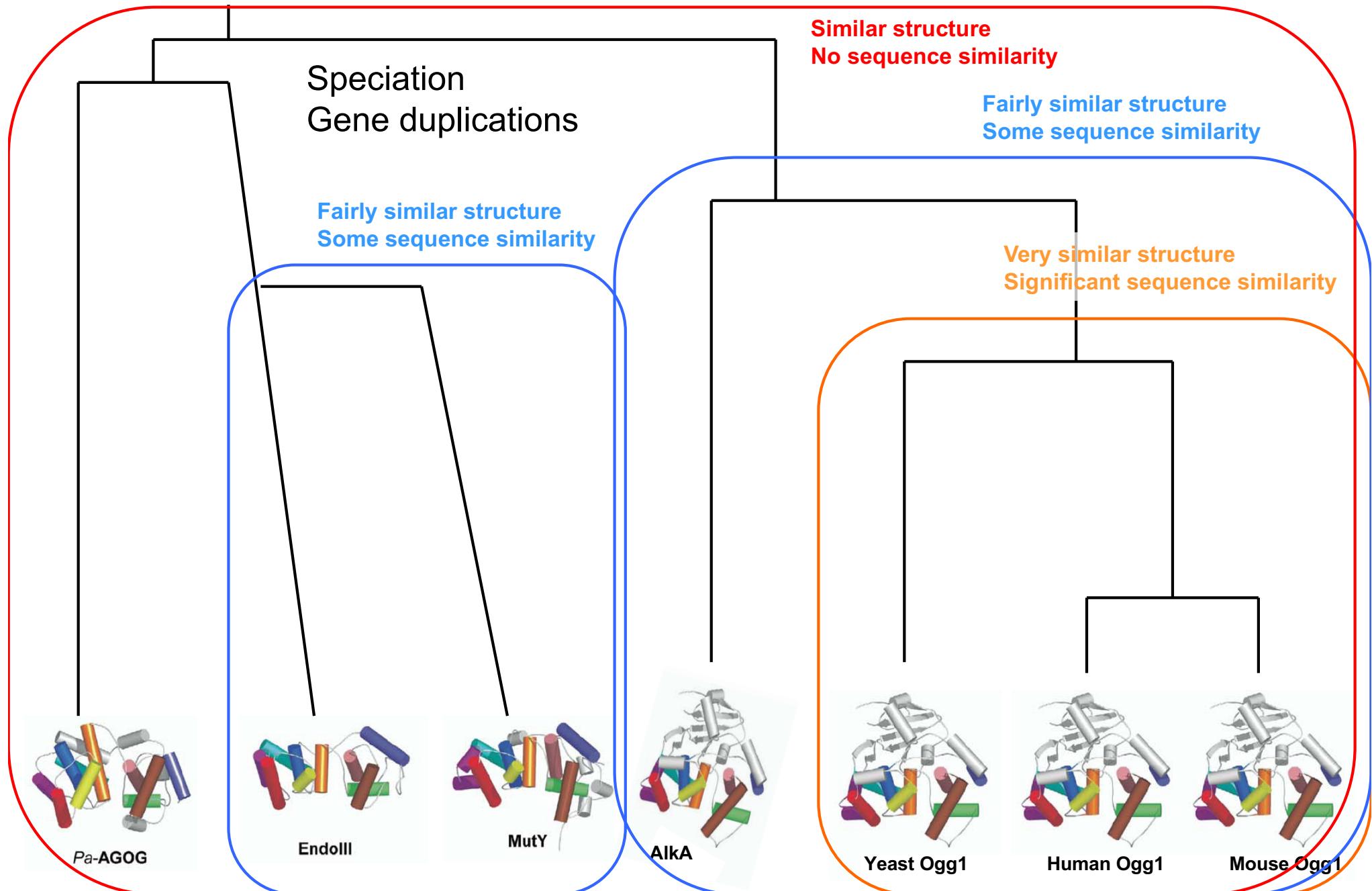




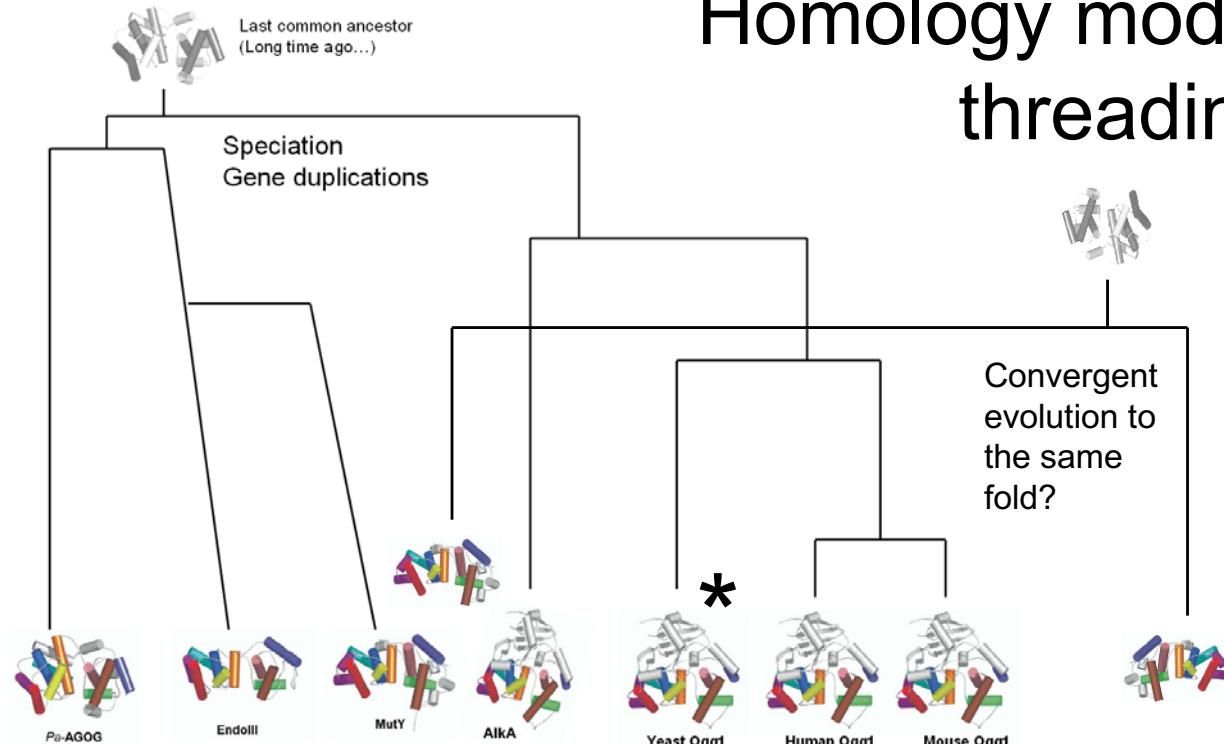
Last common ancestor
(Long time ago...)

Protein structure evolution

Jon K. Lærdahl,
Structural Bioinformatics



Homology modeling and threading



Important goal to have
at least one structure
in all structural
superfamilies!

Structural Genomics
Initiatives

- All proteins (actually domains) in a superfamily have the same overall structure/fold
- If we know (from experiment) the structure of one protein* in a superfamily we may use the information in this structure to model the structure of all other proteins in this superfamily
- Knowledge-based modeling
 - Based on structures in the PDB (*i.e.* they are not *ab initio*)
 - **Homology modeling**
 - When there is significant sequence identity between the protein you want to model (target) and the known structure (template)
 - **Threading**
 - When there is no or little sequence identity between target and template

Structural genomics/The Protein Structure Initiative (PSI)

Jon K. Lærdahl,
Structural Bioinformatics

The screenshot shows the homepage of the PSI Structural Biology Knowledgebase. The top navigation bar includes links for Home, Resource Hubs, Current Focus, Services, About Us, and Contact Us. A search bar allows users to search by sequence, text, PDB ID, or UniProt AC, with a sequence input field containing "MKLTLKNLSMAIMMSTIVMGSSAMAADSNEKIVIAHRGASGYLPEHTLPAKAMAYA". Below the search bar, a sidebar on the left features social media icons (Facebook, Twitter, Email, Print, Plus) and sections for "Current Focus" (Membrane Proteome), "Discoveries", "Membrane Proteins", and "Latest PSI Results" (listing 28 new structures last month, 6528 total structures to date, 5399 total distinct structures, and 564 total community structures). The main content area displays two featured articles: "Membrane Proteome: A Cap on Transport" and "Membrane Proteome: Pumping Out Heavy Metal". To the right, a sidebar titled "Protein Structure Initiative Corner" contains links for the Collaborative Network, Publications (with a "Explore here" button), Latest News, Community Nominations, SBKB Tools, Sequence Analysis, Functional Sleuth, and Visualization Tool.

Traditionally:
solve the structure
of a protein only
after thorough
biological analysis
(years of
research?)

Here: solve
structures of lots
of proteins with
emphasis on
those that are
likely to have a
new fold

Structural genomics/The Protein Structure Initiative (PSI)

Jon K. Lærdahl,
Structural Bioinformatics

Security /
Privacy Notice

MCSG

*Midwest Center
for
Structural Genomics*

• XML Files • Target List • Progress • Statistics • Log in • Site Search: Go

Consortium Project

Investigators Targets

3-D Structures

Related Publications

SG Sites

SG Progress

NIH

MCSG Resources

Job opportunities

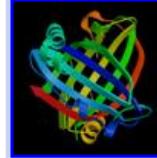
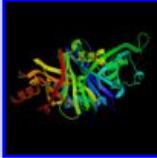
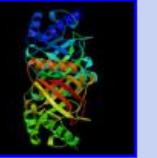
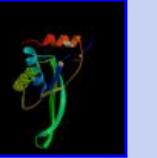
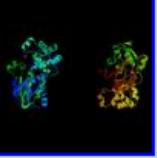
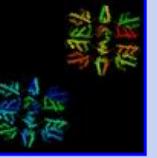
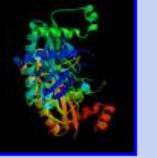
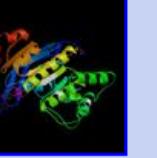
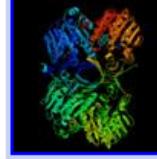
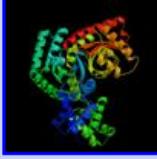
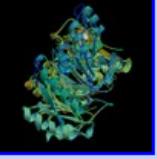
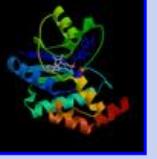
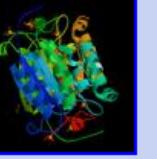
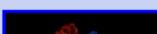
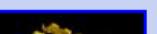
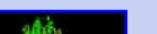
Collaborators

Internals

Technologies

GALLERY OF MCSG STRUCTURES IN PDB

959 targets in PDB (28 new folds)

 APC006 [ref] 1SQE ident: 23.9% annotation	 APC007 1XBW ident: 64.5% annotation	 APC008 2AP3 ident: <20% annotation	 APC009 [ref] 1P99 ident: <20% annotation	 APC010 [ref] 1NG5 New Fold annotation	 APC012 [ref] 1KR4 ident: <20% annotation
 APC014 [ref] 1KYT ident: <20% annotation	 APC037 [ref] 1KXJ ident: 100% annotation	 APC038 [ref] 1M6Y ident: <20% annotation	 APC042 1WPB ident: <20% annotation	 APC043 [ref] 1KUT ident: <20% annotation	 APC046 1J10 ident: 33.5% annotation
 APC047 [ref] 1JQ3 New Fold annotation	 APC048 [ref] 1MKM ident: <20% annotation	 APC049 1T57 ident: <20% annotation	 APC050 [ref] 1EJ2 ident: <20% annotation	 APC063 [ref] 1MKZ ident: 30% annotation	 APC064 [ref] 1M33 ident: 26.2% annotation
 APC009 [ref] 1P99 ident: <20% annotation	 APC010 [ref] 1NG5 New Fold annotation	 APC012 [ref] 1KR4 ident: <20% annotation	 APC042 1WPB ident: <20% annotation	 APC043 [ref] 1KUT ident: <20% annotation	 APC046 1J10 ident: 33.5% annotation

Structural genomics/The Protein Structure Initiative (PSI)

Jon K. Lærdahl,
Structural Bioinformatics

Security /
Privacy Notice

MCSG

Midwest Center
for Structural Genomics

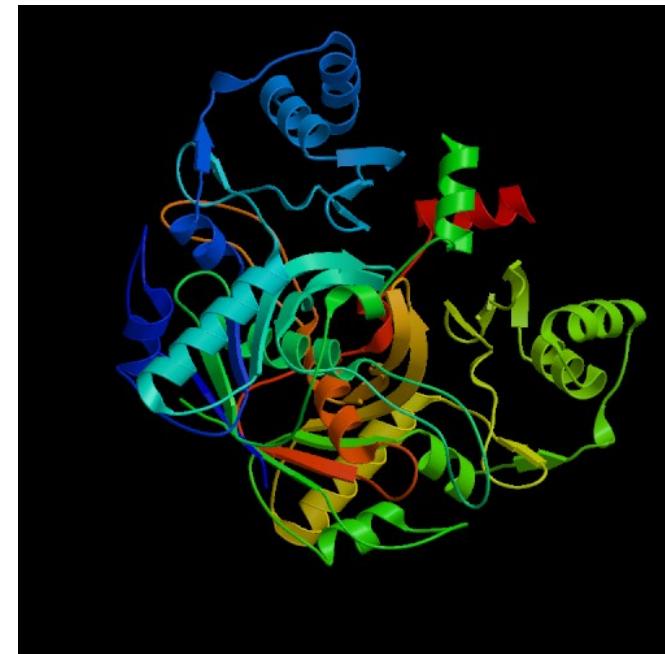
XML Files • Target List • Progress • Statistics • Log in • Site Search: Go

Consortium
Project
Investigators
Targets
3-D Structures
Related Publications
SG Sites
SG Progress
NIH
MCSG Resources
Job opportunities
Collaborators
Internals
Technologies

GALLERY OF MCSG STRUCTURES IN PDB

959 targets in PDB (28 new folds)

APC006 [ref] 1SQE ident: 23.9% annotation	APC007 [ref] 1XBW ident: 64.5% annotation	APC008 [ref] 2AP3 ident: <20% annotation	APC009 [ref] 1P99 ident: <20% annotation	APC010 [ref] 1NG5 New Fold annotation	APC012 [ref] 1KR4 ident: <20% annotation
APC014 [ref] 1KXT ident: <20% annotation	APC037 [ref] 1KXJ ident: 100% annotation	APC038 [ref] 1M6Y ident: <20% annotation	APC042 [ref] 1WPB ident: <20% annotation	APC043 [ref] 1KUT ident: <20% annotation	APC046 [ref] 1JQ0 ident: 33.5% annotation
APC047 [ref] 1IQ3 New Fold annotation	APC048 [ref] 1MKM ident: <20% annotation	APC049 [ref] 1T57 ident: <20% annotation	APC050 [ref] 1EP2 ident: <20% annotation	APC063 [ref] 1MKZ ident: 30% annotation	APC064 [ref] 1M33 ident: 26.2% annotation



Archaeoglobus fulgidus DSM 4304 protein AAB89001.1 has a new fold determined by the MCSG (2PHN/2G9I)

10 yrs ago: “Only” 3D structures for proteins that had been studied a lot

Now: many 3D structures for proteins with unknown function!

PSI concluded in 2015 (7000 structures)

Homology modeling

- Based on: during evolution, structure is more stable and conserved than the associated sequence
- Similar sequences give nearly identical structure
- Distantly related sequences fold into similar structures
- 20-30% identical residues to a known (experimental) structure

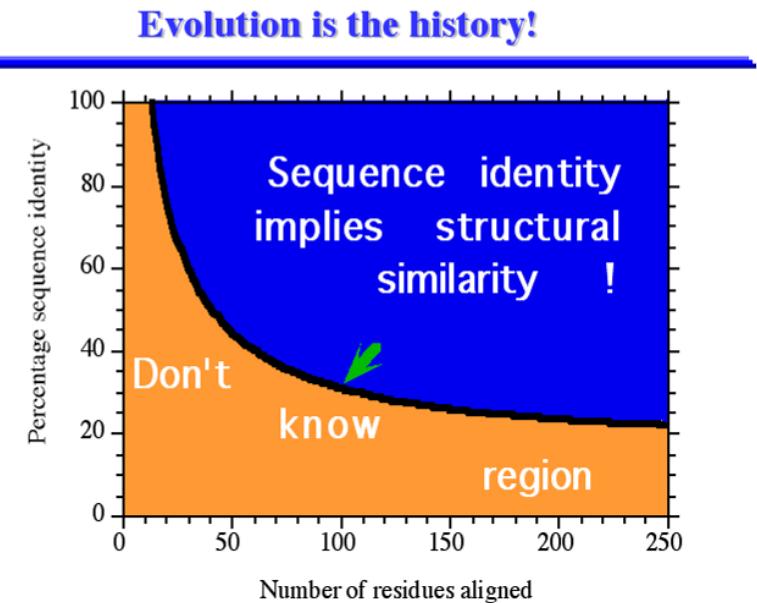
→ Might be able to predict the 3D structure with some confidence

Known (experimental)
structure of protein 1
(*template*)

&
Sequence alignment with
protein 2 (*target*)



Model of
protein 2



B. Rost, *Prot. Engin.* **12**, 85 (1999)

- 30% sequence identity necessary (in textbooks)
- My experience: Might get reasonable results also at 20% or even below
- Depends on
 - Many indels or not?
 - Length of alignment
 - Automatic or manual modeling?

Homology modeling

Start with a protein sequence (target)

1. Template selection:
 - Find template in PDB and align sequences
2. Correct alignments
 - Use the best MSA programs
 - Correct placement of insertions and deletions
3. Backbone model building
4. Model loops and side-chains
 - Rotamer libraries
 - Loop modeling using database or *ab initio* method
5. Refine and optimize model
6. Validate and check model quality!

Homology modeling

Start with a protein sequence (target)

1. Template selection:

- Find template sequences

2. Correct alignment

- Use the best alignment
- Correct placement and deletions

3. Backbone modeling

4. Model loops and

- Rotamer libraries
- Loop modeling using database or *ab initio* method

5. Refine and optimize model

6. Validate and check model quality!

I want to model this!

```
>gi|84618885|emb|CAJ31885.1| methylpurine-DNA  
glycosylase [Bacillus cereus]  
MHPFVKALQEHFIAHKNPEKAEPMARYMKNHFLFIGIQT  
PERRQLLKDVQIHTLPDPKDFRIIVRELWDLPEREFQA  
AALDMMQKYKKYINETHIPFLEELIVTKSWWDTVDSIVP  
TFLGNIFLQHPELISAYIPKWIASDNIWLQRAAILFQLK  
YKQKMDEELLFWVIGQLHSSKEFFIQKAIGWVLREYAKT  
KPDVVWEYVQNNELAPLSRREAIKHIKENYGINNEKIGE  
TLS
```

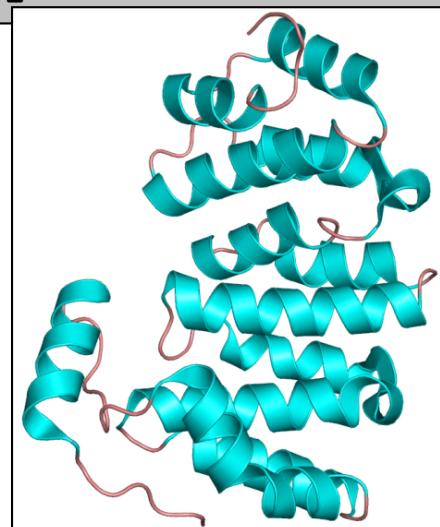
Homology modeling

Start with a protein sequence (target)

1. Template selection:
 - Find template in PDB and align target sequence to template sequences
2. Correct alignments
 - Use the best MSA program
 - Correct placement of insertions and deletions
3. Backbone model building
4. Model loops and side-chains
 - Rotamer libraries
 - Loop modeling using database or *ab initio* method
5. Refine and optimize model
6. Validate and check model quality!

Do sequence search in all "PDB sequences"

- Useful templates have 30% or higher sequence identity to target (but sometimes even lower)
- Several templates?
 - Resolution?
 - Highest sequence identity?
 - Cofactors?
 - Use the structure that best fits your task



Homology modeling

Start with a protein sequence (target)

1. Template selection:

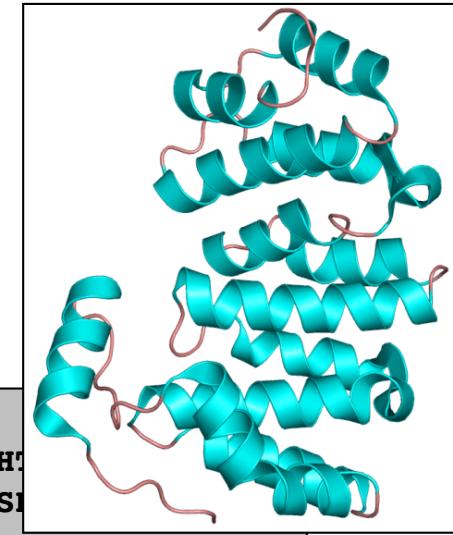
- Find template in PDB and align sequences

2. Correct alignments

Sequence alignment

Bc_AlkD	MHPFVKALQEHFIAHKNPEKAEPMARYMKNHFLFIGIQTPERRQLLKDVIQIH	
EF3068	-----MDTLQFQKNPETAAKMSAYMKHQFVFAGIPAPERQALSQOLLKESI	
	: : :*****.* *: ***: :*: * ** :***: * *: : :	..
Bc_AlkD	FRIIVRELWDLP ERE QAAALDMMQKYKKYINETHIPFLEELIVT K S W D T V D SIVPT F L	120
EF3068	LCQEIEAYYQKT ERE YQVAIDLALQNVQRFSLEEVVAFKAYVPQ K A W D S VDAWRKF F G	122
	: .. : : .****: .*: : : .. : : : : *:****: *: *	
Bc_AlkD	GNIFLQHPELISAYIPKWIASDNIWLQRAAILFQLKY K QMDEELLFWVIGQLHSSKEFF	180
EF3068	SWVALHLTELPT-IFALFYGAENFWNRRVALNLQLMLKEKTNQDLLKKAIYDRTEEFF	171
	. : *: **: .. : .. : :*: : *: : *: : :** . * : : : ***	
Bc_AlkD	IQKAIGWVLREYAKTPDVVWEYVQNNELAPLSRREA I KHIKENYGINNEKIGETLS	237
EF3068	IQKAIGWSLRQYSKTNPQWVEELMKELVLSPLAQREGSKYLAKASE-----	217
	***** *: :*: :*: * * : : *: :*: :*. * : :	

Alignment of the sequences of *B. cereus* AlkD (target) and *E. faecalis* hypothetical protein EF3068 (template from MCG).



3.

4.

5. Refine and optimize model

6. Validate and check model quality!

Homology modeling

Start with a protein sequence (target)

1. Template selection:

- Find template in PDB and a large number of homologous sequences

2. Correct alignments

- Use the best MSA program

Correct placement of insertions

Check indels/deletions

3. Backbone tracing

4. Model building

- Rely on loops

5. Refinement

6. Validation

Obtaining the correct alignment is the most important step!! in homology modeling

FIRST: Align target, template and a large number (50-100?) of homologs with Praline, T-Coffee, Muscle or a different good MSA program

Use target/template alignment from this MSA

SECOND: Look at the template structure and move all indels

- to loops
- out of helices/sheets

Sequence alignment		
Bc_AlkD	MHPFVKALQEHFIAH KNPEKAEPMARY	
EF3068	-----MDTLQFQ KNPETAAKMSAY	
	: : :*****.* *: *	
Bc_AlkD	FRIIVRELWDLP EREFQAAALDMMQKYKKYINETHIPFLEELIVTKSWWDTVDSIVPTFL	120
EF3068	LCQRIEAYYQKT EREYQVVAIDLALQNVQRFSLEEVVAFKAYVPQKAWWDSVDAWRKFFG	122
	: : .. :: .*****.* .*: : : : .. : : : : *:*****: *	
Bc_AlkD	GNIFLQHPELISAYIPKWIASNWLQRAAILFQLKYKQKMDEELLFWVIGQLHSSKEFF	180
EF3068	SWVALHLTELPT-IFALFYGAENFWNRRVALNLQLMLKEKTNQDLLKKAIYDRTTEEFF	171
	. : *: **: .. : .:*****: :*. : ** *: : :** .*. : :*****	
Bc_AlkD	IQKAIGWVLREVAKTPDVVWEYVNNEAPLSRREAIIKHIKENYGINNEKIGETLS	237
EF3068	IQKAIGWSLRQYSKTNPQWVEELMKELVLSPLAQREGSKYLAKASE-----	217
	***** *:*****: * * : : *:*****. *: :	

Alignment of the sequences of *B. cereus* AlkD (target) and *E. faecalis* hypothetical protein EF3068 (template from MCGS).

Homology modeling

Obtaining the correct alignment
ant step!! in

Start

1.

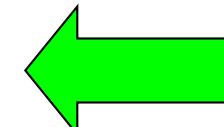
Where is the correct position of the gap?

XP_968170_Tribolium_castaneum	10	20	30
KTLGTGSFGRVMIVQHKPT-KEYYAMKILDKOK			
Q4JIV3_Lymnaea_stagnalis	KTLGTGSFGRVMLVQHKGENDKAYYAMKILDKOK		



2.

XP_968170_Tribolium_castaneum	10	20	30
KTLGTGSFGRVMIVQHK-PTKEYYAMKILDKOK			
Q4JIV3_Lymnaea_stagnalis	KTLGTGSFGRVMLVQHKGENDKAYYAMKILDKOK		



The MSA gives the answer!!

3.

Homolog2_Petromyzon_marinus	10	20	30
KTLGTGSFGRVMLVKHK-ATDRYFAMKILDKOK			
ENSCJAP0000040924_Callithrix_jacchus	KTLGIGSFGRVVLVSHR-ESGSHYAMKILNKEK		
P22612_Homo_sapiens	RTLGGMGSFGRVMLVRHQ-ETGGHYAMKILNKOK		
XP_968170_Tribolium_castaneum	KTLGTGSFGRVMIVQHK-PTKEYYAMKILDKOK		
Q4JIV3_Lymnaea_stagnalis	KTLGTGSFGRVMLVQHKGENDKAYYAMKILDKOK		
ENSTRUP0000015108_Takifugu_rubripes	KTLGTGSFGRVMLVKHK-ETNQFYAMKILDKOK		

4.

5.

BC_AlkD	IQKAIGWVIREVAKTPDVVWEYVQNNELAPLSRREAIKHIKENYGINNEKIGETLS	237
EF3068	IQKAIGWSLRQYSKTNPQWVEELMKELVLSPLAQREGSKYLAKASE-----	217
		***** * : * : * : * : * * : : * : * : * : * : * : :

Alignment of the sequences of *B. cereus* AlkD (target) and *E. faecalis* hypothetical protein EF3068 (template from MCSG).

Valid

t, template
(50-100?) of
ine, T-
a different

e alignment

e template
all indels

heets

DSIVPTFL	120
DAWRKFFG	122
*: *	
LHSSKEFF	180
DRTTEEFF	171
: : : : ***	

BC_AlkD	IQKAIGWVIREVAKTPDVVWEYVQNNELAPLSRREAIKHIKENYGINNEKIGETLS	237
EF3068	IQKAIGWSLRQYSKTNPQWVEELMKELVLSPLAQREGSKYLAKASE-----	217
		***** * : * : * : * : * * : : * : * : * : * : :

Homology modeling

Start with

1. Temp

- Find sequence

2. Corre

- Use
- Correct and detections

3. Back

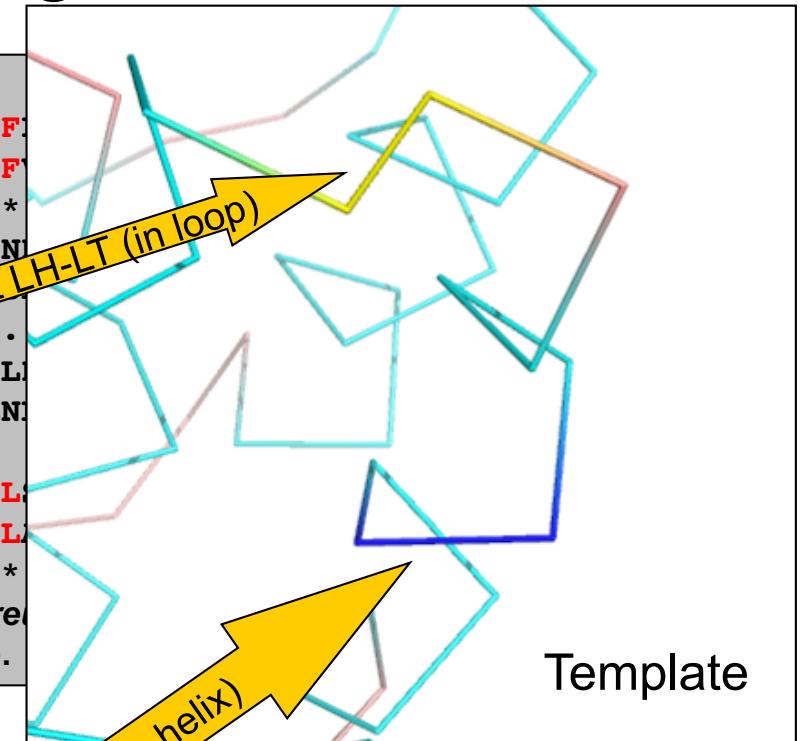
4. Mod

- Replace
- Loop

5. Refin

6. Valid

```
Sequence alignment
Bc_AlkD MHPFKALQEHFIAHKNPEKAEPMARYMKNHF
EF3068 -----MDTLQFQKNPETAAKMSAYMKHQF
: : :****.* *: ***::*
Bc_AlkD FRIIVRELWDLPEREFQAAALDMMQKYKKYIN
EF3068 LCQEIEAYIYQKTEREYQVVAIDLALONI
: .. :: .***:* .**.*: : : ..
Bc_AlkD GNIFLQHPELISAYIPKWIASDNIWLQRAAIL
EF3068 SWVALH-LTELPTIFALFYGAENFWNRRVALN
. : *: .. : . : .***: *.*: *
Bc_AlkD IQKAIGWVLPREVAKTKEPDVVWEYVQNNELAPL
EF3068 IQKAIGWSLRQYSKTNPQWVEELMKELVLSPL
***** * :***: * : * : : * : *: *
CORRECTED Alignment of the sequences of B. cereus AlkD (target) and E. faecalis hypothetical protein EF3068 (template from MCGS).
```



```
Sequence alignment
Bc_AlkD MHPFKALQEHFIAHKNPEKAEPMARYMKNHF 60
EF3068 -----MDTLQFQKNPETAAKMSAYMKHQF 52
: : :****.* *: ***::* ** :***: * *::: .. :
Bc_AlkD FRIIVRELWDLPEREFQAAALDMMQKYKKYINETHIPFLEELIVTKSWWDTVDSIVPTFL 120
EF3068 LCQEIEAYIYQKTEREYQVVAIDLALONIQVRFSLEEVVAFKAYVPQKAWWDSDVAWRKFFG 122
: .. :: .***:* .**.*: : : .. : : *:***: *: *
Bc_AlkD GNIFLQHPELISAYIPKWIASDNIWLQRAAILFQLKYKQMDEELLFWVIGQLHSSKEFF 180
EF3068 SWVALHLTELPT-IFALFYGAENFWNRRVALNLQLMLKEKTNQDLLKKAIYDRTTEEFF 171
. : *: **: .. : . : .***: *.*: :***: * : : : *
Bc_AlkD IQKAIGWVLPREVAKTKEPDVVWEYVQNNELAPLSRREAIKHIKENYGINNEKIGETLS 237
EF3068 IQKAIGWSLRQYSKTNPQWVEELMKELVLSPLAQREGSKYLAKASE----- 217
***** * :***: * : * : : * : : *:***: *. *: :
```

Alignment of the sequences of *B. cereus* AlkD (target) and *E. faecalis* hypothetical protein EF3068 (template from MCGS).

Homology modeling

Start with a protein sequence (target)

1. Template selection:

- Find template in PDB and align sequences

2. Correct alignments

- Use the best MSA programs
- Correct placement of insertions and deletions

3. Backbone model building

4. Model loops and side-chains

- Rotamer libraries
- Loop modeling using database or *ab initio* method

5. Refine and optimize model

6. Validate and check model quality!

The most
important step in
homology
modeling!

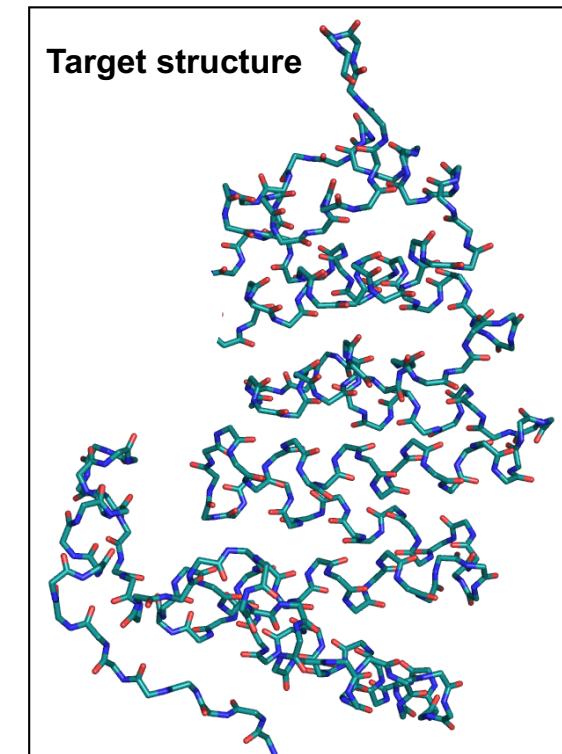
Homology modeling

Start with a protein sequence (target)

1. Template selection:
 - Find template in PDB and sequences
2. Correct alignments
 - Use the best MSA programs
 - Correct placement of insertions and deletions
3. Backbone model building
4. Model loops and side-chains
 - Rotamer libraries
 - Loop modeling using database or *ab initio* method
5. Refine and optimize model
6. Validate and check model quality!

For all aligned residues in template and target:

- Take coordinates for template backbone atoms and use for target
 - If residues are identical: Use all atom coordinates from template in target
 - Indels: Nothing to copy



Homology modeling

Start

1.

Short loops (3-5 residues) :
Reliable results with both
methods

2.

Long loops (more than 10-15
residues) : Highly unlikely
that you get a correct
result!!

3.

- Use the best MSA programs
- Correct placement of insertions
and deletions

4.

Model loops and side-chains

5.

– Rotamer libraries

6.

– Loop modeling using database
or *ab initio* method

7.

Refine and optimize model

8.

Validate and check model quality!

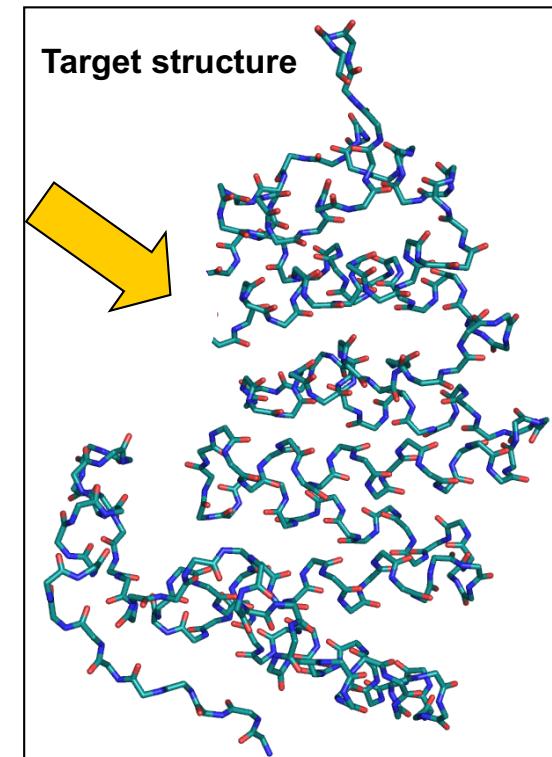
(tar

nd

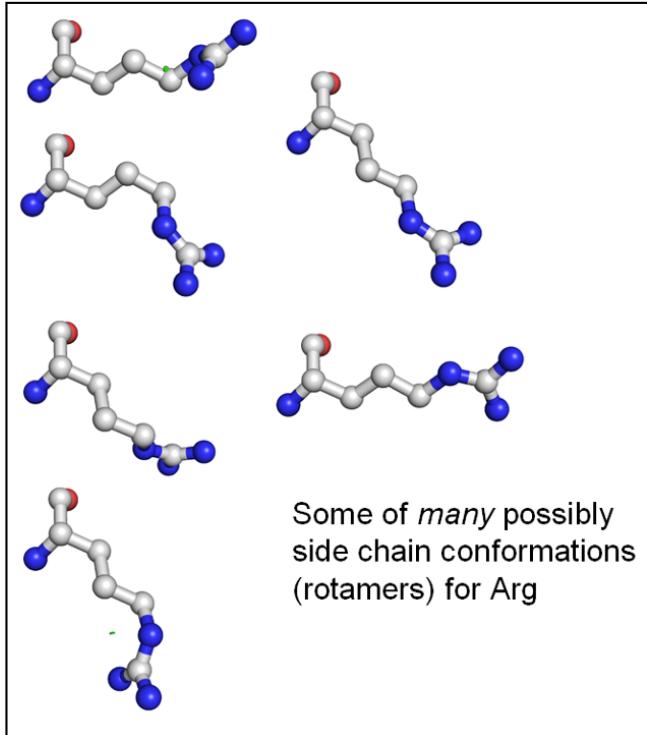
Ab initio: Generates random
loops and chooses the one with

- Lowest energy scores
- Ok Ramachandran plot
- No clashes

Database method: Try loops
taken from a "loop-library"
extracted from the PDB



Homology modeling



Some of *many* possible side chain conformations (rotamers) for Arg

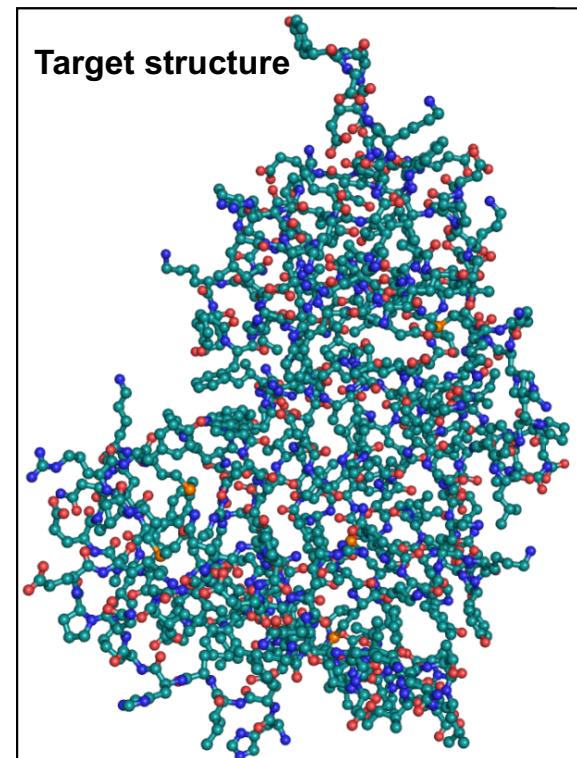
Sequence (target) +
on:
e in PDB and
nts
MSA programs
ment of insertions

Get side chain conformations from rotamer libraries generated from known structures

Use those that give

- Lowest energy score
- No clashes with backbone/other side chains

3. Backbone model building
4. Model loops and side-chains
 - Rotamer libraries
 - Loop modeling using database or *ab initio* method
5. Refine and optimize model
6. Validate and check model quality!



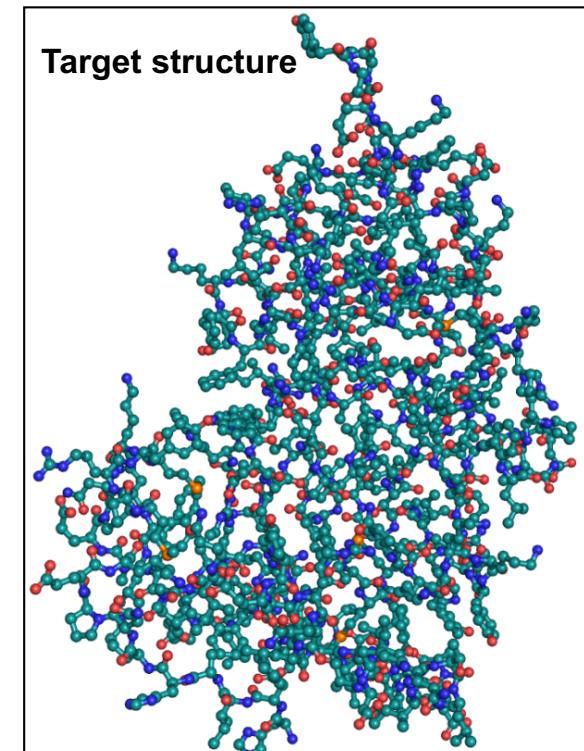
Homology modeling

Start with a protein sequence (target)

1. Template selection:
 - Find template in PDB and search against sequences
2. Correct alignments
 - Use the best MSA programs
 - Correct placement of insertions and deletions
3. Backbone model building
4. Model loops and side-chains
 - Rotamer libraries
 - Loop modeling using database or *ab initio* method
5. Refine and optimize model
6. Validate and check model quality!

Do a few hundred iterations of energy minimization?

- Will hopefully remove clashes and very unfavorable conformations
- Too many iterations will most likely destroy structure
- Not always necessary (depends on the program)



Homology modeling

Check if model makes sense?

- Ramachandran plot ok?

- No clashes?

- No funny bond

lengths/angles/conformations?

- Use programs such as:

- Procheck

- WHAT IF

- ANOLEA

- Verify3D

- These can only check if the chemical/physical properties are ok

- The model might still be 100% meaningless biologically and completely wrong!

- Rotamer library
- Loop modeling or *ab initio* met

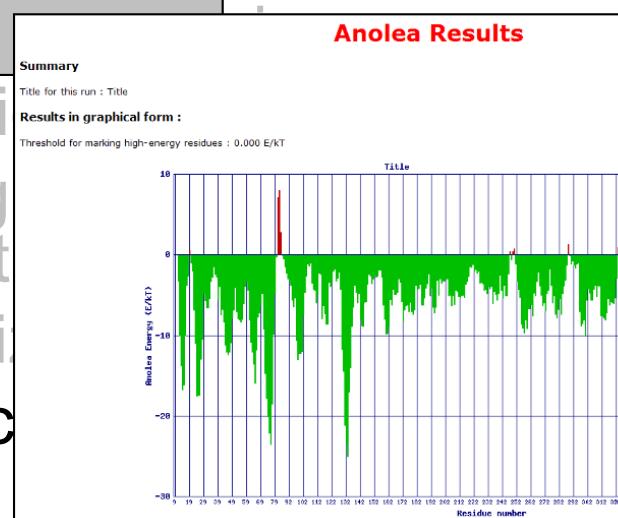
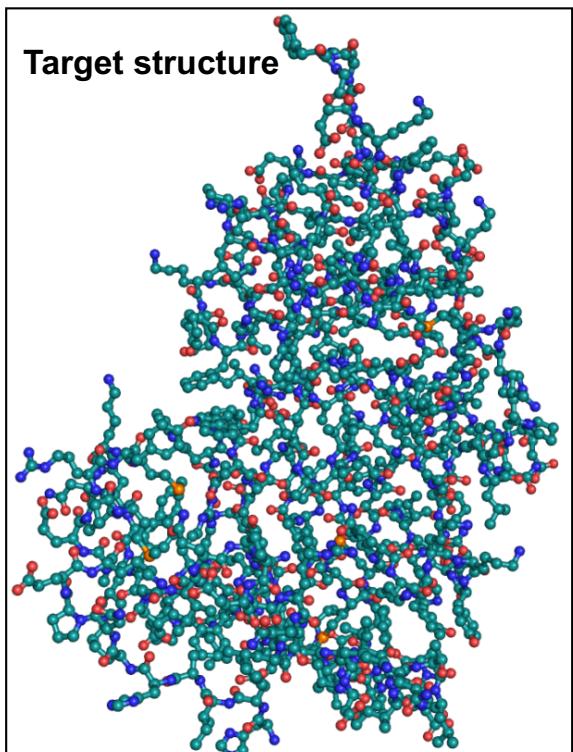
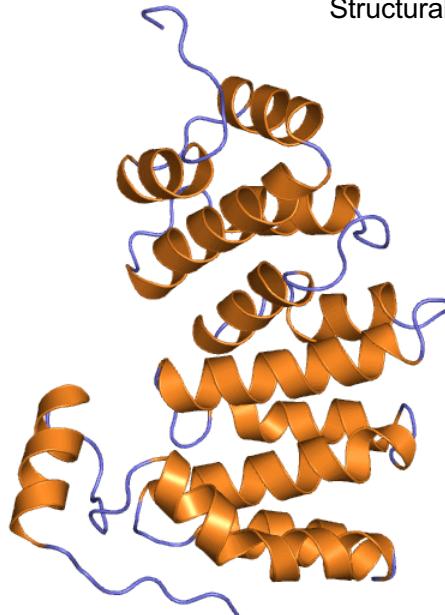
5. Refine and optimize

6. Validate and check

target

and align

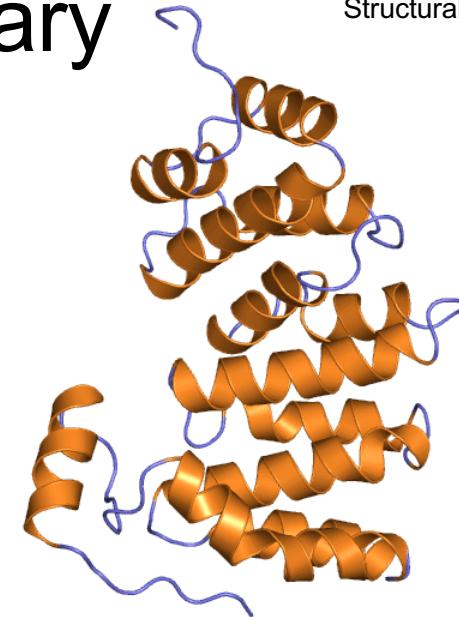
grams
n insertions



Homology modeling summary

1. Template selection:
 - Find template in PDB and align sequences
2. Correct alignments
 - IMPORTANT!
3. Backbone model building
4. Model loops and side-chains
5. Refine and optimize model(?)
6. Validate and check model quality!

Automatic models usually less accurate than manually generated models (if the modeler knows what she is doing...)



Tools:

- Modeller
- Swiss-Model
- 3D-JIGSAW

Homology model databases:

- Modbase (automatic modeling with Modeller)
- SWISS-MODEL Repository (automatic modeling with Swiss-Model)

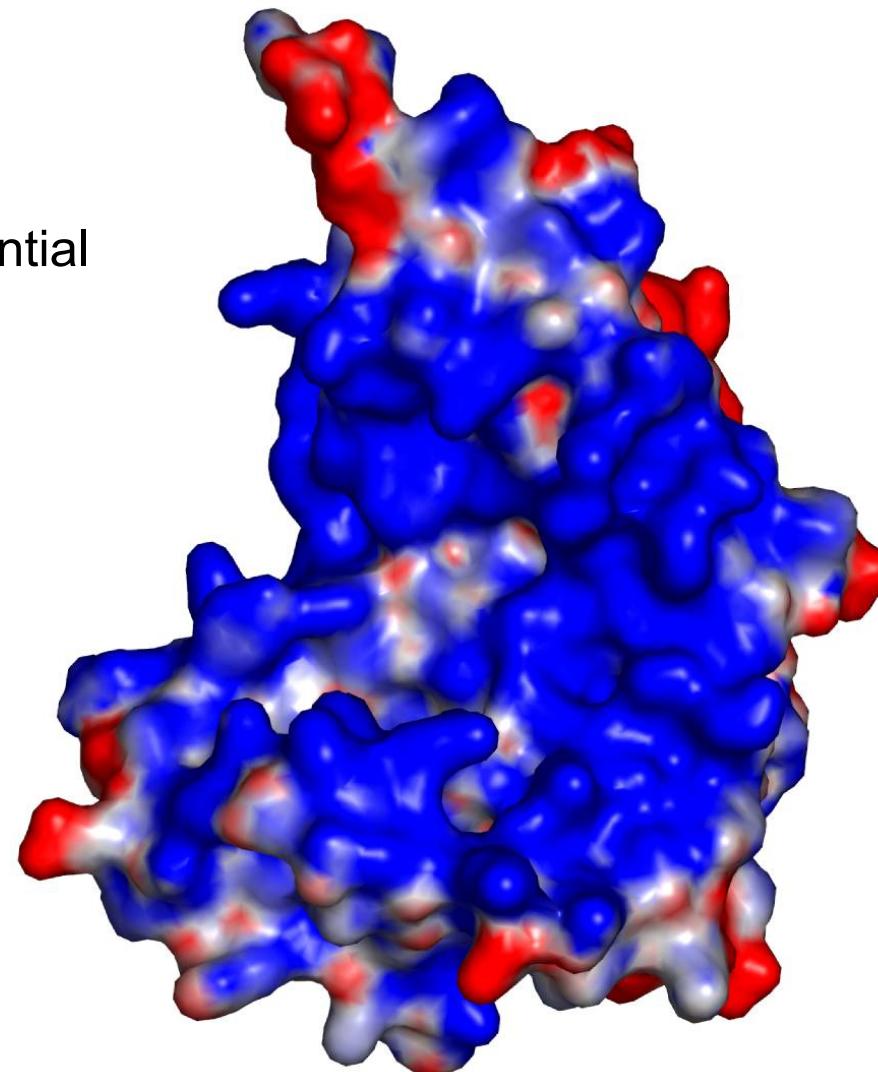
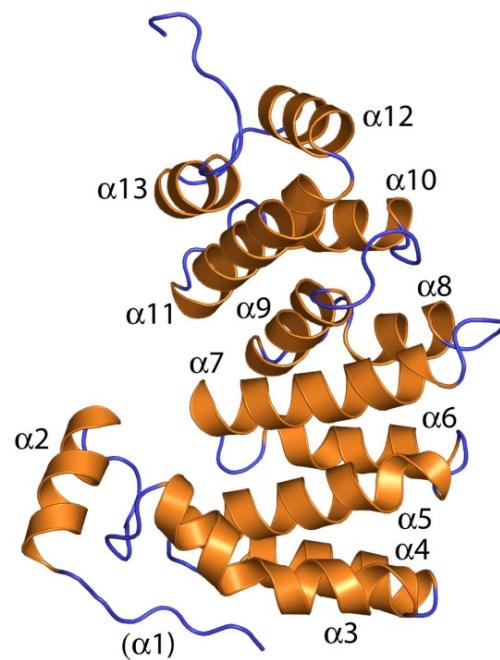
Structural bioinformatics

Jon K. Lærdahl,
Structural Bioinformatics

When the structure (experimental or model) is available, there are many more possibilities to obtain understanding

Some examples:

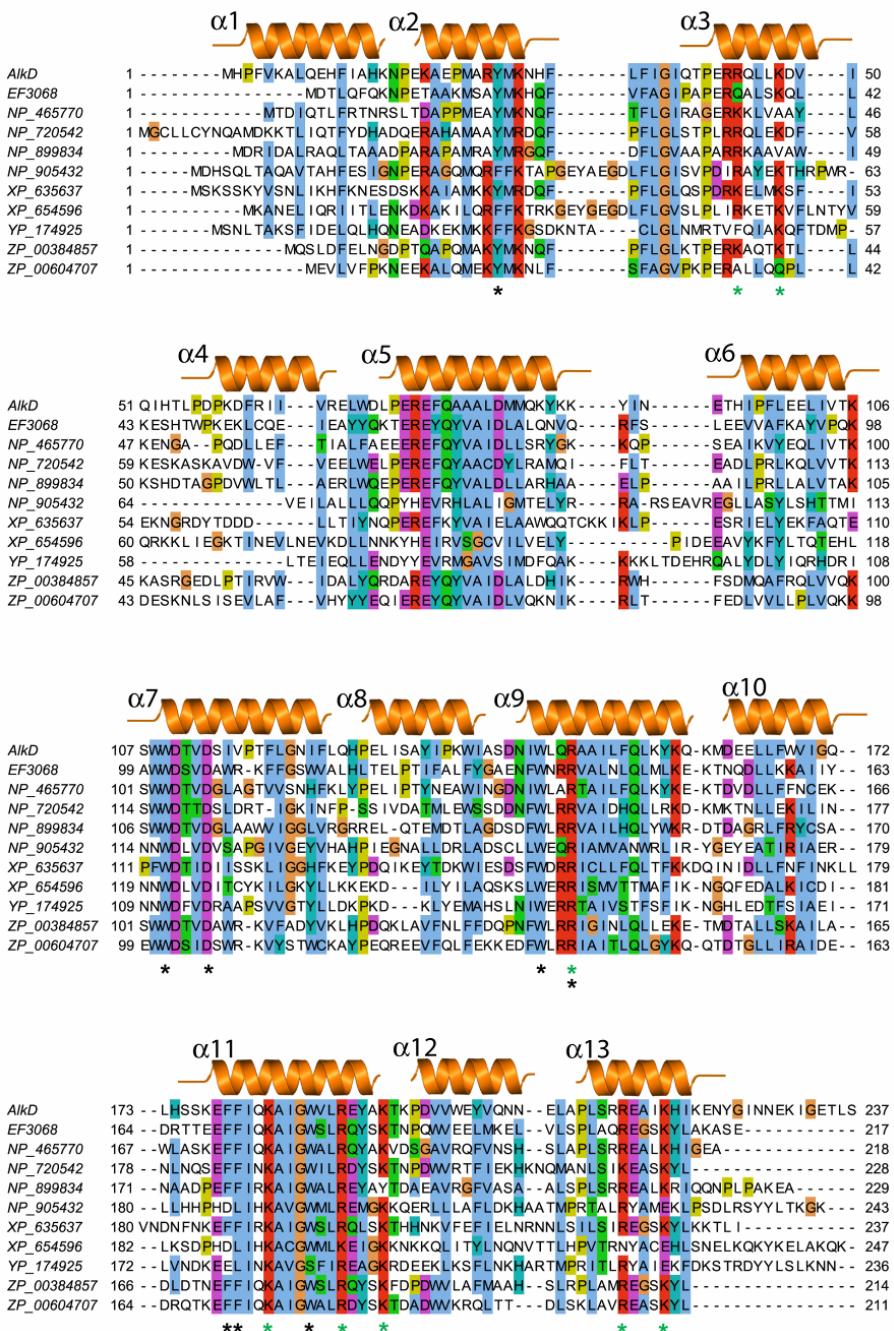
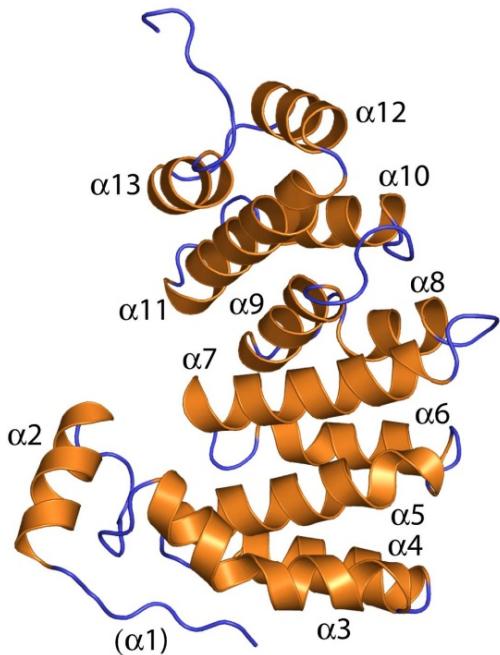
B. cereus AlkD electrostatic potential



Structural bioinformatics

Jon K. Lærdahl,
Structural Bioinformatics

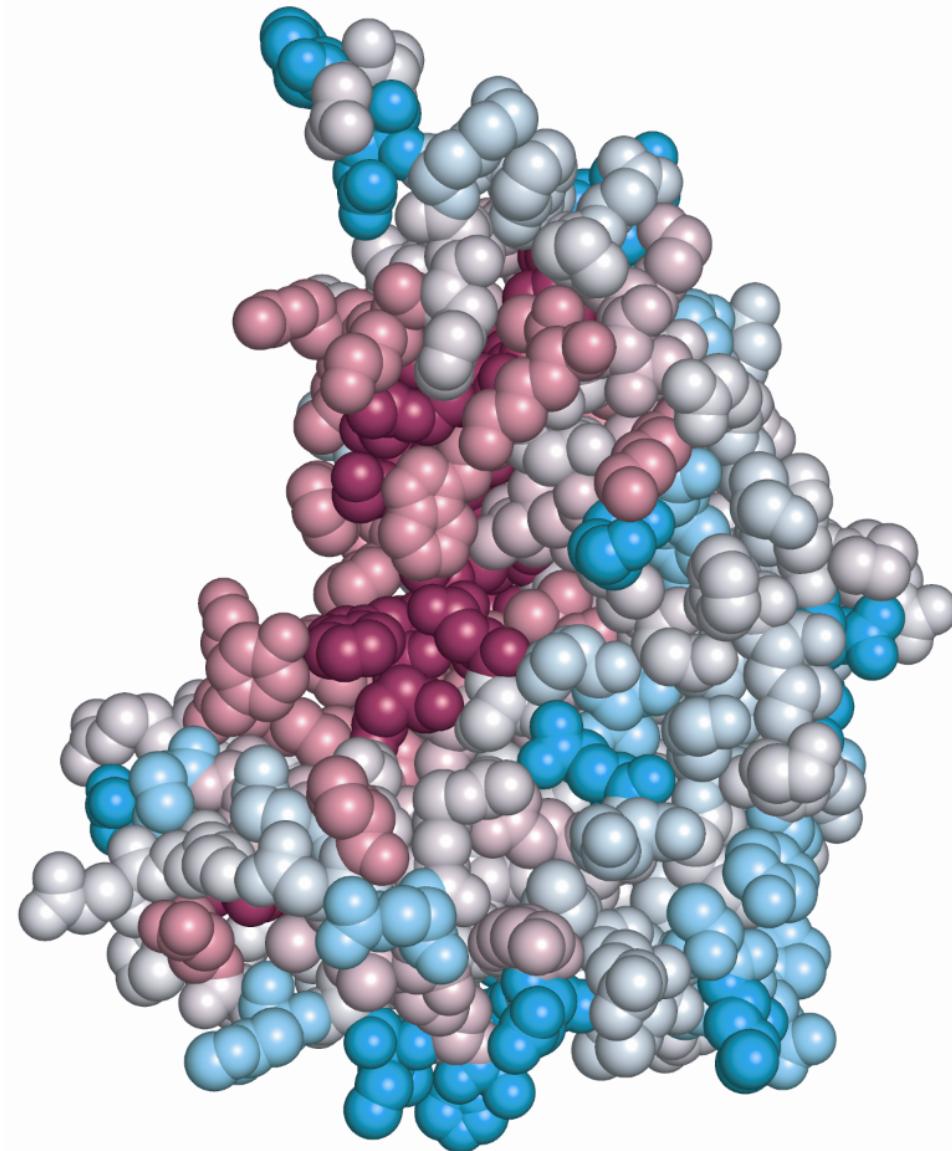
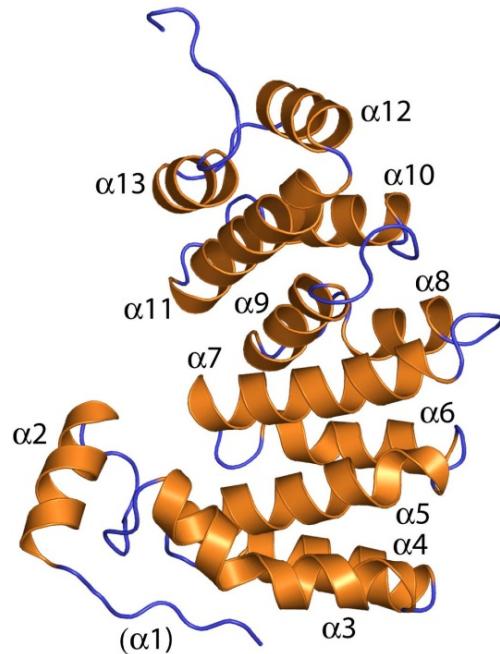
B. cereus AlkD sequence conservation from ConSurf:



Structural bioinformatics

Jon K. Lærdahl,
Structural Bioinformatics

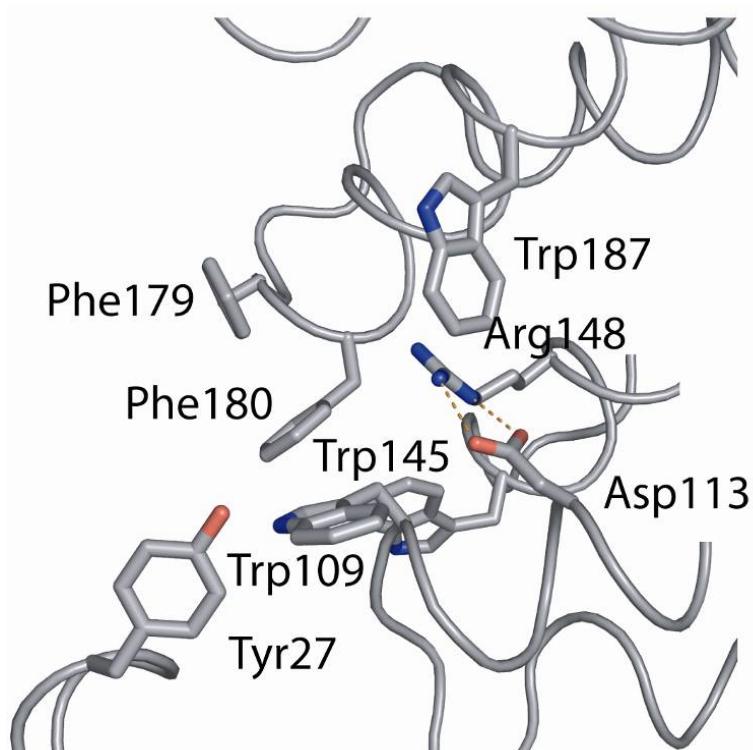
B. cereus sequence conservation from ConSurf:



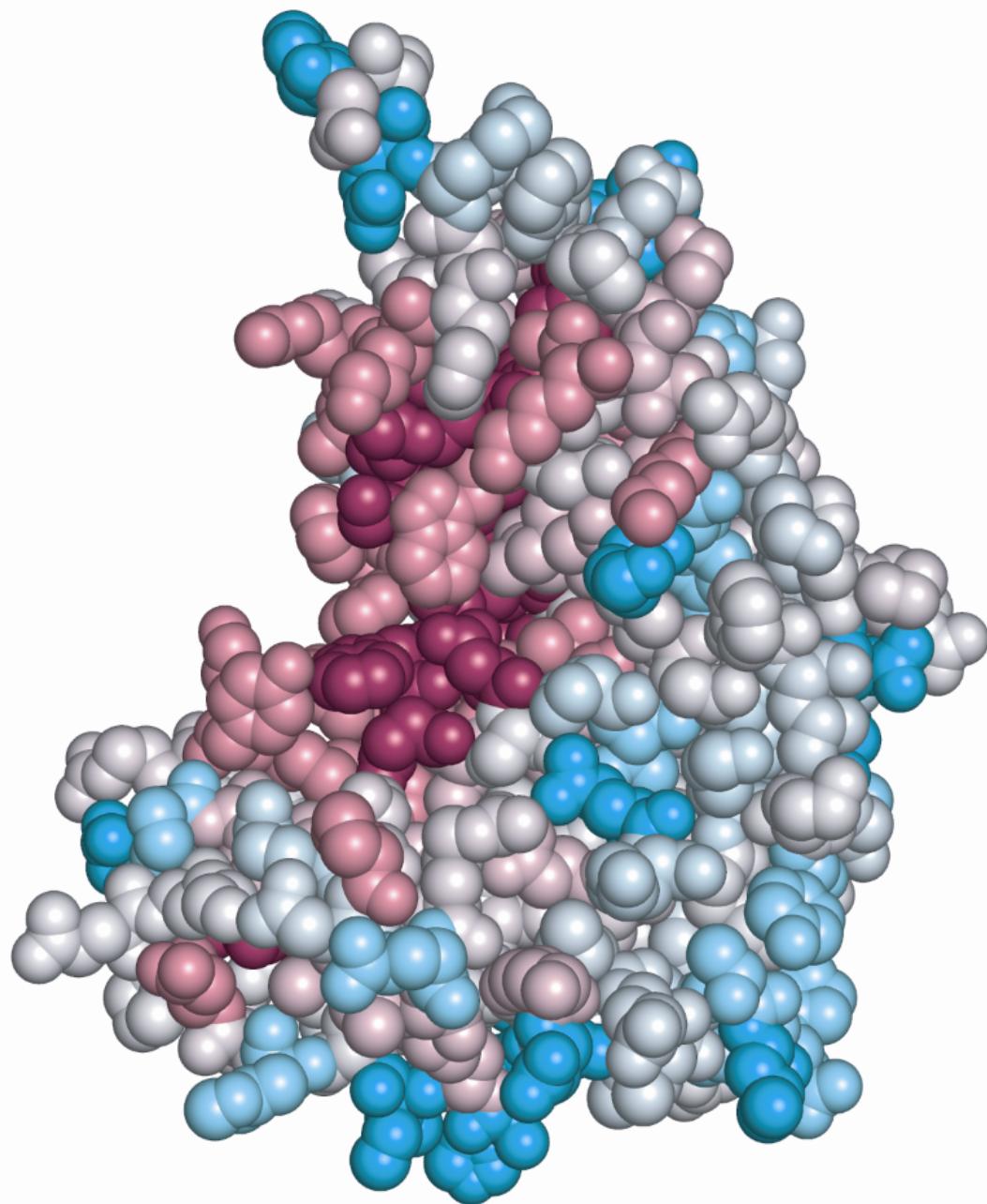
Structural bioinformatics

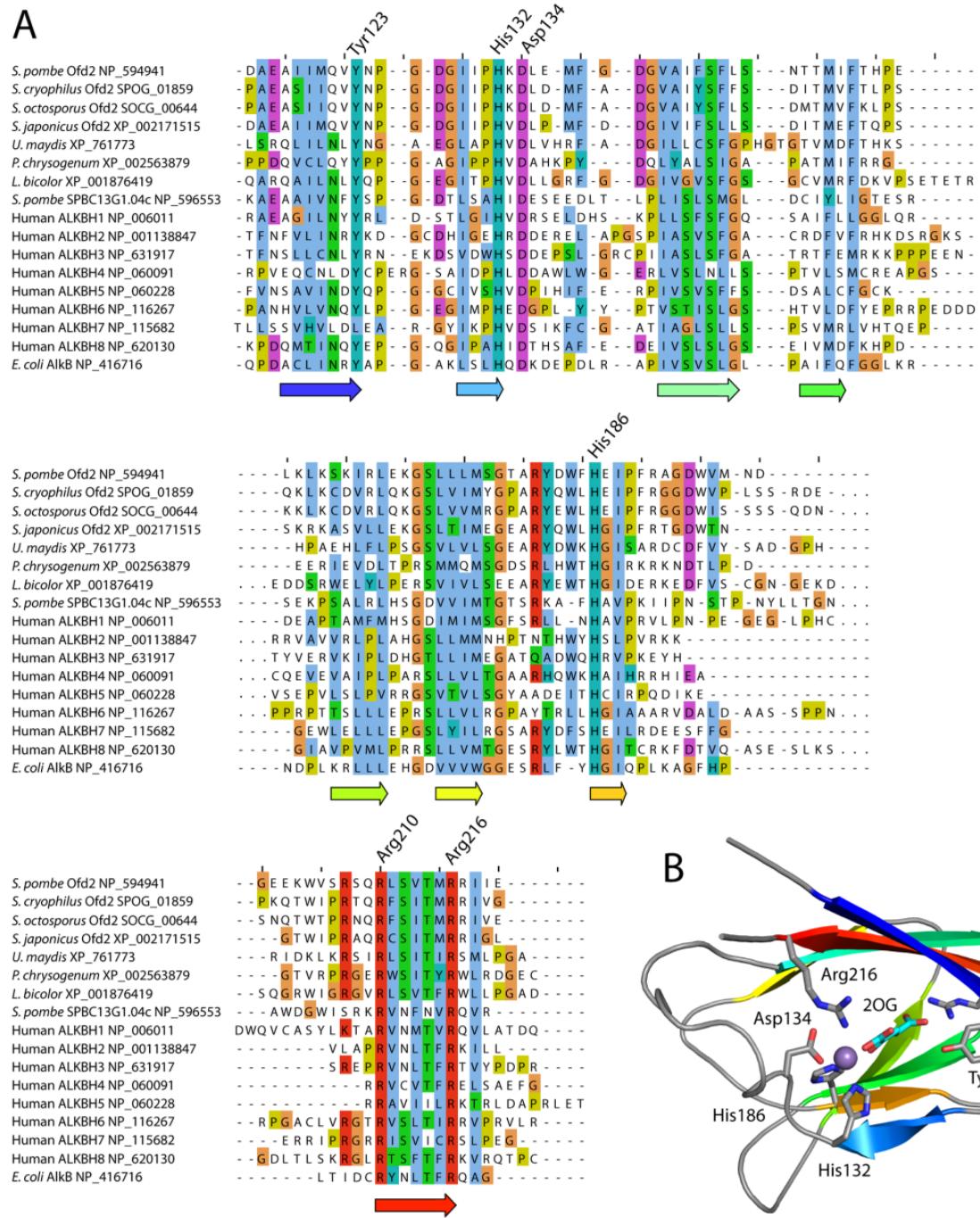
Jon K. Lærdahl,
Structural Bioinformatics

B. cereus sequence conservation
from ConSurf:



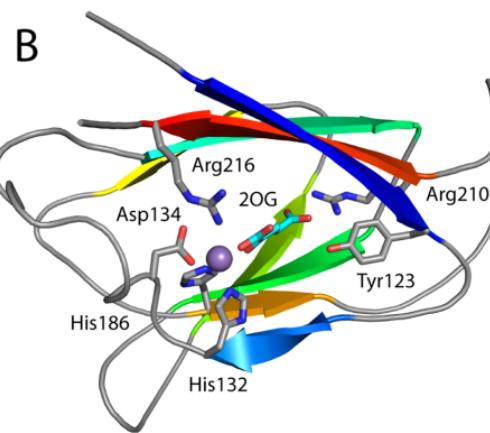
Dalhus et al., Nucleic Acids Res. 35, 2451 (2007).



A

Use MANY homologs to align two (or a few) homologs!

Korvald *et al.* PLOS One 6, e25188 (2011)

**B**

Use MANY homologs to align some homologs!

Bioinformatics analysis

The sequences of *S. pombe* Ofd2 and 63 homologous sequences from fungal and metazoan species were obtained from the NCBI protein sequence databases [43]. The sequences for two additional *Schizosaccharomyces* species, *S. cryophilus* and *S. octosporus*, were retrieved from the Broad Institute Schizosaccharomyces group database [44] (http://www.broadinstitute.org/annotation/genome/schizosaccharomyces_group). The sequences were aligned with

Korvald et al. PLOS
One 6, e25188 (2011)

Ofd2 Dioxygenase Interacts with Histones

Expresso [45] and manipulated in Jalview [46]. The main bulk of sequences were subsequently removed in order to give a reliable alignment of Ofd2, AlkB and human and fungal homologs. Structural disorder predictions were performed with the VSL1 algorithm [47] and DISOPRED2 [48]. The structural model of the Ofd2 core domain was derived from an *E. coli* AlkB template from Yu and Hunt [32] (Protein Databank identifier 3I3Q) and the illustration was generated with PyMOL [49].