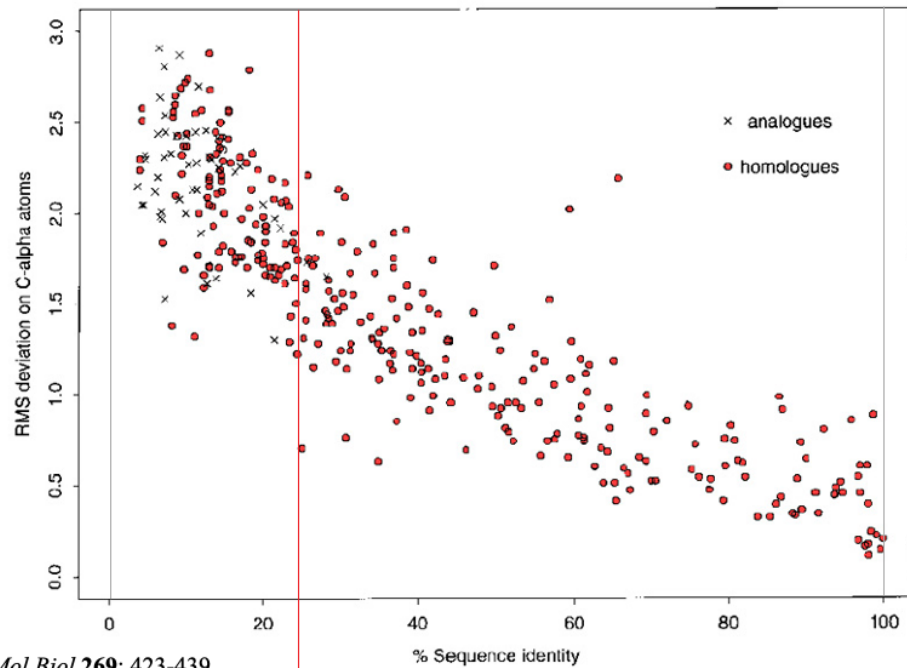
1BM8/2XFV: No significant sequence similarity

Structures can be similar despite sequences being dissimilar.

In the case of the *Mbp1* APSES domain compared to its homologue *Swi6*, an indel after the DNA binding helix has shifted the orientation of two helices (red and magenta) and destroyed the ability of the *Swi6* ancestor to bind DNA, the domain took on different functions and has diverged beyond any detectable sequence similarity. Nevertheless, the structures at the termini of the domains can be perfectly superimposed (white).

Homologous proteins retain structural similarity even if they are highly diverged. As a corollary, once homology can be established between two sequences, based on sequence similarity, it is virtually certain that the structure of one protein can be modelled from knowledge of the structure of the other

STRUCTURAL SIMILARITY



Russell *et al.* (1997) *J Mol Biol* **269**: 423-439

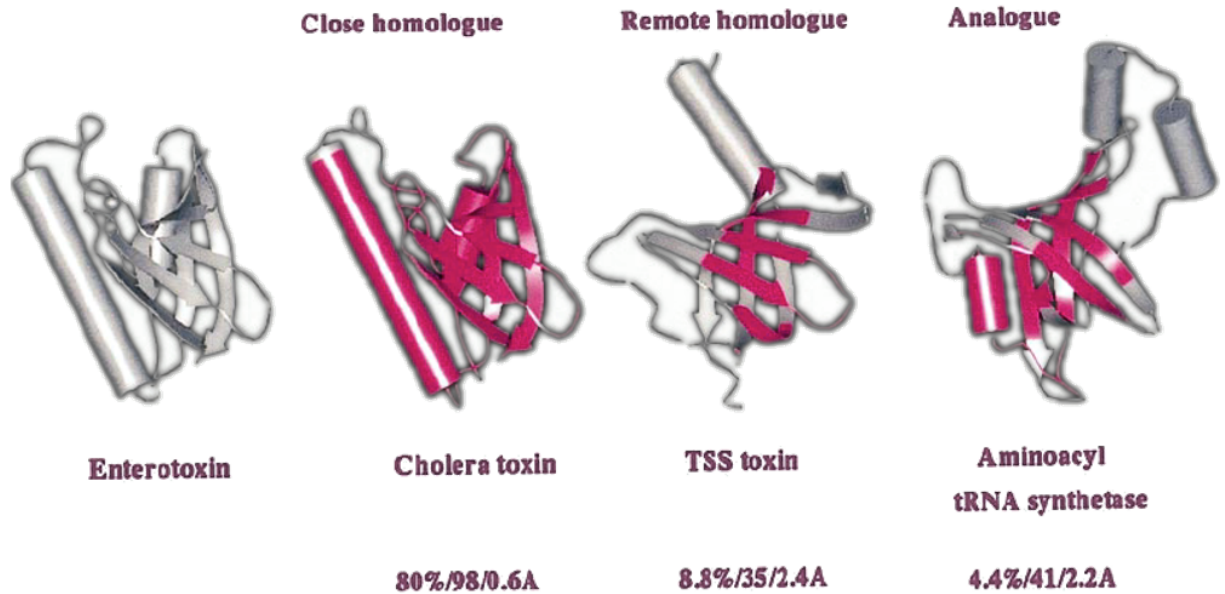
25%ID

Proteins that diverge in evolution maintain their global fold !

Note the 25% ID level, which we usually take as the cutoff for inferring homology from sequence similarity.

Even at 100% sequence identity (*i.e.* the structure of the *same* protein observed under different conditions) structures can vary by 0.3Å RMSD or more. But up to ~50% identity, structure is conserved at less than 1.0Å RMSD – less than the length of a carbon-hydrogen bond! At the extreme of dissimilarity, even in the complete absence of sequence identity, structural similarity does not drop below 3.0Å for homologous structures. However, at the left hand of this distribution, homologous and analogous folds cannot be distinguished by sequence or structure similarity alone.

HOMOLOGY AND ANALOGY



Fold space is finite, and structures that have a similarly folded core can have arisen by convergent evolution. These are not homologous folds, they are **analogous**.

This example shows close and remote homologues of enterotoxin, and a structurally similar aminoacyl tRNA synthetase structure. Numbers in brackets are %ID / number of matched residues / RMSD of superposition. By these metrics, the tRNA synthetase is a better match to enterotoxin than the TSS toxin. But it is not a homologue.

cf. Russell *et al.* (1997) Recognition of analogous and homologous protein folds: analysis of sequence and structure conservation. *J Mol Biol.* **269**:423-439.

HOMOLOGY MODELLING

How to build a homology model (in three easy steps).

- 1: Align a **target** sequence with the sequence of a suitable structure **template**.
- 2: **Replace** the coordinates of sidechains in the template structure file with coordinates for target residue sidechains. This is the **model**.
- 3: That's all. *(except you can still mess it up...)*

Target: the sequence for which you want to obtain a modelled structure.

Template: a sequence that is homologous to the target and for which the structure is known.

Model: Structural model of the target, obtained by replacing template sidechains with target sidechains

Sometimes researchers add a third step to homology modelling: energy refinement of the model to adopt more reasonable atom-atom contacts, bond-lengths and angles *etc.* While this practice is popular, it is likely to do more harm than good: usually the RMSD between the model and the true structure gets worse.

Thus energy refinement of homology models is in general an example of Cargo Cult science.

STRUCTURE CONSERVATION

	E-E.coli	[...]	IKTRFAPSP	TYLHVGGART	[...]	EQMAKGE---	KPRYDGRC	[...]	AHVS	MINGDDGK	KL	SKRH					
Target	E-P.putida	[...]	VRTR	IAPSP	TGDP	PHVGTAYIA	[...]	EQQARGE---	T	PRYDGRA	[...]	CYMP	LLRN	PK	SKLSKRK		
Template	Q-E.coli	[...]	VHTR	FPPE	PN	GYLH	IGHAKSI	[...]	TLTQ	PGKNSPYR	DRSVEEN	[...]	YEF	SRL	NLEY	TVMSKRK	
Guiding sequences	Q-Fly	[...]	VHTR	FPPE	PN	GILH	IGHAKAI	[...]	FNP	KPS---	PWRER	PIEES	[...]	WEY	GRL	NMNY	ALVSKRK
	Q-Human	[...]	VRTR	FPPE	PN	GILH	IGHAKAI	[...]	HNT	LPS---	PWRDR	PMEES	[...]	WEY	GRL	NLHY	AVVSKRK
	E-Fly	[...]	VVVR	FPPE	AS	GYLH	IGHAKAA	[...]	QRVE---	SANRS	NSVEKN	[...]	WSY	SRL	NMTN	TVLSKRK	
	E-Human	[...]	VTVR	FPPE	AS	GYLH	IGHAKAA	[...]	QRIE---	SKHRK	NPIEKN	[...]	WEY	SRL	NLN	NTVLSKRK	
	E-Yeast	[...]	VVTR	FPPE	PS	GYLH	IGHAKAA	[...]	DGVA---	SARR	DRSVEEN	[...]	WDF	FARI	NFV	RTLLSKRK	

Ligand binding sites | | | | |

QRS *E. coli* vs. ERS *P. putida*: ~ 19% pairwise ID

- Many regions are expected to be highly conserved in structure.
- Some changes are straightforward to model.

Alignment of aminoacyl t-RNA synthetase sequences, and consequences of the types of changes we observe for modeling 3D-structure. In this example, we consider the target of *P. putida* glutamyl tRNA ligase model on the structure of *E. coli* glutaminyl tRNA ligase.

Side chains which can be changed by deleting atoms (I→V, I→A, N→A, T→S), or which merely require changing chemical elements (N→D, even L→Q), are straightforward to model and don't require changing coordinates.

STRUCTURE CONSERVATION

	E-E.coli	[...]	IKTRFAPSPG YLHVGGART A	[...]	EQMAKGE----KPRYDGR C	[...]	AHVSMINGDDGK LSKR H
Target	E-P.putida	[...]	VRTR IAPSPGDPHVGTAYIA	[...]	EQQARGE----TPRYDGR A	[...]	CYMP LLRNPKSK LSKR K
Template	Q-E.coli	[...]	VHTRFPPEP NGYLHIGHAK SI	[...]	TLTQPGKNSPYRDRSVEEN	[...]	YEF SRL-NLHYTVMSKR K S
Guiding sequences	Q-Fly	[...]	VHTRFPPEP NGILHIGHAKAI	[...]	FNP KPS---PWRERPIEES	[...]	WEY GRL-NMNYALVSKR K
	Q-Human	[...]	VRTRFPPEP NGILHIGHAKAI	[...]	HNTLPS---PWRDRPMEES	[...]	WEY GRL-NLHYAVVSKR K
	E-Fly	[...]	VVVRFPPEASG YLHIGHAKAA	[...]	QRVE----SANRSNSVEKN	[...]	WSYSRL-NMTNTVLSKR K
	E-Human	[...]	VTVRFPPEASG YLHIGHAKAA	[...]	QRIE----SKHRKNPIEKN	[...]	WEYSRL-NLNNTVLSKR K
	E-Yeast	[...]	VVTRFPPEPSG YLHIGHAKAA	[...]	DGVA----SARRDRSVEEN	[...]	WDFARI-NFVRTLLSKR K

Ligand binding sites | || | ||

QRS *E. coli* vs. ERS *P. putida*: ~ 19% pairwise ID

■ How would sidechain rotamers be modeled?

- conserved dihedral angles
- preferred rotamers
- DEE (Dead End Elimination theorem) for global consistency.

For model sidechains that are larger than template sidechains (H→R, S→I, T→Q, V→K ...), we need to decide on the correct geometry. Looking up preferred conformations in “rotamer” dictionaries helps make the decision.

STRUCTURE CONSERVATION

	E-E.coli	[...]	IKTRFAPSPTGYLHVGGARTA	[...]	EQMAKGE----	KPRYDGRC	[...]	AHVSMINGDDGKKLSKRH
Target	E-P.putida	[...]	VRTRIAPSPTGDPHVGTAYIA	[...]	EQQARGE----	TPRYDGRA	[...]	CYMPLLRNPKSKLSKRK
Template	Q-E.coli	[...]	VHTRFPPEPNGYLHIGHAKSI	[...]	TLTQFGKNSPYRDRSVEEN	[...]	YEF SRL-NLEYTVMSKRK	S
Guiding sequences	Q-Fly	[...]	VHTRFPPEPNGILHIGHAKAI	[...]	FNP KPS---	PWRERPIEES	[...]	WEYGRL-NMNYALVSKRK
	Q-Human	[...]	VRTRFPPEPNGILHIGHAKAI	[...]	HNTLPS---	PWRDRPMEES	[...]	WEYGRL-NLHYAVVSKRK
	E-Fly	[...]	VVVRFPPEASGYLHIGHAKAA	[...]	QRVE----	SANRSNSVEKN	[...]	WSYSRL-NMTNTVLSKRK
	E-Human	[...]	VTVRFPPEASGYLHIGHAKAA	[...]	QRIE----	SKHRKNPIEKN	[...]	WEYSRL-NLNNTVLSKRK
	E-Yeast	[...]	VVTRFPPEPSGYLHIGHAKAA	[...]	DGVA----	SARRDRSVEEN	[...]	WDFARI-NFVRTLLSKRK
			Ligand binding sites					

How would you (or should you even) model indels?

- Where should the insertion be placed?
- What is the conformation of the new residues?
- Which residues should be deleted?
- How many additional residues need to change conformation?

Indels cannot be modelled with any certainty.

Homology modelling accuracy
is determined by sequence
alignment accuracy.

The goal for an alignment to be used in homology modelling is a sequence to structure mapping, not sequence to sequence!

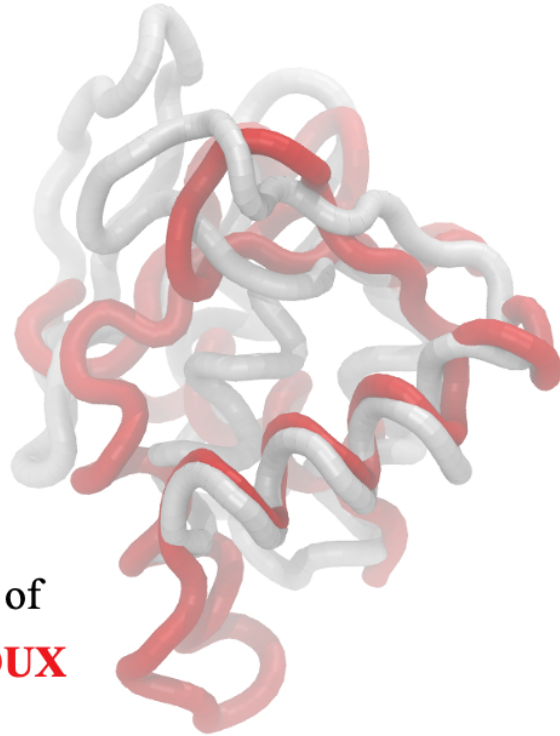
A sequence to structure mapping asks: which residue of the target should replace which residue of the model – this is not necessarily the same as asking which residue of the target is related by evolution to which residue of the model.

SUPERPOSITION AND ALIGNMENT

The coordinates of two proteins can be superimposed in space.

An alignment may be derived from a **superposition**, by correlating residues that are close in space.

An **optimal sequence alignment** may lead to a different correlation ...



Superposition of
1BM8 and **1DUX**

Alignment is not superposition. These are different procedures and they have different objectives. One can *derive* an alignment from a superposition, by aligning e.g. all residue pairs for which the C^α atom distance fall below a threshold, such as 1.9\AA , *i.e.* half the average C^α - C^α separation in proteins.

INDELS

Indels may be structurally accommodated some distance away from the site of their sequence insertion.

Example - a two residue deletion in GFP

Sequence alignment (shows what happened)

```
GFP  DGSVQLADHYQQNTPIGDGPVLLP
      ||...|.:. . : : : : : . |||.||
RFP  DGGHYLVEF--KSIYMAKKPVQLP
```

Structure superposition (shows how it's accommodated)

```
GFP  DGSVQLADHYQQNTPIGDGPVLLP
      ||...|.:. . . |||.||
RFP  DGGHYLVEFKSIYMAK..KPVQLP
```



Sequence alignment and superposition have different objectives. Alignment – based on an evolutionary model – recovers information on an evolutionary event. Superposition shows how the event has been structurally accommodated.

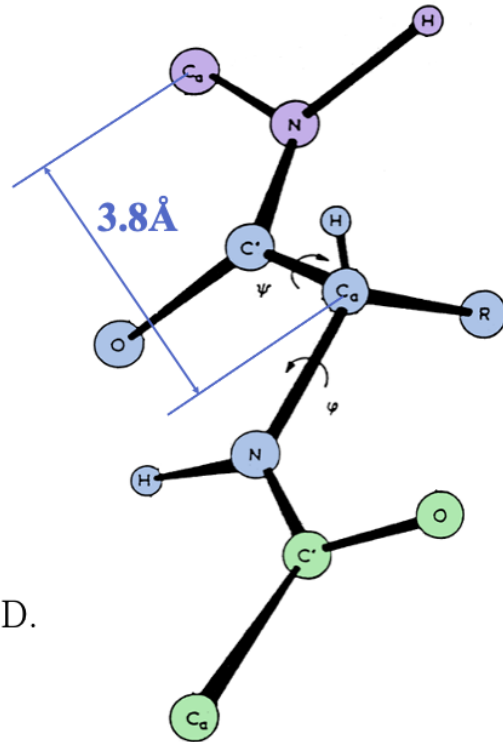
To recover an alignment from a superposition, align each residue with the one it is closest to in the superimposed pair of structures.

ALIGNMENT ERROR AND MODEL ACCURACY

A shift in alignment of just *one* residue corresponds to an error in the modelled structure of about 4 Å.

Nothing you can do AFTER the alignment will fix this error. Energy refinements (in general) can't correct errors as large as that.

Remember that RMSDs of homologous proteins are less than 3Å, regardless of %ID.



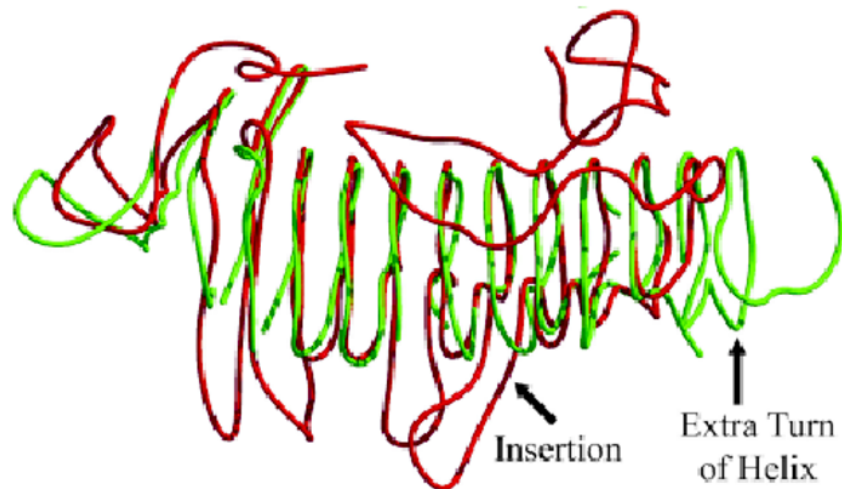
Indels can cause alignment errors. Alignment errors cannot be corrected through modelling.

ALIGNMENT IS THE MOST IMPORTANT FACTOR FOR ACCURACY

Sequence alignment is the limiting step for homology model accuracy. No amount of energy minimization or other computational procedure will put a misaligned residue in the right place just because it may be more plausible there!
(And sometimes *plausible* still means *wrong*.)

TRUE structure

Homology model



From:

HOMSTRAD @ CASP4:

Williams MG et al. (2001) *Proteins Suppl.5*: 92-97

In this example, the alignment "slipped" on an inserted loop in the beta-helix ("Insertion") and added an additional loop at its end. This makes more than half of the model "wrong", and there is *no way* to resolve this error by computational means, short of outright *ab initio* structure prediction.

It's absolutely crucial to get the alignment right.

ALIGNMENT IS THE MOST IMPORTANT FACTOR FOR ACCURACY

Consequently:

use the best alignment that is available, in situations where there is any doubt (i.e. **always** in the presence of indels)

Use a carefully computed **multiple sequence alignment**

Include your **target** and your selected **template** sequence

Include many relevant **homologues** to resolve ambiguities

Manually edit results for plausibility of the target/template alignment, consider template secondary structure

(cf. SAS: <https://www.ebi.ac.uk/thornton-srv/databases/sas/> or
PDBSUM: <https://www.ebi.ac.uk/pdbsum>)

Extract the input pairwise alignment from the multiple alignment

Your alignment should strive for the most plausible sequence-to-structure mapping.

After producing an MSA from carefully selected sequences, extract the pairwise alignment of target and template by copying the two rows, and removing all gap characters that are present in **both** sequences.

Note that PSI-BLAST is probably not the best tool to search for related sequences for homology modeling: the goal is not to have a comprehensive set of homologues, but to include sequences that improve the accuracy of the MSA; that will generally require high levels of sequence similarity. In contrast, the use-case for PSI-BLAST is detecting very distant relationships.

As a rule: **homologous sequences that are not descendants of the Last Common Ancestor of target and template, are unlikely to be useful for the alignment.**

MODELLING INDELS

- Comparisons of alignments and structures demonstrate that *uniform gap penalty assumptions* are ***not biological***.
- Indels are most often observed in loops, less often in secondary structure elements
- When they do not occur in loops, there is frequently a maintenance of helical or strand properties and the indel is structurally accommodated in a position different from where the insertion occurred.
- In that case, the homology model should be based on the superposition, not on the optimal sequence alignment.

Accordingly:

- Use structure-weighted gap penalties
- Use a multiple sequence alignment
- Use manual improvement of alignment

STRUCTURE-AWARE GAP PENALTIES

- Should we use **position specific, structure aware** gap penalties for general alignment, or more specifically to align for homology modeling?
- Some MSA algorithms support the use of **secondary structure masks**.
- Alternatively: use a manual **sequence alignment editor** to move gaps out of secondary structure regions, which you can identify from the *template* structure(s).

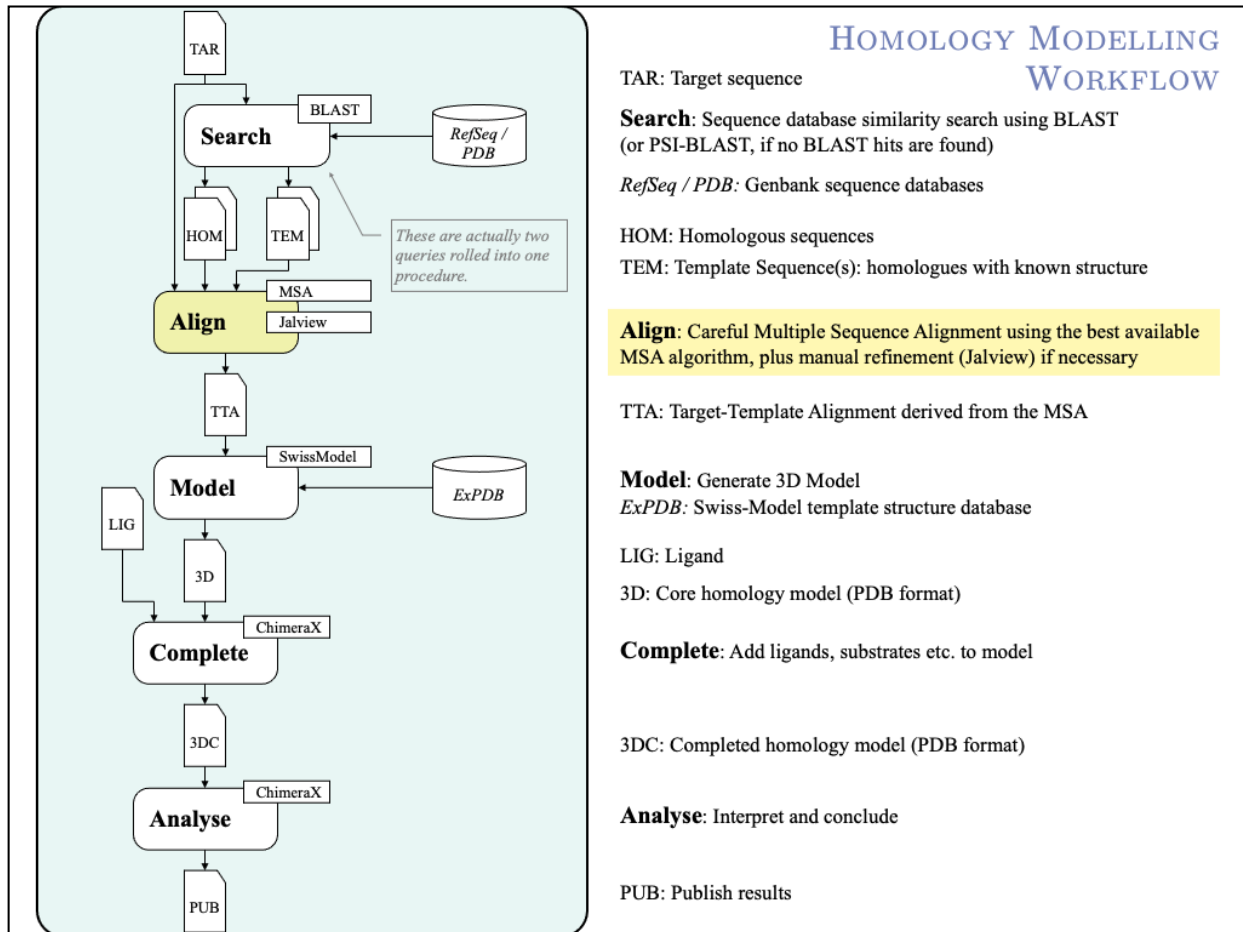
However: This may be implicitly achieved by modern (consistency based) multiple sequence alignment programs.

How to choose a template:

- 1: Choose the structure with the highest sequence similarity (smallest number of indels!)
- 2: Choose a structure with bound ligands, cofactors, or other biologically relevant modifications
- 3: Choose the highest resolution X-ray structure, (except if the resolution is worse than $\sim 3.0\text{\AA}$, then NMR structures might be just as accurate). Don't trust loops to yield correct coordinates.

If 1-3 don't all indicate the same structure, you have to balance the requirements intelligently, based on the purpose of the model. Or create several models and then intelligently cut and paste coordinates.

Include more than one template sequence in the MSA if possible: comparing a template **superposition** with the template **alignment** will be very informative in terms of alignment details.



Workflow for homology modeling sketched in SPN (Structured Process Notation).

HOMOLOGY MODELLING: ONLINE SERVERS PERFORM WELL

The screenshot shows the SWISS-MODEL web interface for a template alignment project. The main heading is "Template Alignment: 4ux5.1" created today at 04:19. The interface is divided into several sections:

- Model Results:** Shows a 3D ribbon diagram of the protein structure. Below it, a table of quality metrics is displayed:

Metric	Value
GMQE	0.71
QMEAN	-1.98
Global Quality	
QMEAN	-1.98
C β	-1.31
All Atom	-0.27
Solvation	0.03
Torsion	-1.74
- Local Quality:** A line graph showing the local quality estimate across the protein sequence.
- Comparison:** A scatter plot comparing the model's quality to a set of 500 structures.
- Sequence Alignment:** A table showing the alignment between the template (4ux5.1.A) and the model (Model_01). The alignment shows a 59.26% sequence identity and coverage. The description of the template is "TRANSCRIPTION FACTOR MBP1".

Online servers such as Swiss-Model provide convenient, high-quality homology modelling services for free. Here the APSES domain of *Cryptococcus Neoformans* Mbp1 has been modelled on the orthologous domain of *Magnaporthe oryzae*.

How to assess accuracy:

- 1: Assume all indels to be wrong.
- 2: Assume disordered portions of the template have lead to wrongly modelled sections.
- 3: Structure analysis (“threading”, “solvent accessibility”, compatibility with ligands) can point out possible alignment errors.
- 4: Homology modelling programs have tools to assess local reliability of structure.

But: there is no point in “repairing” stereochemistry, you can only review the alignment and try again once you have improved it.

Assume all indels to be wrong: *SwissModel* flags “invented” conformations (loops) by giving them a value of 0.00 in the *occupancy* column of the PDB file.

Solvent exposed sidechain conformations are unreliable because of the lack of packing constraints – except if engaged in conserved ligand interactions.

Ligands may include structurally conserved water molecules – check the template. Water molecules will not have been included by the modeling algorithm.

HOMOLOGY MODELLING: USE OF RESULTS

Biochemical inference based on 3D similarity – *i.e.* homology models

- **Bonds**
 - **Angles, plain and dihedral**
 - **Surfaces, solvent accessibility**
- } ————— **Don't bother!**
- Amino acid functions, presence in structure patterns
 - **Spatial relationship of residues to active site**
 - **Spatial relationship to other residues**
 - **Participation in function / mechanism**
 - Static and dynamic disorder
 - Electrostatics
 - *Conservation patterns (structural and functional)*
 - *Posttranslational modification sites (but not structural consequences!)*
 - *Suitability as drug target*

Just as in structure analysis, treat your model as a spatial map of features and annotations, not necessarily as a representation of accurate coordinates.

Is it possible to predict function from models? Usually not, however some functions may be incompatible with the model and thus can be excluded from consideration.

HOMOLOGY MODELLING: CARGO CULT

Some types of analysis amount to Cargo Cult:

- Modelling properties that cannot / will not be verified
- Analysing detailed geometry of the model
- Interpreting loop structures near indels
- Inferring relative domain arrangement
- Inferring structures of complexes

You need to consider whether there actually IS a use for the model coordinates – or whether the model is just intended as a pretty picture.

Prototype 1: **Analytical**

Explain mechanistic aspects of protein.

(e.g. in terms of)

- residues involved in catalysis
- global properties (like electrostatics)
- shape, relative orientation and distances of domains or subdomains
- flexibility and dynamics - e.g. hypothesizing about the rate limiting step

HOMOLOGY MODELLING: USE OF RESULTS

Prototype 2: **Comparative**

Bring conservation patterns into a spatial context in order to infer causality from (database) correlations –

- describe context specific conservation patterns and analyze these according to conserved properties;
- analyze the predicted effect of sequence variation (e.g. for engineering changes, fusing domains or predicting SNP effects);
- distinguish physiological vs. nonphysiological interactions.

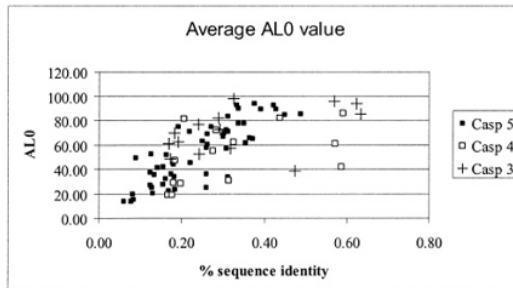
HOMOLOGY MODELLING: ENERGY MINIMIZATION

- Brings protein to lowest energy in about 1-2 minutes CPU time
- Removes atomic overlaps and unnatural strains in the structure
- Stabilizes or reinforces strong hydrogen bonds, breaks weak ones
- Efficient way of “polishing” your protein model
- but is it “**true**”???

ENERGY MINIMIZATION OFTEN MAKES MODELS WORSE

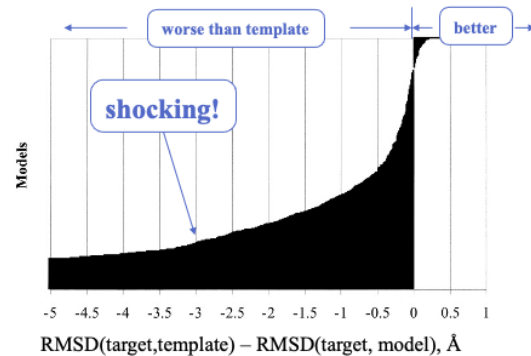
Homology modeling is getting better ...

... because:



Remote sequence similarity detection methods have improved.

... not because:



Coordinate manipulations introduce inaccuracies where there would not need to be any!

Tramontano A & Morea V (2003) Assessment of homology based predictions in CASP5
Proteins **S6**:352-368

Comparison of $\text{RMSD}(\text{target}, \text{template}) - \text{RMSD}(\text{target}, \text{model})$:

target is the true structure of the target protein that has been independently solved and used for comparison to validate the modelling procedure. If $\text{RMSD}(\text{target}, \text{template})$ is smaller than $\text{RMSD}(\text{target}, \text{model})$, this means that the true structure is more similar to the original template than to the homology model that was produced.

Structure prediction assessments have shown that often the template structure is more similar to the true structure of the target, than the model structure. This is troubling.

Something was done to the template backbone (*i.e.* energy refinement) that actually made the model more wrong than simply keeping the template as-is would have been.

The number of cases where such manipulations **improved** the model – if anything, by a tiny amount – is vanishingly small.

MODEL DATABASES MAY ALREADY CONTAIN THE STRUCTURE YOU NEED

ModBase: Database of Comparative Protein Structure Models

ModBase is a database of comparative protein structure models, calculated by our modeling pipeline ModPipe.

All available datasets are selected.

Model Details Example

Database ID: Q12121

This page makes all models and model details for one sequence available.

Sequence Model Coverage Summary for all Models of the Sequence

433 Models

- The current model is shown prominently.
- Additional models are accessible through their thumbnails.
- Use the Perform action pull-down menu for:
 - Coordinates
 - Alignments
 - Various ModBase representations
 - links to UCSF Chimera

<https://modbase.compbio.ucsf.edu>

SWISS-MODEL Repository

The SWISS-MODEL Repository is a database of annotated 3D protein structure models generated by the SWISS-MODEL homology-modelling pipeline.

The aim of the SWISS-MODEL Repository is to provide access to an up-to-date collection of annotated 3D protein models generated by automated homology modelling for relevant model organisms and experimental structures information for all sequences in UniProtKB. Regular updates ensure that target coverage is complete, that models are built using the most recent sequence and template structure databases, and that improvements in the underlying modelling pipeline are fully utilized. It also allows users to assess the quality of the models using the latest QMEAN results. If a sequence has not been modelled, the user can build models interactively via the SWISS-MODEL workspace.

Currently the repository contains 1,045,107 models from SWISS-MODEL for UniProtKB targets as well as 125,067 structures from PDB with mapping to UniProtKB.

We currently provide models for the reference proteomes of the following model organisms, based on UniProtKB release 2017_10, if you want to download a large number of models, please contact us.

	Proteome Size (Canonical sequences)	Sequences Modelled	Models	Seq Coverage	Download Metadata (Models and structures)	Download Coordinates (Models mapping to SwissProt only)
<i>Homo sapiens</i>	20,984	18,046	43,400		↓ 9.2 MB	↓ 3.4 GB
<i>Mus musculus</i>	22,249	18,683	43,743		↓ 6.4 MB	↓ 2.0 GB
<i>Caenorhabditis elegans</i>	20,006	11,673	22,858		↓ 2.2 MB	↓ 397.2 MB
<i>Escherichia coli</i>	4,306	3,382	5,980		↓ 950.7 KB	↓ 321.7 MB
<i>Arabidopsis thaliana</i>	27,561	19,211	38,296		↓ 3.3 MB	↓ 1.2 GB
<i>Drosophila melanogaster</i>	13,776	9,319	19,340		↓ 2.0 MB	↓ 416.4 MB
<i>Saccharomyces cerevisiae</i>	6,049	4,323	7,449		↓ 1.1 MB	↓ 333.3 MB
<i>Caulobacter crescentus</i>	3,720	2,831	5,188		↓ 474.7 KB	↓ 51.8 MB
<i>Mycobacterium tuberculosis</i>	3,993	3,081	5,194		↓ 665.1 KB	↓ 157.4 MB

<https://swissmodel.expasy.org/repository>

These databases provide high-quality automated models on a genome scale. As you can infer from the coloured bars of the Swiss-Model repository, approximately 50% of a given model organism proteome can be modelled with high confidence based on existing protein structures.

<http://steipe.biochemistry.utoronto.ca/abc>

B O R I S . S T E I P E @ U T O R O N T O . C A

DEPARTMENT OF BIOCHEMISTRY & DEPARTMENT OF MOLECULAR GENETICS
UNIVERSITY OF TORONTO, CANADA