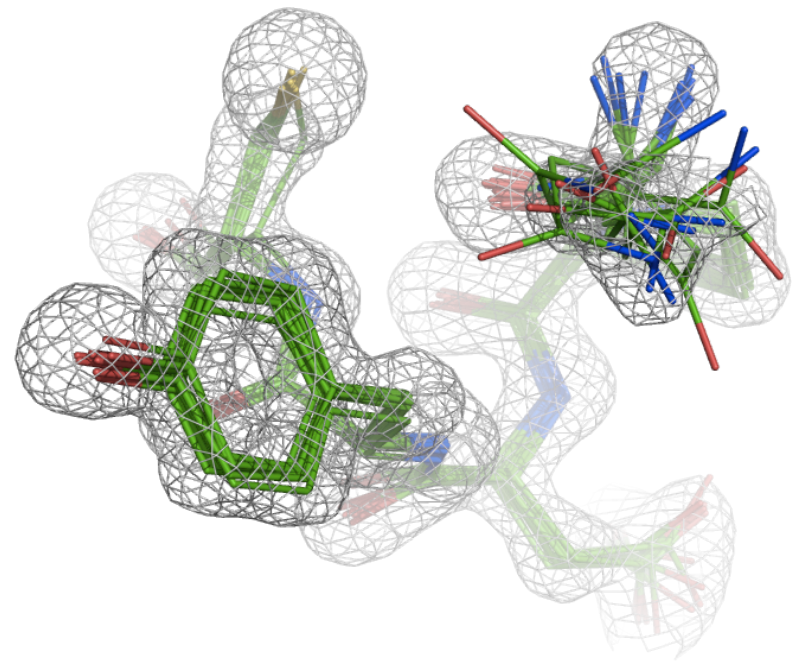


# *Ensemble refinement of protein crystal structures in PHENIX*



Tom Burnley | Piet Gros



Universiteit Utrecht



Incomplete modelling of disorder contributes to R factor gap

Only ~5% of residues in the PDB are modelled with more than one conformation (x-ray structures)

Multiple discrete models restricted due to increase in number of model parameters

Incomplete modelling of disorder contributes to R factor gap

Only ~5% of residues in the PDB are modelled with more than one conformation (x-ray structures)

Multiple discrete models restricted due to increase in number of model parameters

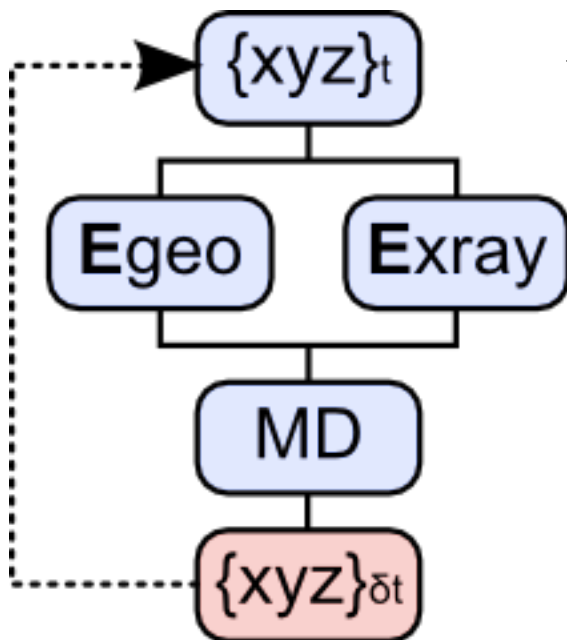
Molecular dynamics simulations produce a Boltzmann-weighted population of inter-related structures

MD simulations can be restrained with x-ray data

# Simulated Annealing/MD refinement

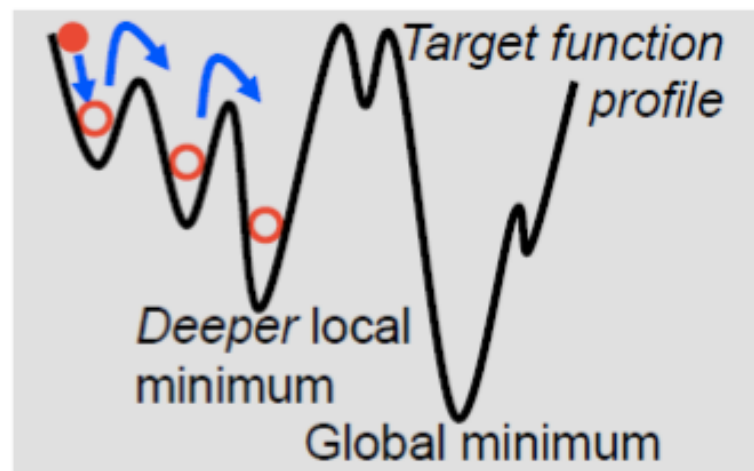
Geometric restraints:

$E_{bond}$   
 $E_{angle}$   
 $E_{dih}$   
...



X-ray restraint:

$$E_{X\text{-ray}} = \sum_{hkl} w_{hkl} (|F_{\text{obs}}(hkl)| - k |F_{\text{calc}}(hkl)|)^2$$



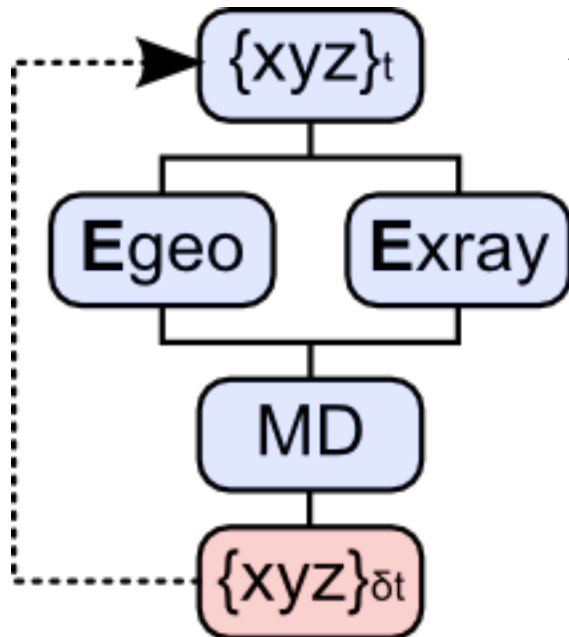
*Simulation temperature >2000K*  
*Trajectory resolves local minima*  
*'Final model' = end structure*



# “Time-averaged” MD refinement

Geometric restraints:

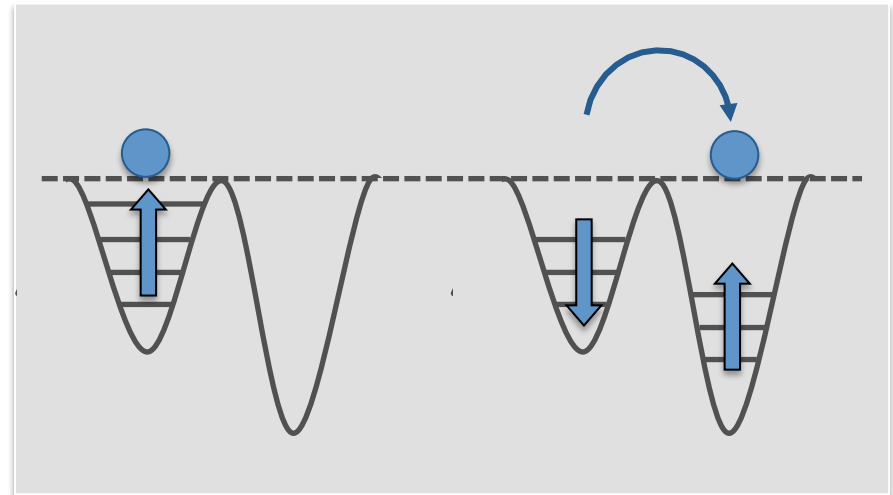
$E_{bond}$   
 $E_{angle}$   
 $E_{dih}$   
 ...



X-ray restraint:

$$E_{X\text{-ray}} = \sum_{hkl} w_{hkl} (|F_{\text{obs}}(hkl)| - k | \langle F_{\text{calc}}(hkl) \rangle |)^2$$

Sampling

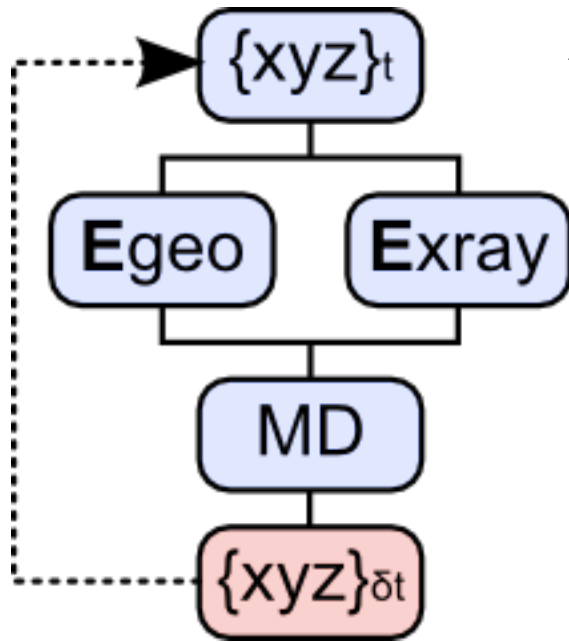


$$\langle F_{\text{calc}} \rangle_t = (1 - e^{-\Delta t / \tau_x}) F_{\text{calc}}^t + e^{-\Delta t / \tau_x} \langle F_{\text{calc}} \rangle_{t-\Delta t}$$

# “Time-averaged” MD refinement

Geometric restraints:

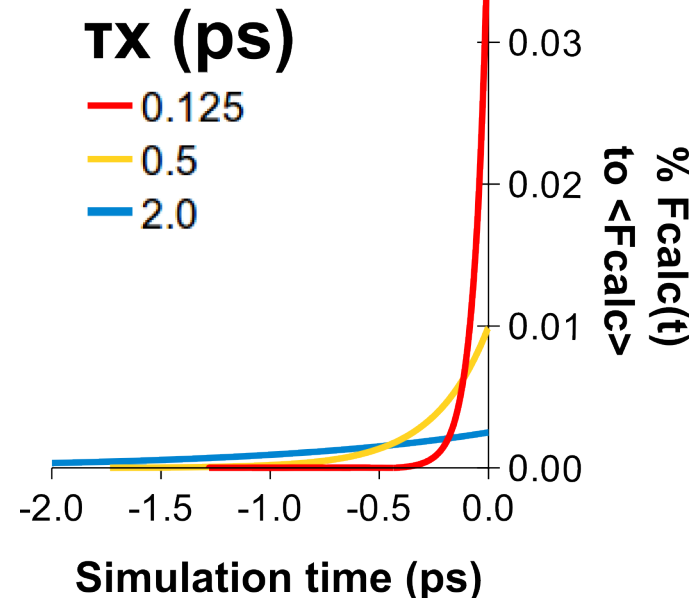
$E_{bond}$   
 $E_{angle}$   
 $E_{dih}$   
 ...



X-ray restraint:

$$E_{X\text{-ray}} = \sum_{hkl} w_{hkl} (|F_{obs}(hkl)| - k | \langle F_{calc}(hkl) \rangle |)^2$$

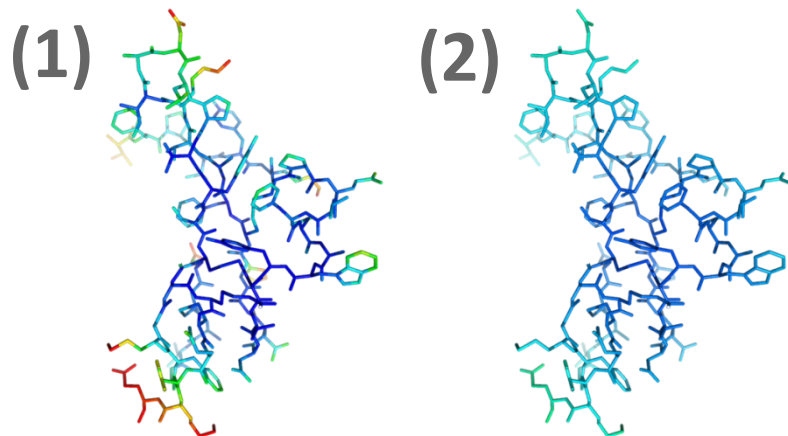
Running average



Simulation temperature = 300K  
 ‘Final model’ = trajectory ensemble

$$\langle F_{calc} \rangle_t = (1 - e^{-\Delta t / \tau_x}) F_{calc}^t + e^{-\Delta t / \tau_x} \langle F_{calc} \rangle_{t-\Delta t}$$

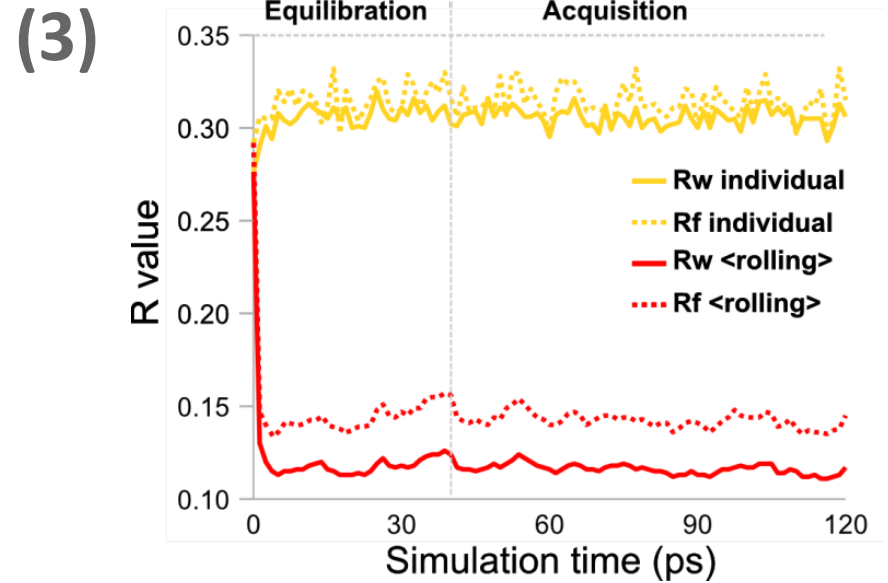
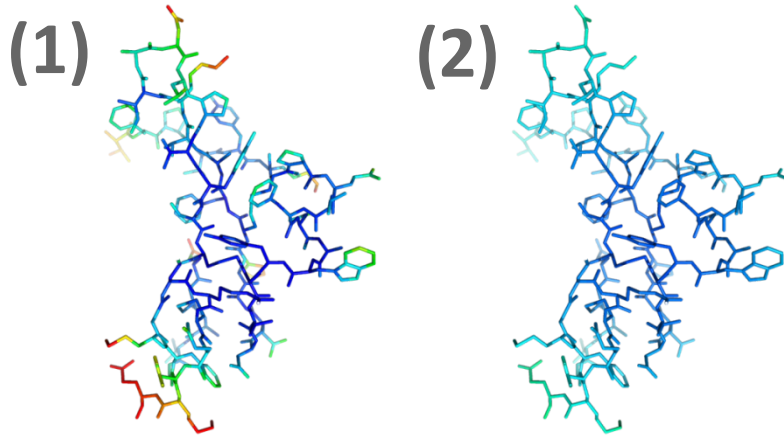
# ***Ensemble refinement with TA restraints***



**(1) Start: 'Traditional' structure**

**(2) Fit TLS / remove alt. conf.**

# Ensemble refinement with TA restraints



(1) Start: 'Traditional' structure

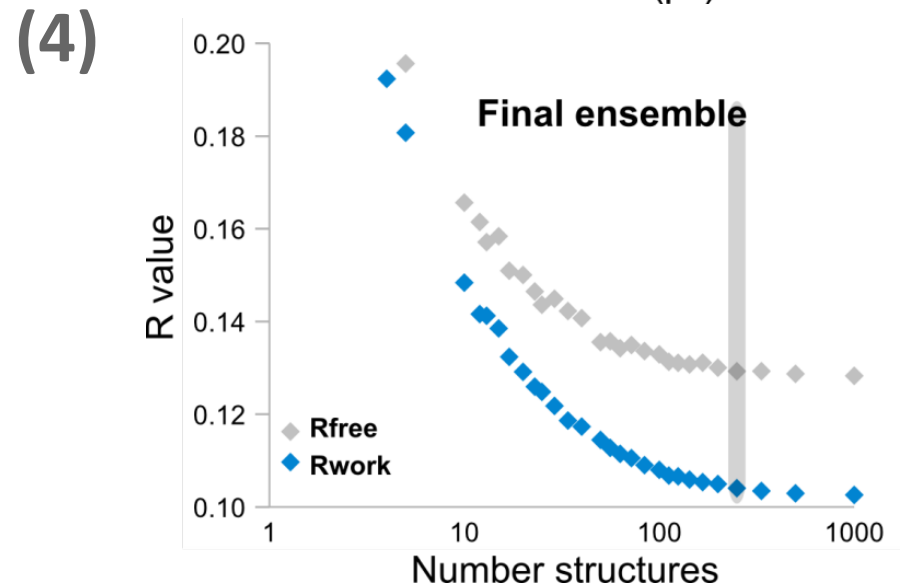
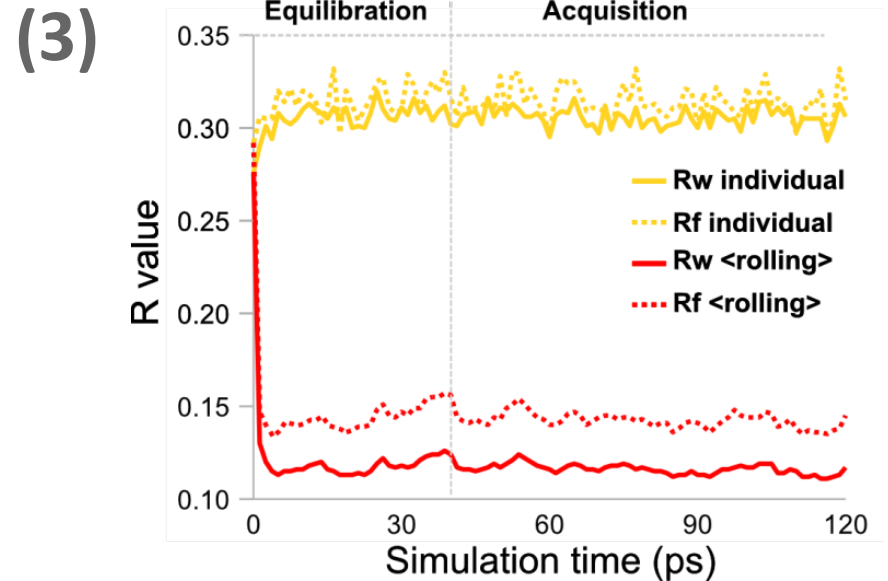
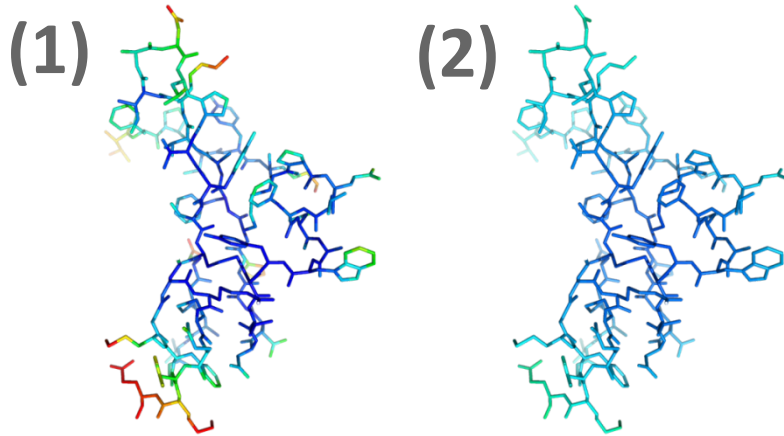
(2) Fit TLS / remove alt. conf.

(3) TA restrained MD simulation

Collect structure / 0.04ps

X-ray restraints accelerate sampling

# Ensemble refinement with TA restraints



(1) Start: 'Traditional' structure

(2) Fit TLS / remove alt. conf.

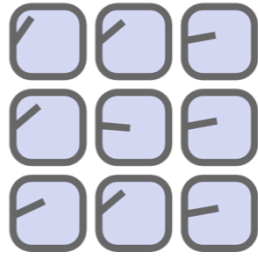
(3) TA restrained MD simulation

(4) Final ensemble

# *Molecular disorder / lattice disorder*



**Protein with  
local disorder**



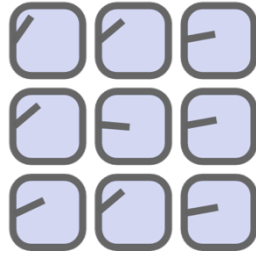
**Perfect crystal  
lattice**



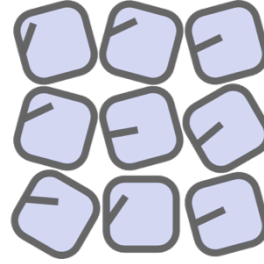
# *Molecular disorder / lattice disorder*



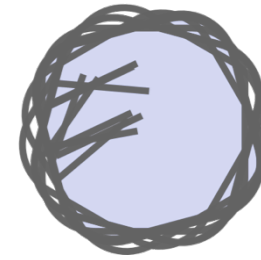
**Protein with  
local disorder**



**Perfect crystal  
lattice**



**'Real' crystal  
lattice**



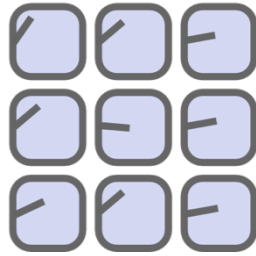
**Averaged  
data**



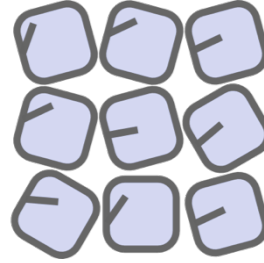
# *Molecular disorder / lattice disorder*



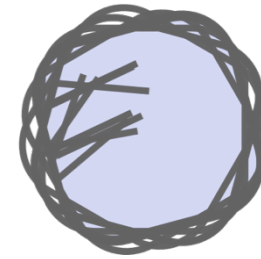
**Protein with  
local disorder**



**Perfect crystal  
lattice**



**'Real' crystal  
lattice**



**Averaged  
data**

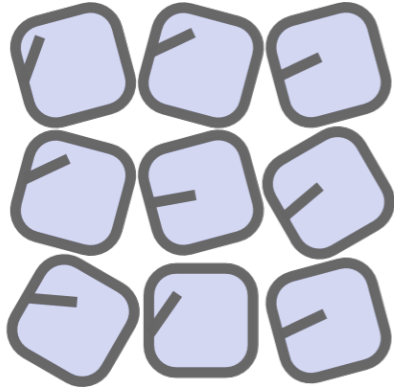


**Deconvolute:  
molecular disorder from  
lattice disorder**



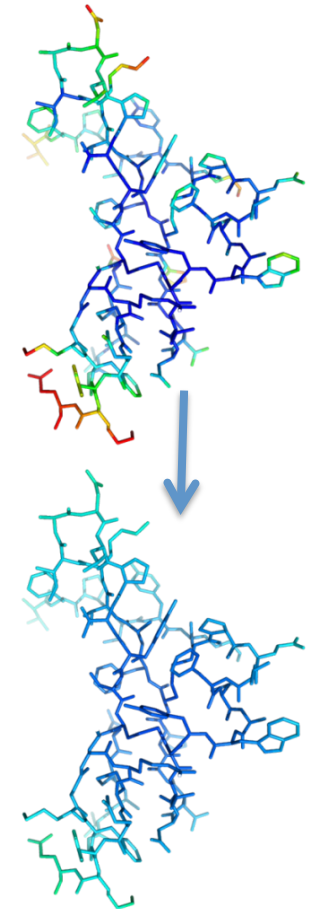
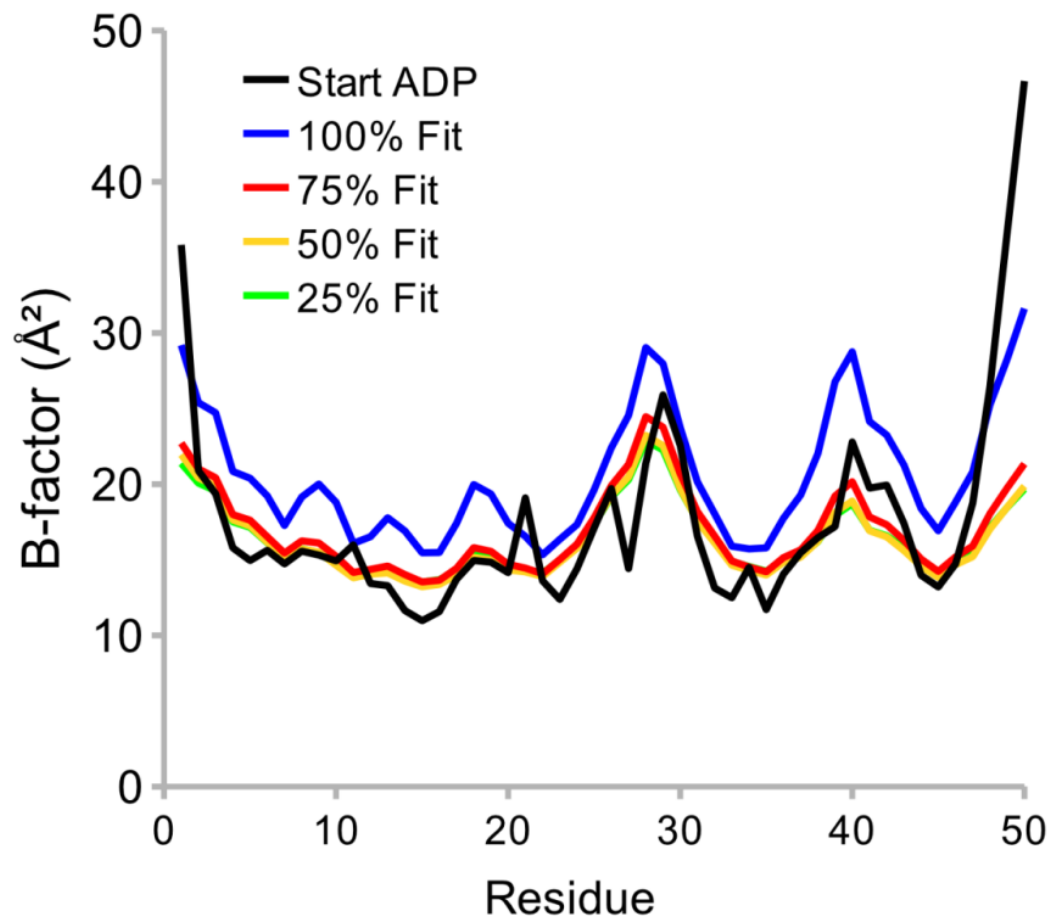


# ***Rigid body disorder modelled with TLS***



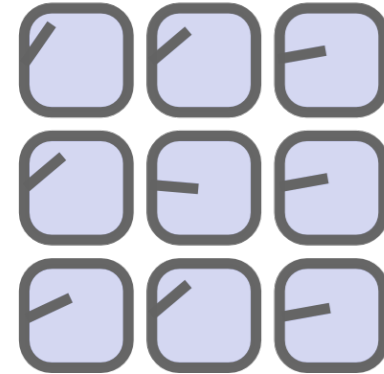
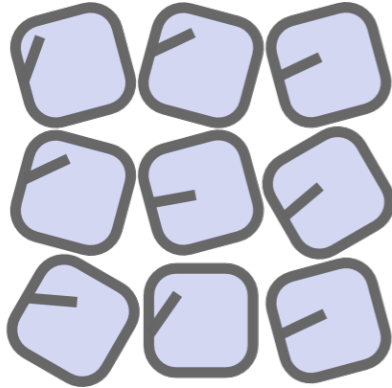
- **Lattice distortions, domain motions...**
- Model rigid body motions with ***TLS*** model

# *Rigid body disorder modelled with TLS*



Extract core rigid body motion by excluding atoms with large local fluctuations (defined as deviations from  $B_{\text{TLS}}$ ).

# *Local disorder sampled within restrained MD*

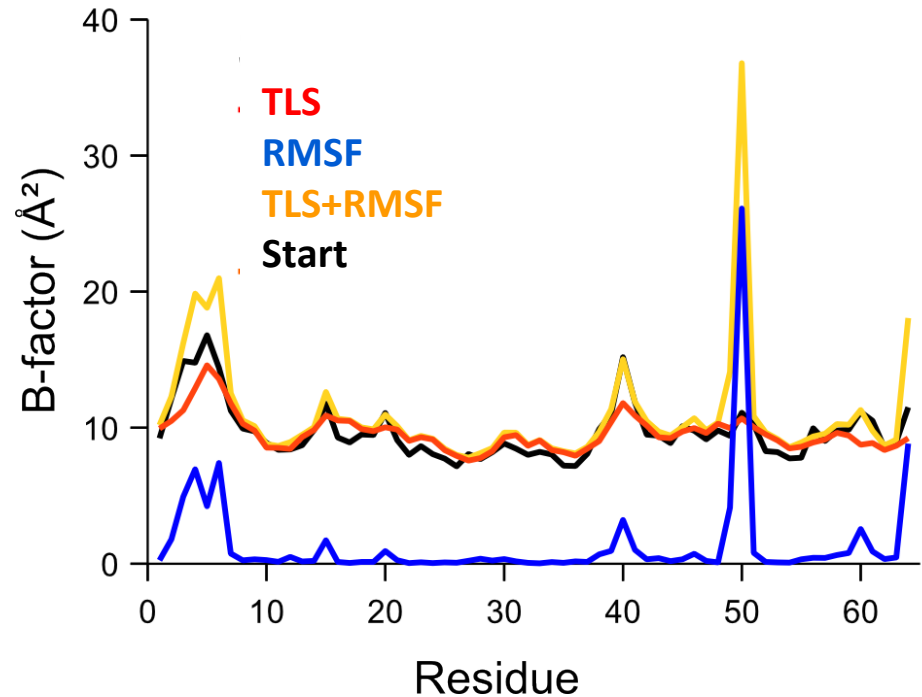


- **Lattice distortions, domain motions...**
- Model rigid body motions with TLS model
  - 20 parameters per group
  - 1 group / domain
- Fit TLS to starting structure B-factors

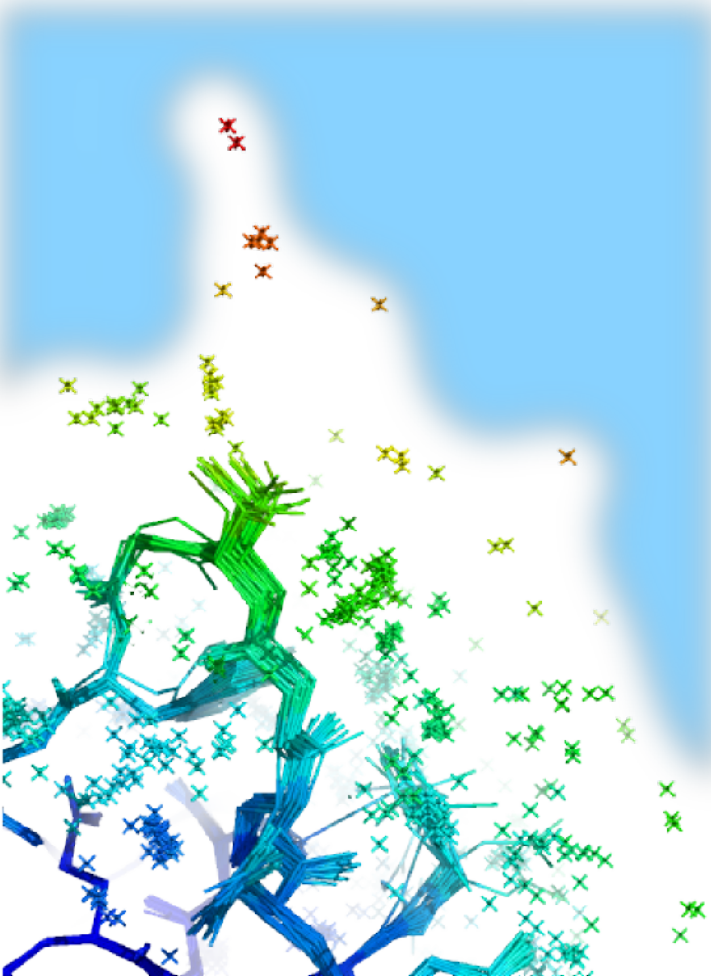
- **Side-chains, loops...**
- Local disorder sampled in MD
- MD simulation restrained with X-ray data

# Composite B-factor sum of TLS and atomic fluctuations

- **TLS**
  - Core TLS model
- **RMSF**
  - Atomic fluctuations in ensemble
- **TLS+RMSF**
  - Total disorder in ensemble
- **Start**
  - Input single structure ADPs

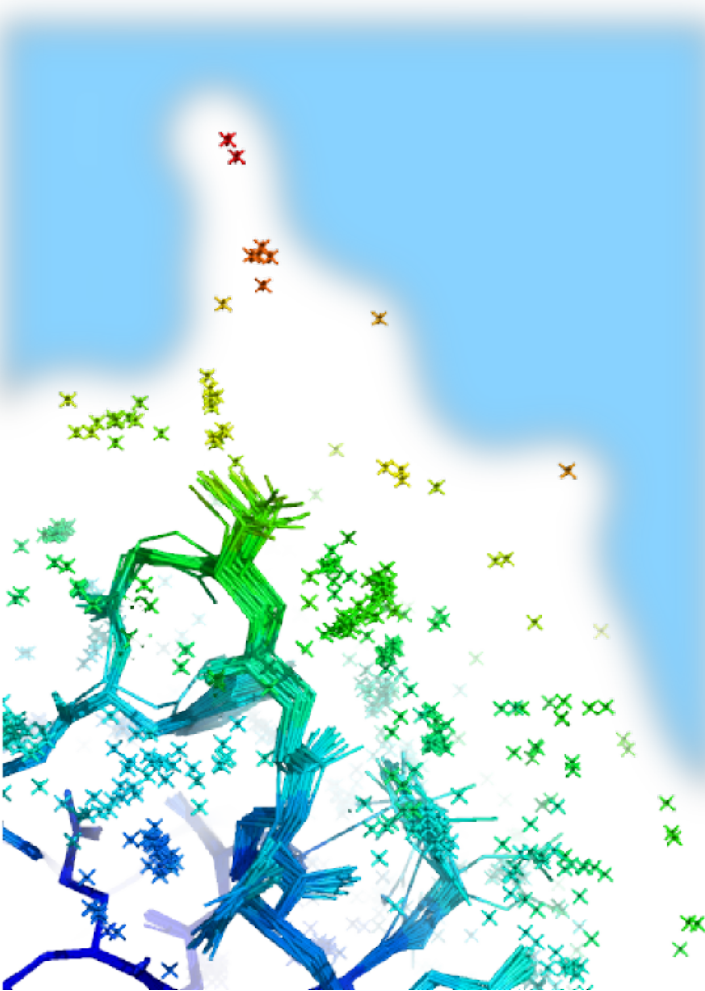


# *Dual explicit-bulk solvent model*



- Explicit solvent
  - Model with explicit atoms
  - Water picked every 250 steps
    - “standard” rules:
      - >  $3\sigma$  in difference map
      - < 3 Å distances
  - B-factor from nearest TLS group

# Dual explicit-bulk solvent model

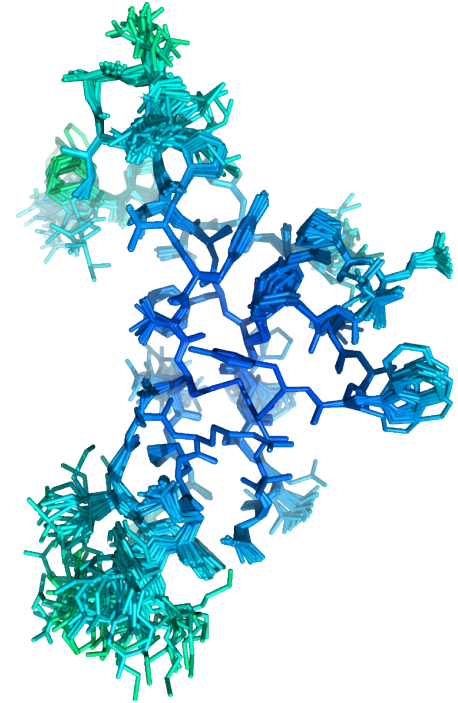


- Explicit solvent
  - Model with explicit atoms
  - Water picked every 250 steps
    - “standard” rules:
      - > 3  $\sigma$  in difference map
      - < 3 Å distances
  - B-factor from nearest TLS group
- Bulk solvent
  - Model with ‘density mask’

$$\langle F_{mask} \rangle_t = (1 - e^{-\Delta t / \tau_x}) F_{mask}^t + e^{-\Delta t / \tau_x} \langle F_{mask} \rangle_{t-\Delta t}$$

# *Development of ensemble refinement*

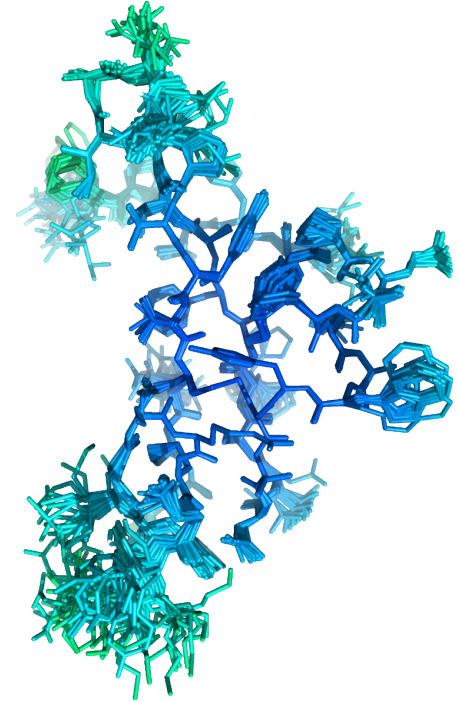
- PHENIX
- Time-averaged x-ray restrained MD
- TLS fitting
- Explicit- and bulk- solvent model
- Maximum-likelihood target function



**Phenix**

# *Development of ensemble refinement*

- Tested with 20 datasets
- Resolution: 1 - 3 Å
- ASU size: 50-1000 residues
- CPU time: 7 - 100 hours
- 50 – 500 models / ensemble



**Phenix**



# Ensemble refinement reduces $R_{\text{free}}$

- *$R_{\text{free}}$ : ensemble vs phenix.refine*

- $R_{\text{free}}$  reduced in all cases

- 4.9% (max)

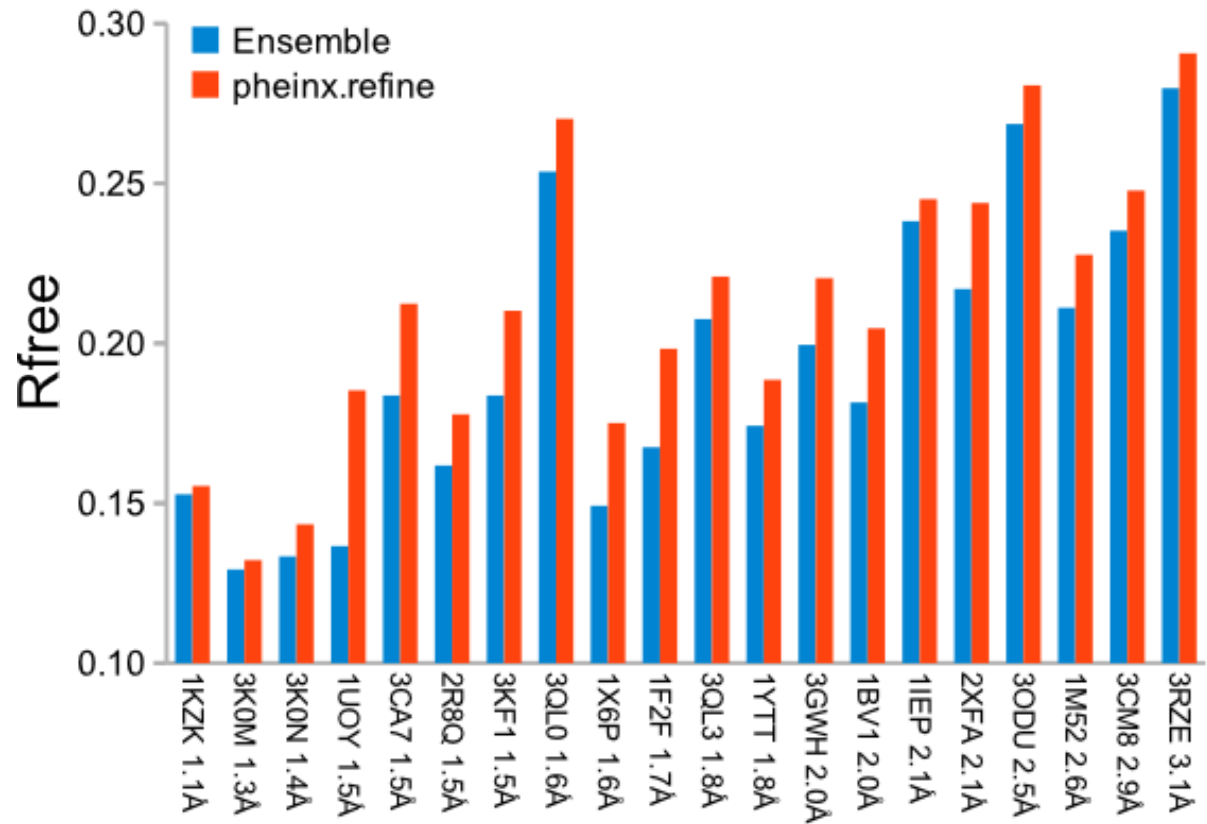
- 0.3% (min)

- 1.8% (mean)

- $R_{\text{f}}/R_{\text{w}}$  ratio (mean):

- = 1.23 phenix.refine

- = 1.25 ensemble



# Ensemble refinement reduces $R_{\text{free}}$

- **$R_{\text{free}}$ : ensemble vs phenix.refine**

- $R_{\text{free}}$  reduced in all cases

- 4.9% (max)

- 0.3% (min)

- 1.8% (mean)

- $R_f/R_w$  ratio (mean):

- = 1.23 phenix.refine

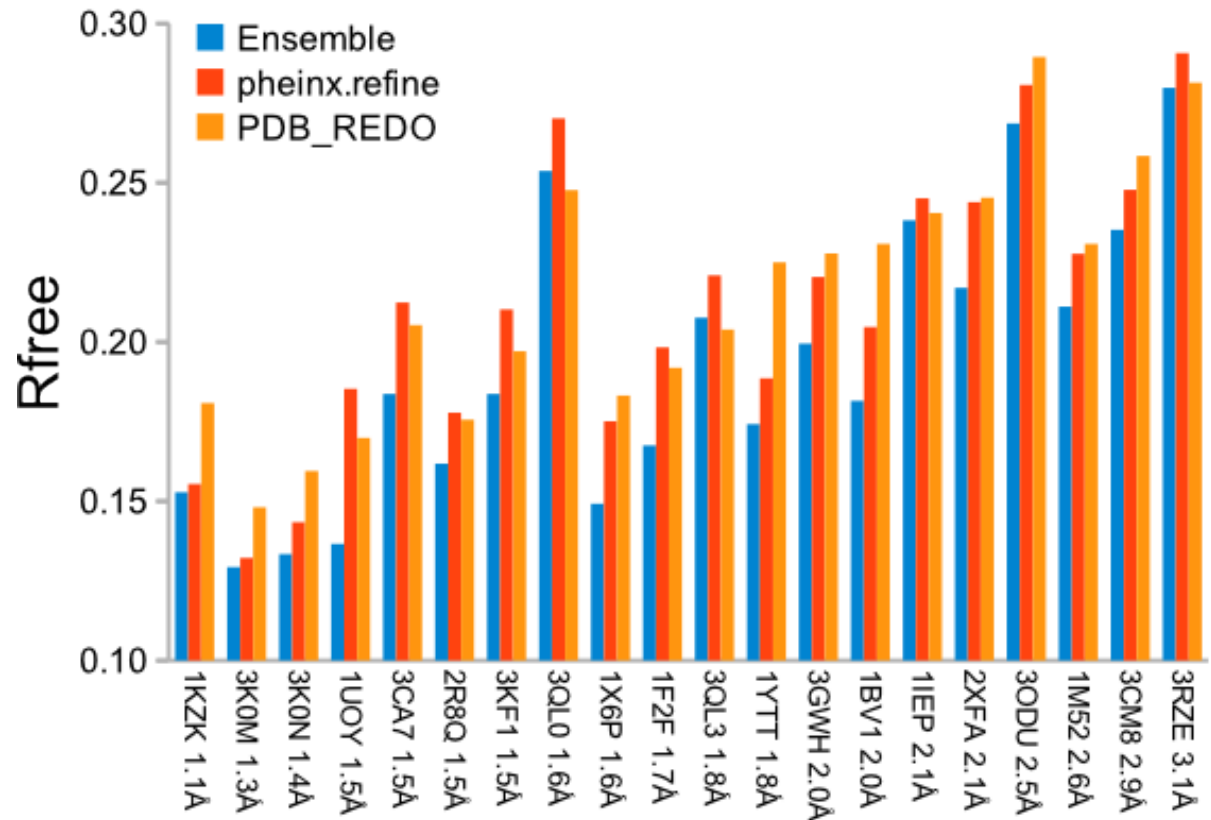
- = 1.25 ensemble

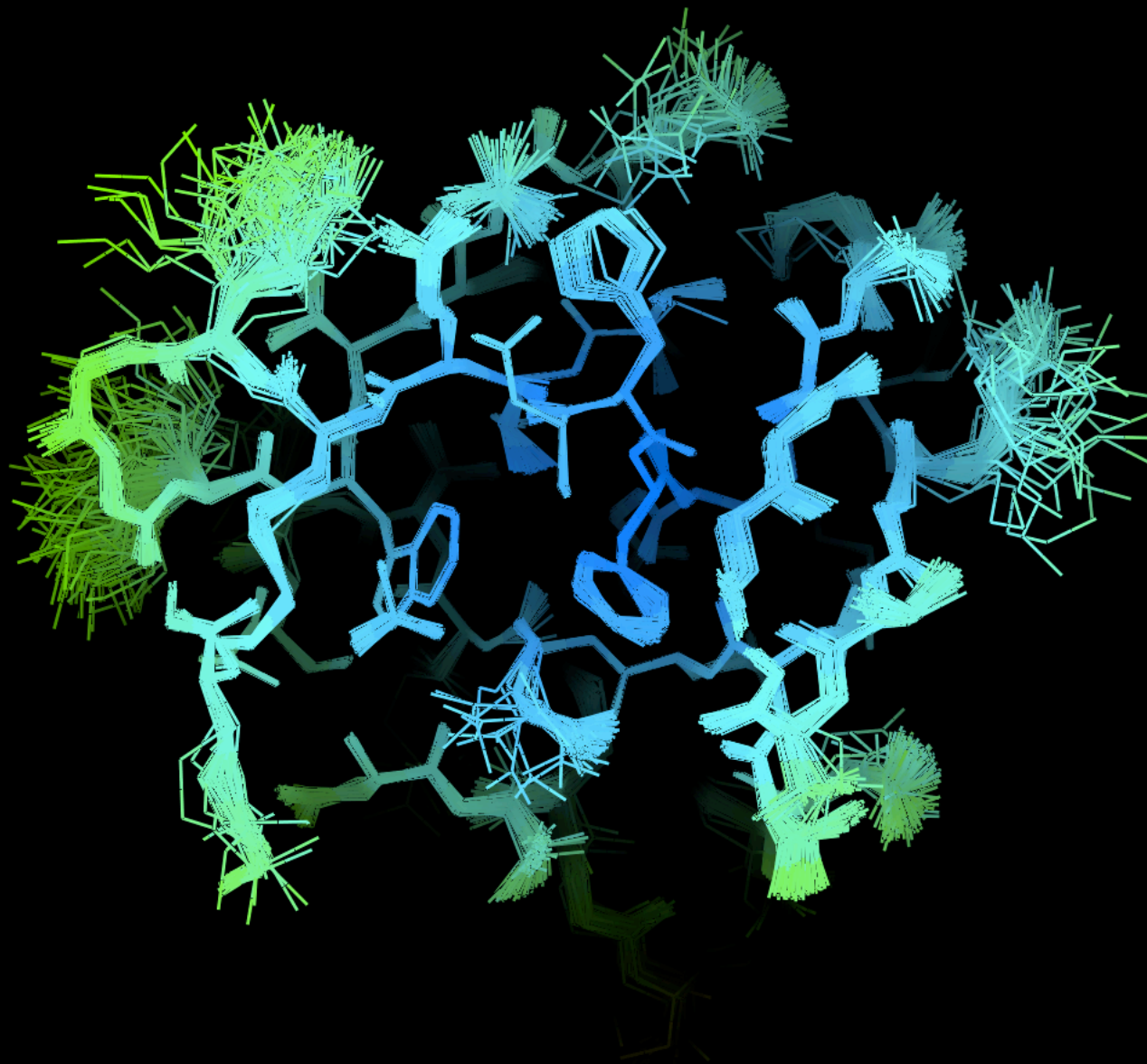
- **$R_{\text{free}}$ : ensemble vs PDB\_REDO**

- 5.1% (max)

- + 0.6% (min)

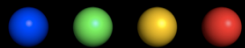
- 2.1% (mean)





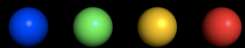
**B factor**

5 - 20Å<sup>2</sup>

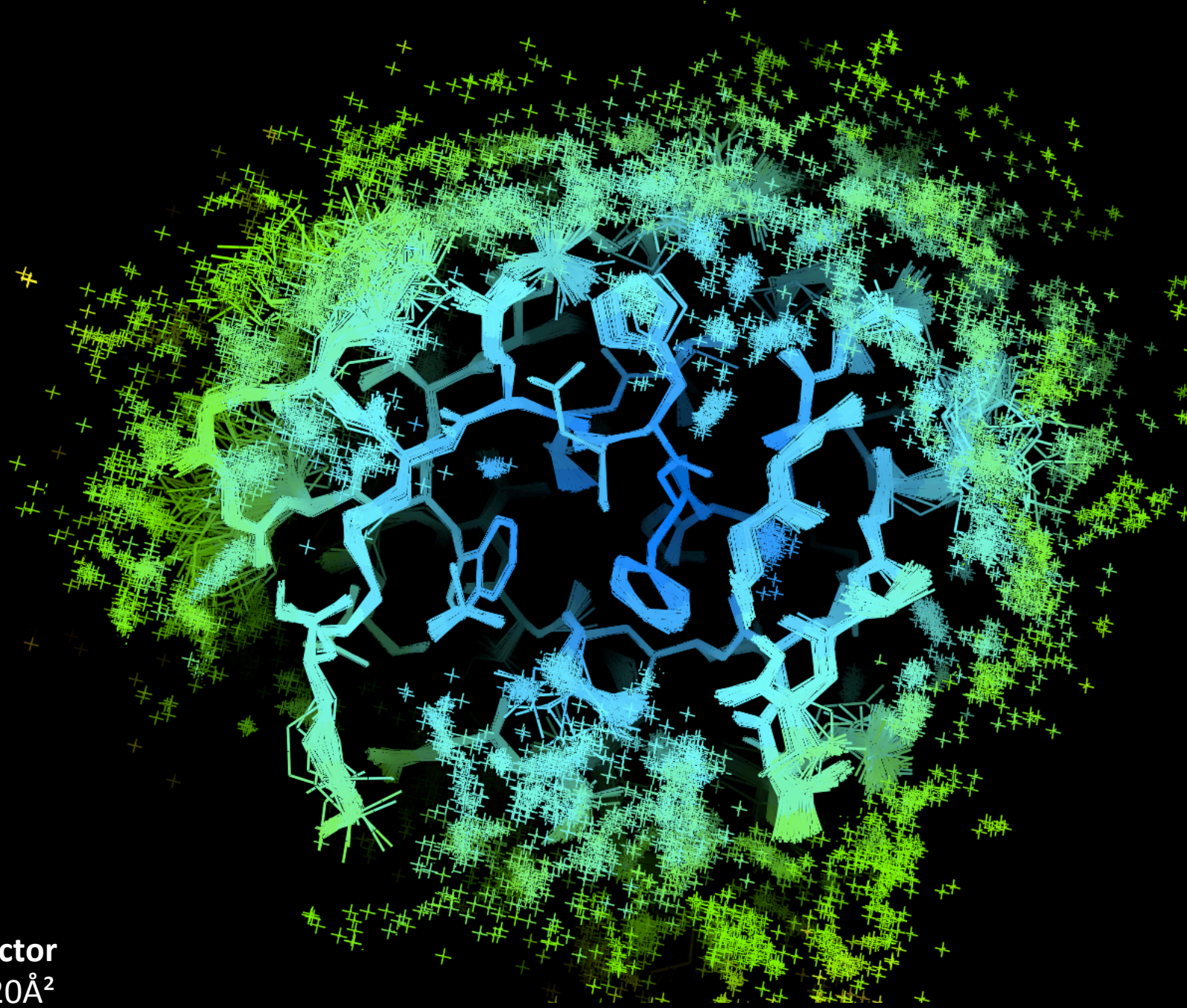


1uoy.pdb | 188 ensemble | 40ps acquisition time | 1 tls group

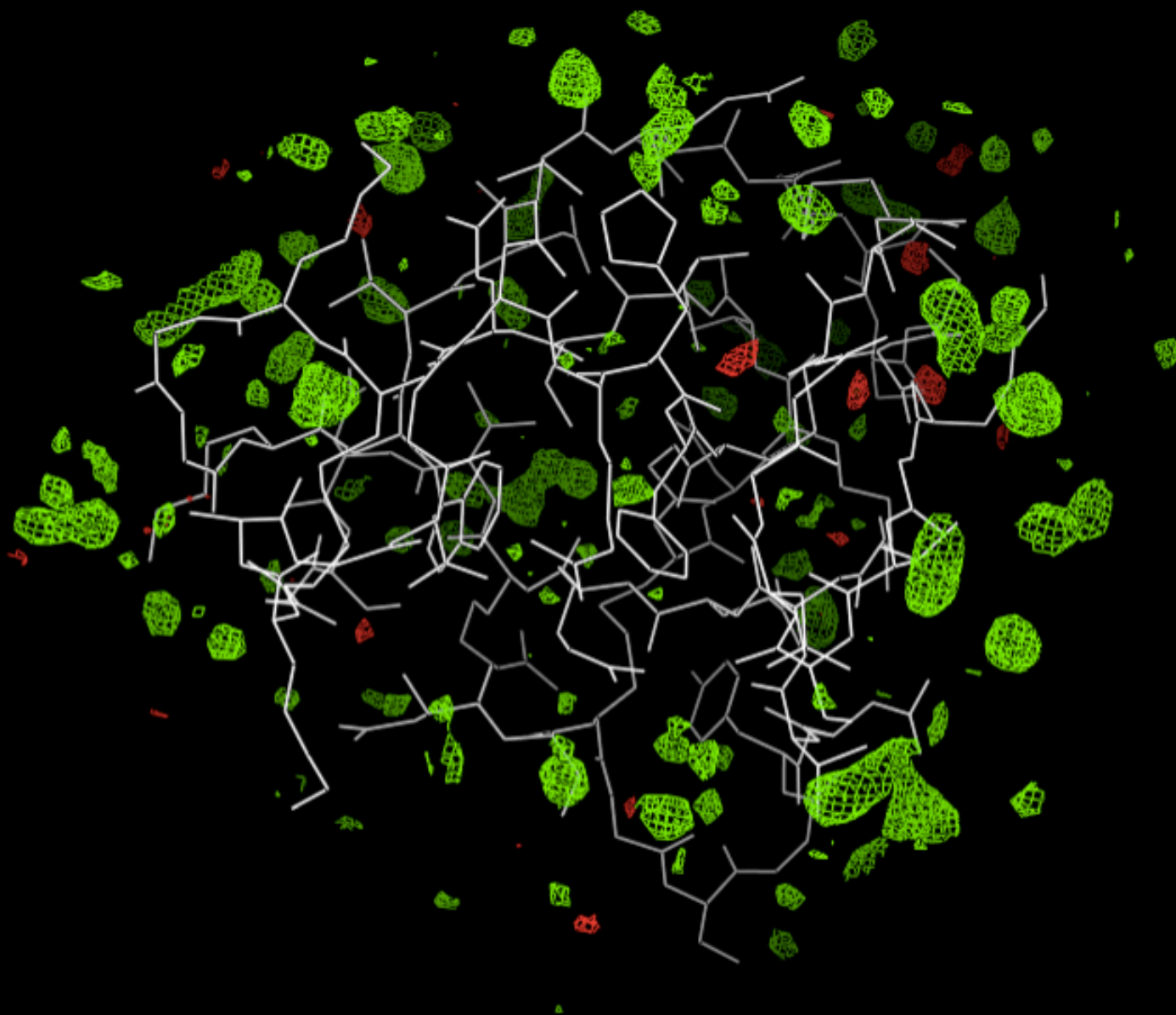
**B factor**  
5 - 20Å<sup>2</sup>



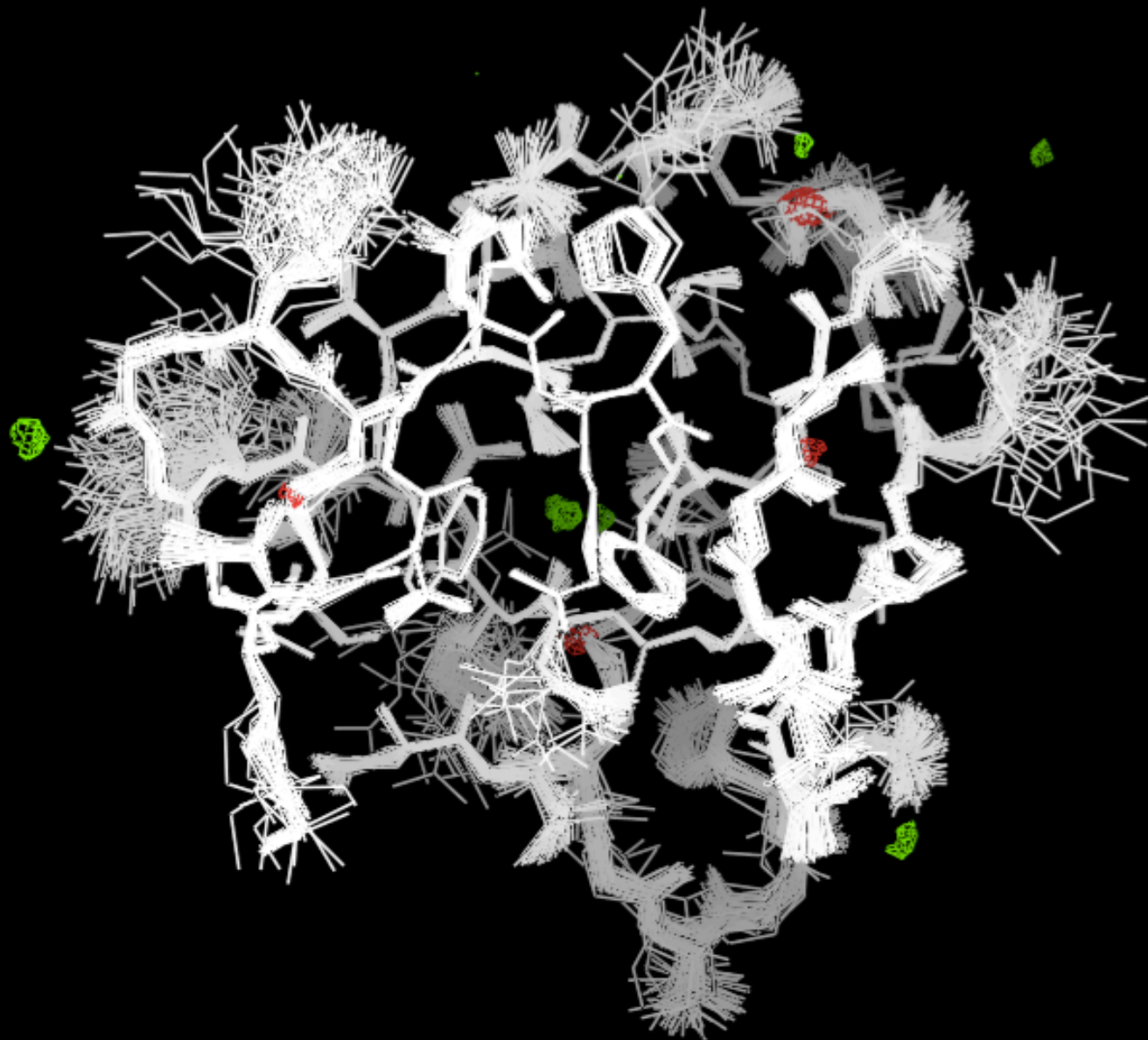
**1uoy.pdb | 188 ensemble | 40ps acquisition time | 1 tls group**







1uoy.pdb | phenix.refine | 1 tls group | mFo-DFc  $\pm 0.49$  e/ $\text{\AA}^3$  (3.00  $\sigma$ )



1uoy.pdb | 188 ensemble | 1 tls group | mFo-DFc  $\pm 0.49$  e/Å<sup>3</sup> (4.27  $\sigma$ )

# Improved real-space correlation

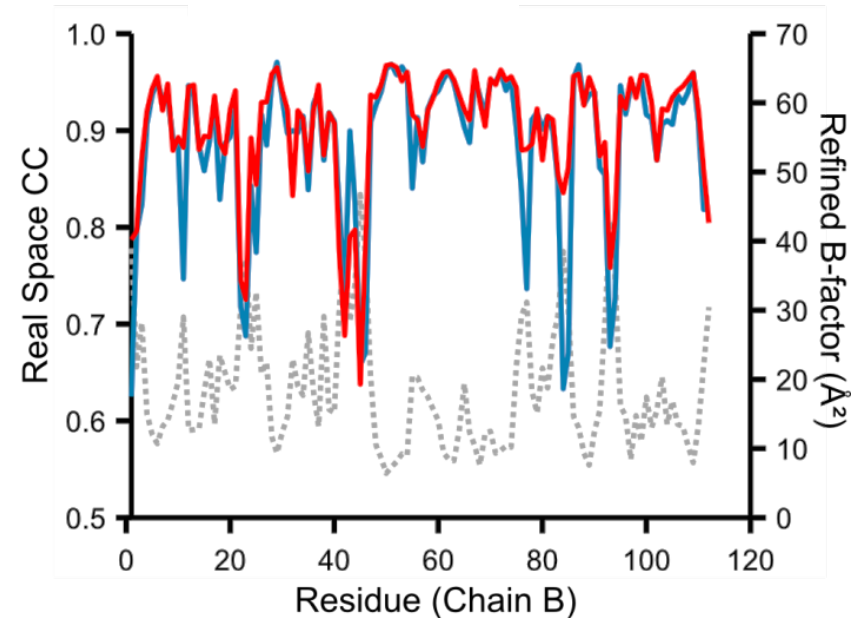
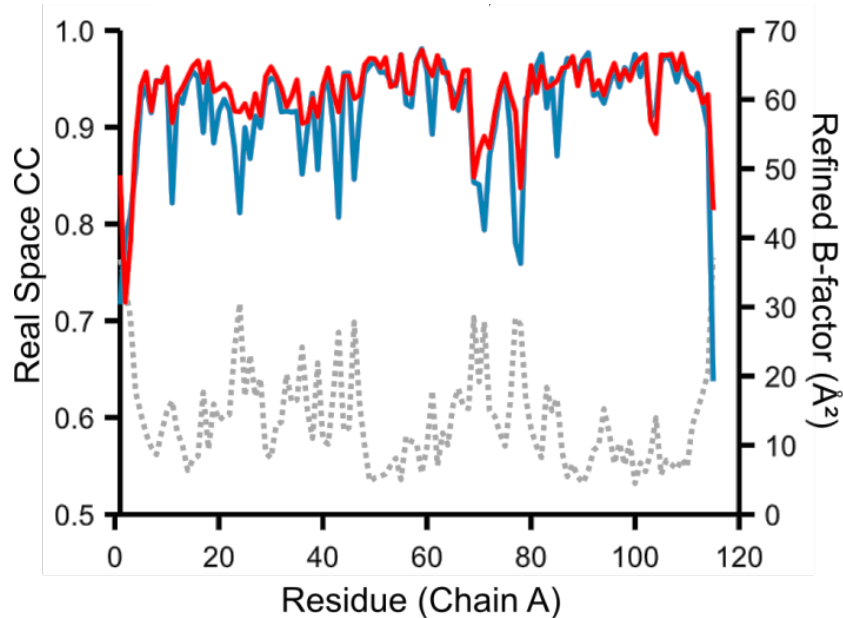
Burling *et al.* (1996):

Excellent experimentally phased data for MBP: 1YTT (1.8-Å res.)

Ensemble

Phenix.refine

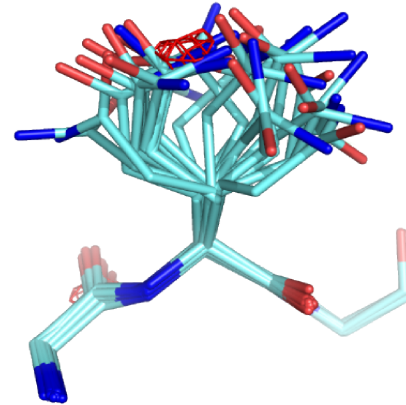
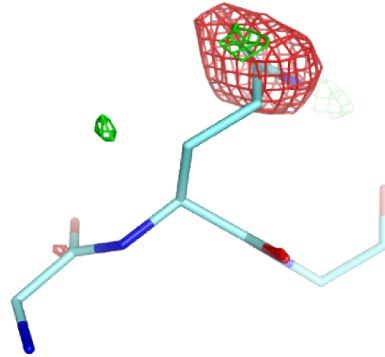
B-factor



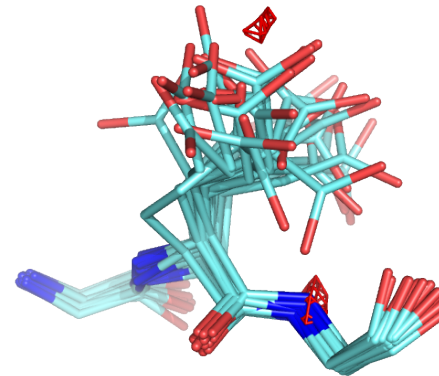
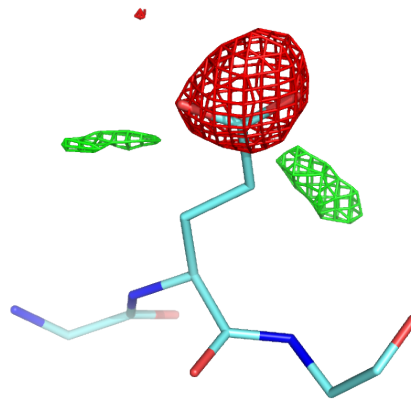
	<b>Rwork</b>	<b>Rfree</b>	<b>Real Space CC</b>
Ensemble	0.139	0.174	0.903
phenix.refine	0.166	0.189	0.895
PDB	0.185	0.206	0.873

# Disordered side chain in MBP (1YTT)

Gln167 (A)



Glu117 (A)



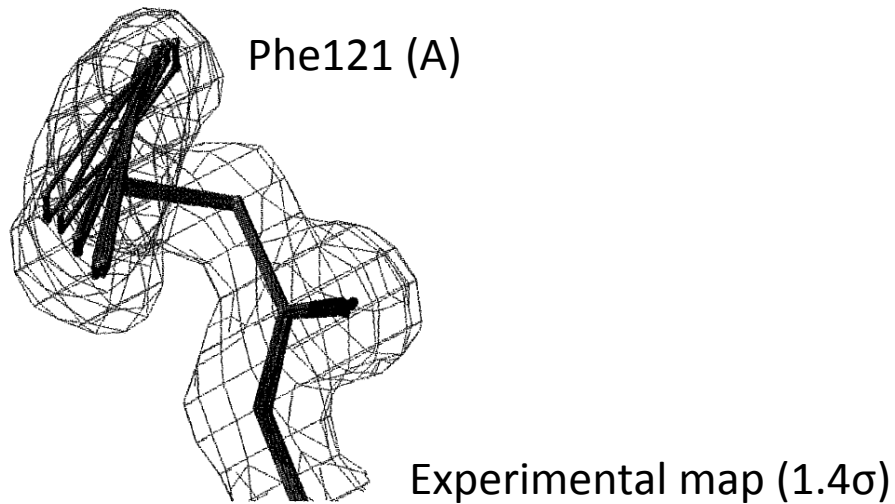
Phenix.refine  
Diff. vector map ( $3\sigma$ )

Ensemble  
Diff. vector map ( $3\sigma$ )

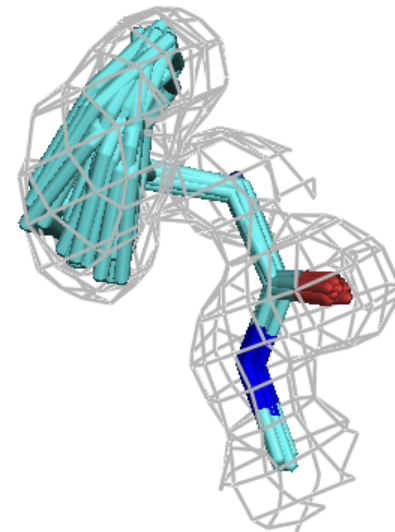


# Anisotropic side chain in MBP (1YTT)

Burling *et al.* (1996)



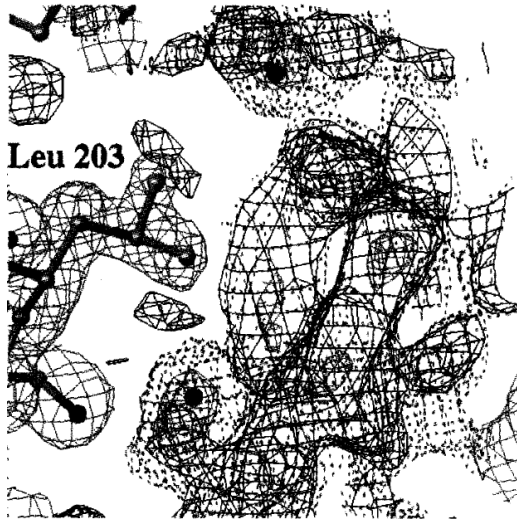
**Multi-conformer**  
(Rfree 20.3%)



**Ensemble**  
(Rfree 17.4%)

# Diffuse solvent in MBP (1YTT)

Burling *et al.* (1996)



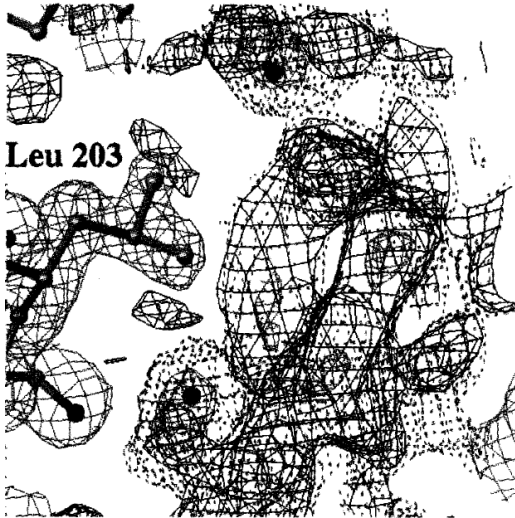
(Fig. 2B). The electron density around Leu<sup>203</sup> in protomer A suggests a network of four to five partially disordered water molecules (23). The disorder is presum-

Published:

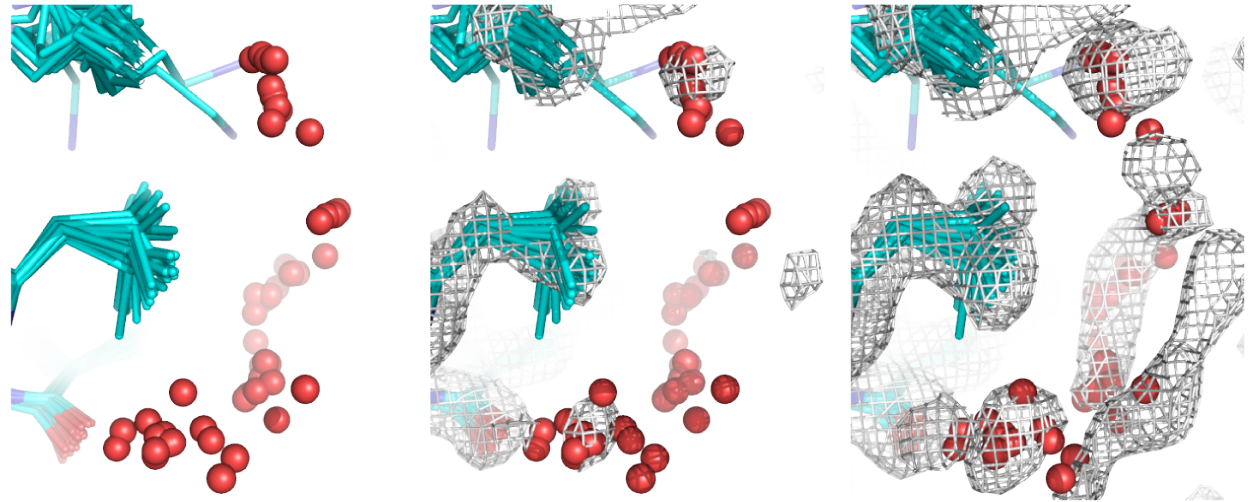
Experimental map  
(1.4 $\sigma$  and 0.7 $\sigma$ )

# Diffuse solvent in MBP (1YTT)

Burling *et al.* (1996)

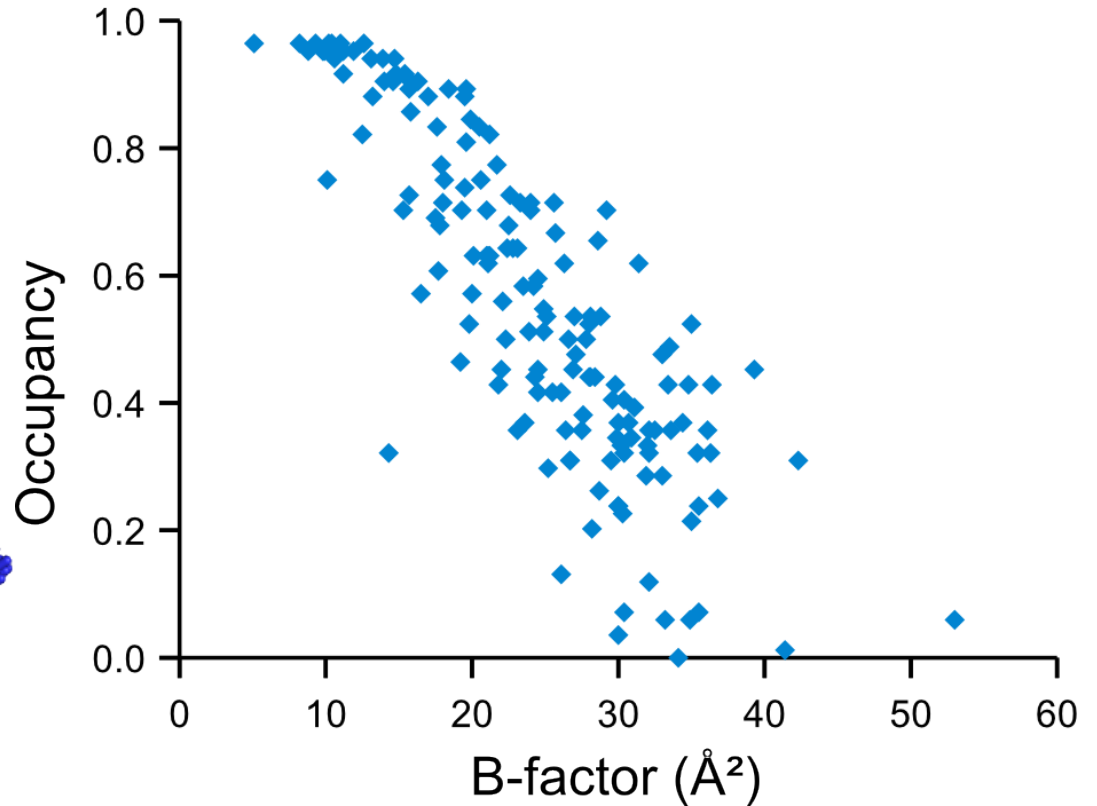
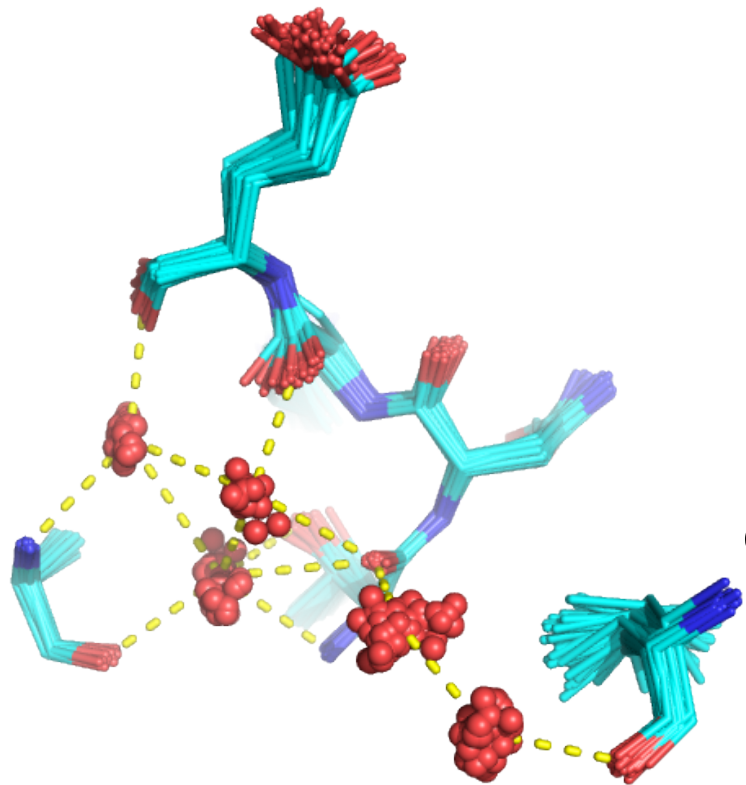


Published:  
Experimental map  
( $1.4\sigma$  and  $0.7\sigma$ )



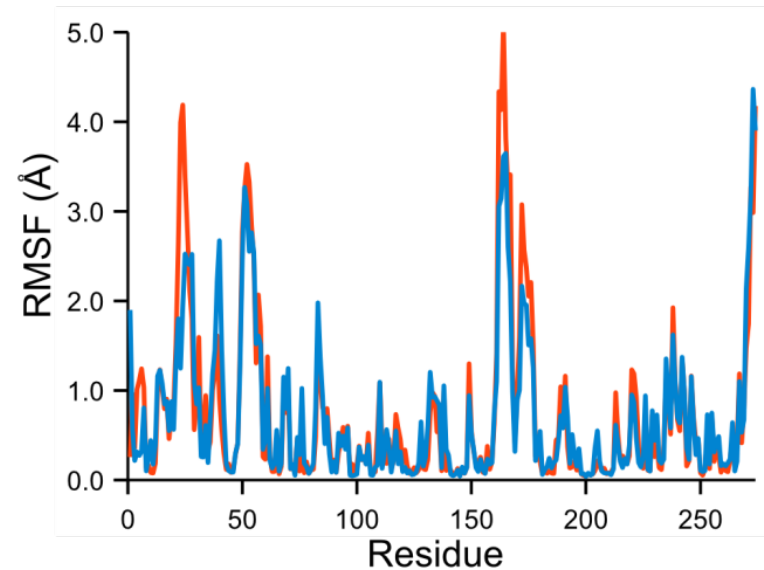
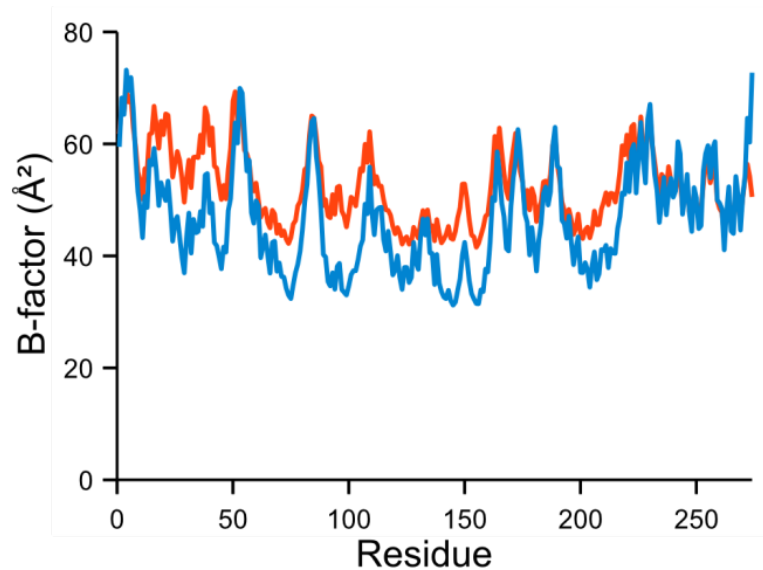
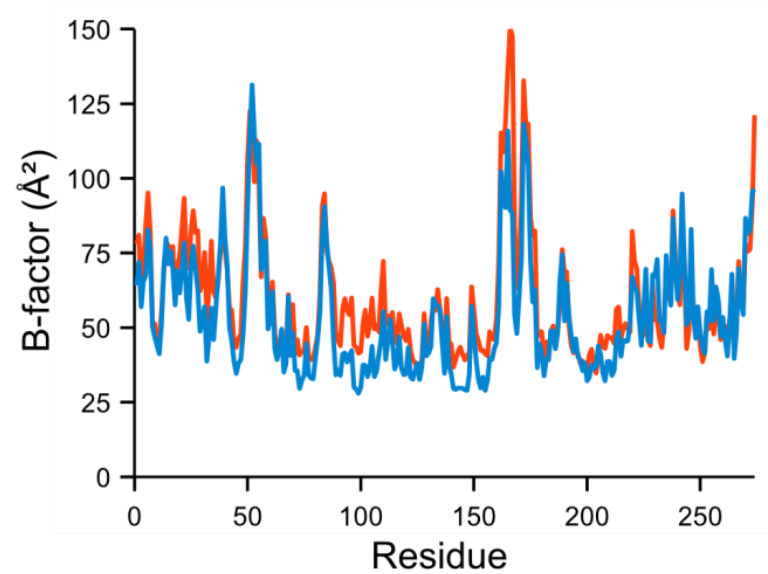
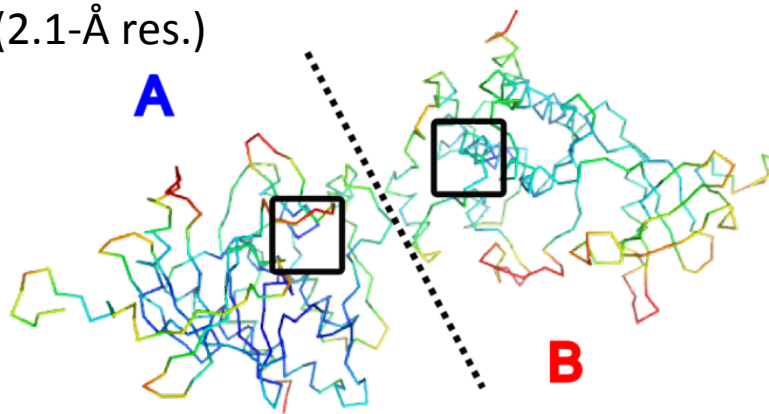
Ensemble:  
Experimental map      ( $1.4\sigma$ )      ( $0.7\sigma$ )

# Explicit waters in MBP (1YTT)

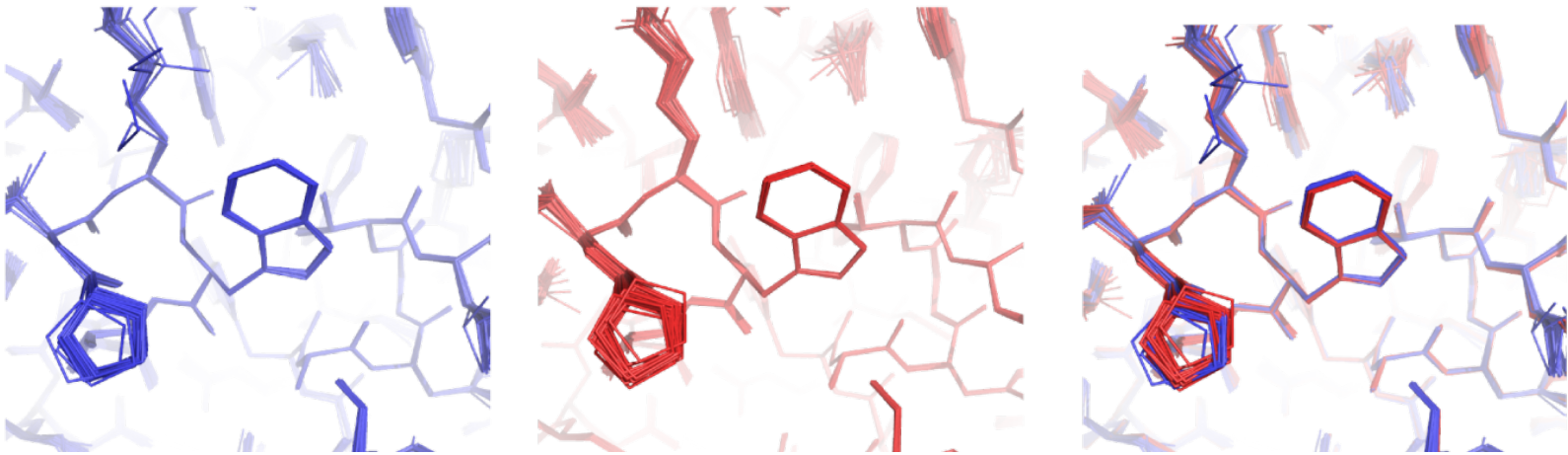


# How similar are NCS copies?

1IEP (2.1-Å res.)



# NCS copies show similar distributions

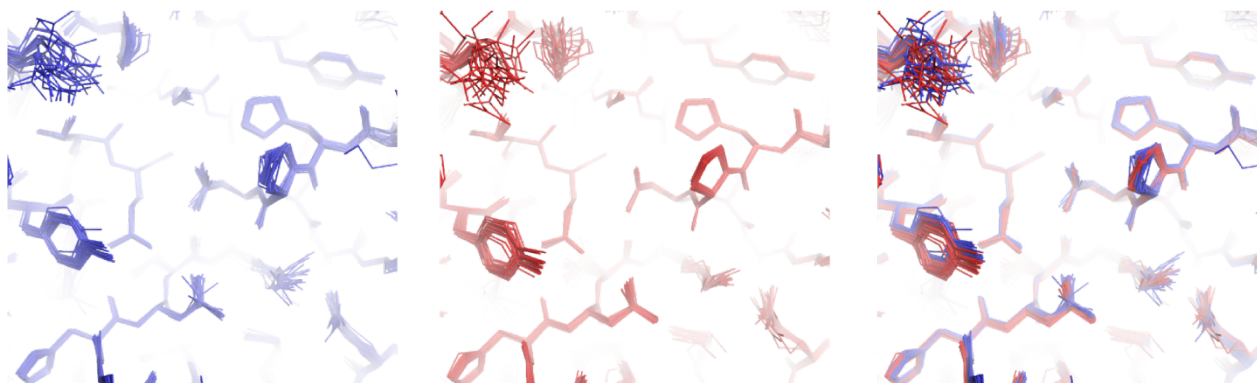
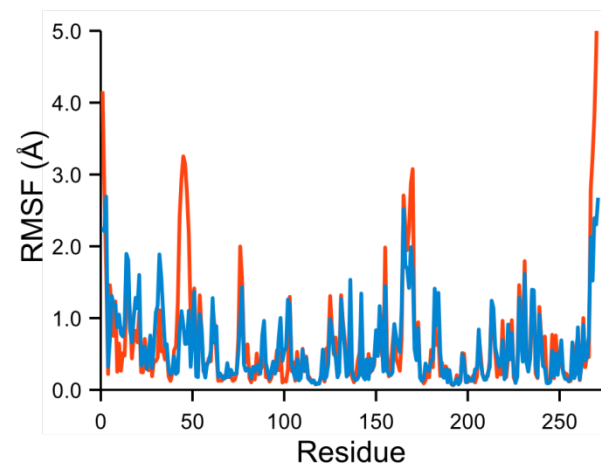
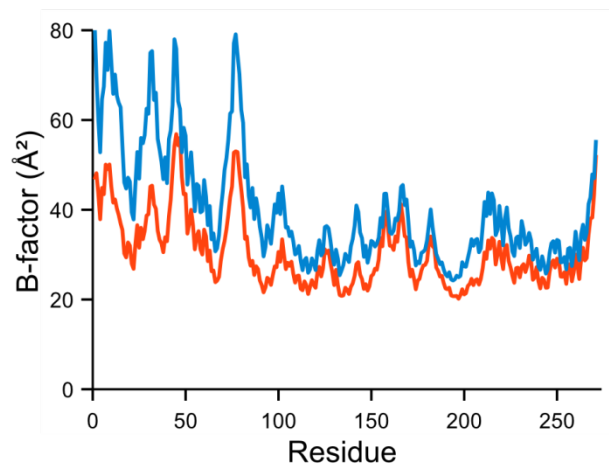
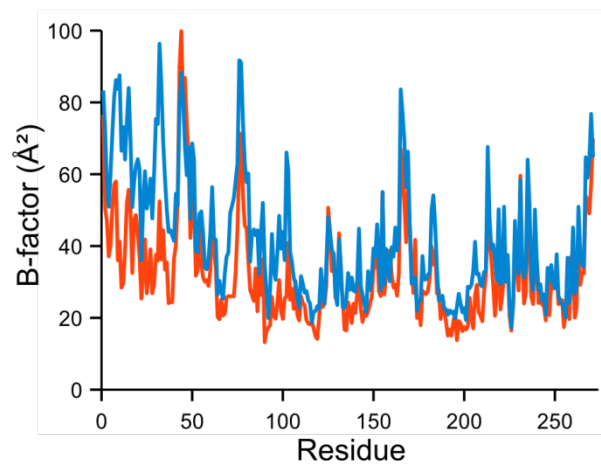
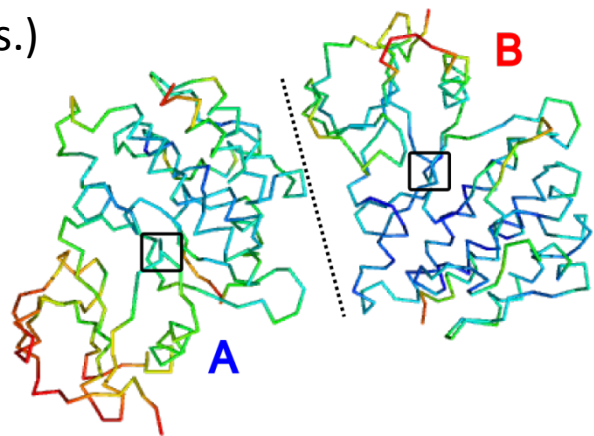


1IEP (2.1-Å res.)

- Global TLS model accounts for differences in packing of copies
- Local fluctuations are similar between the two copies

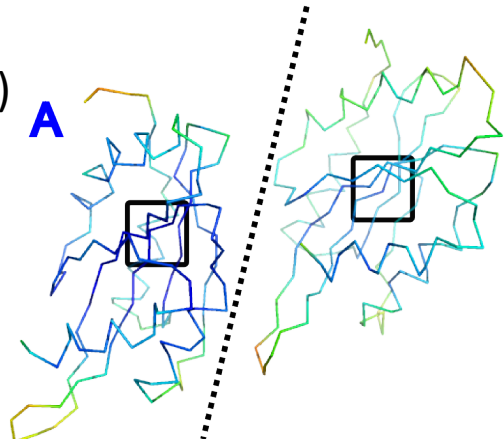


1M52 (2.6-Å res.)

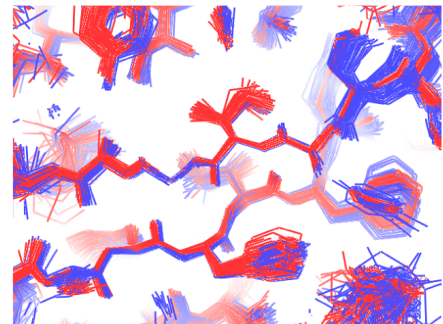
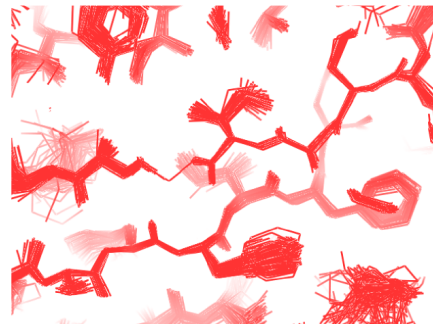
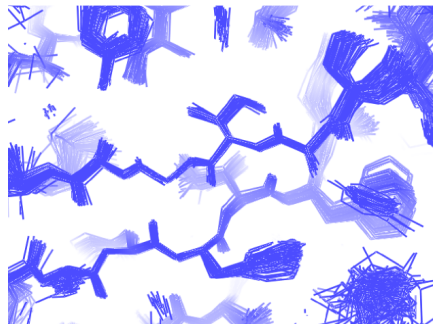
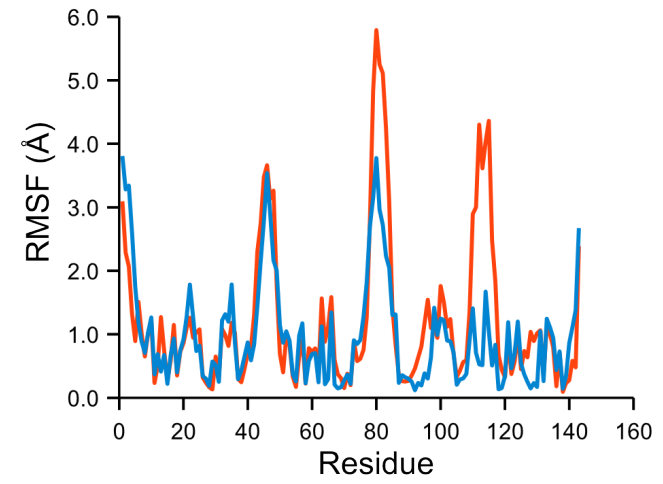
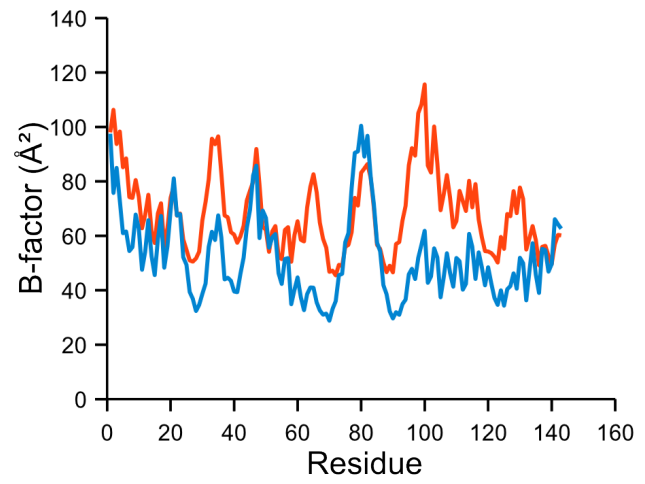
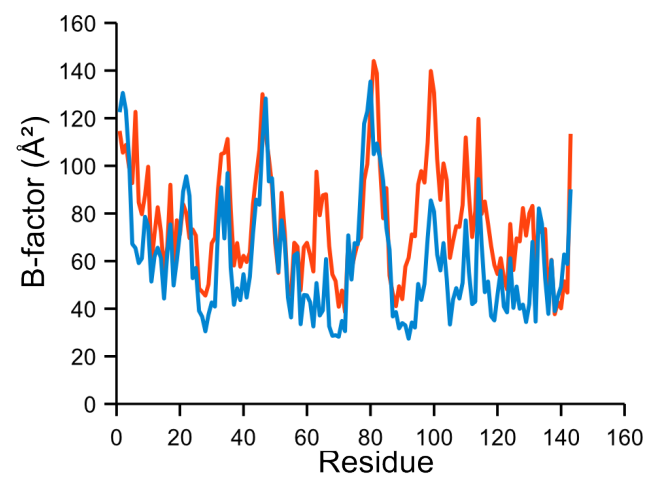


2XFA (2.1-Å res.)

**A**

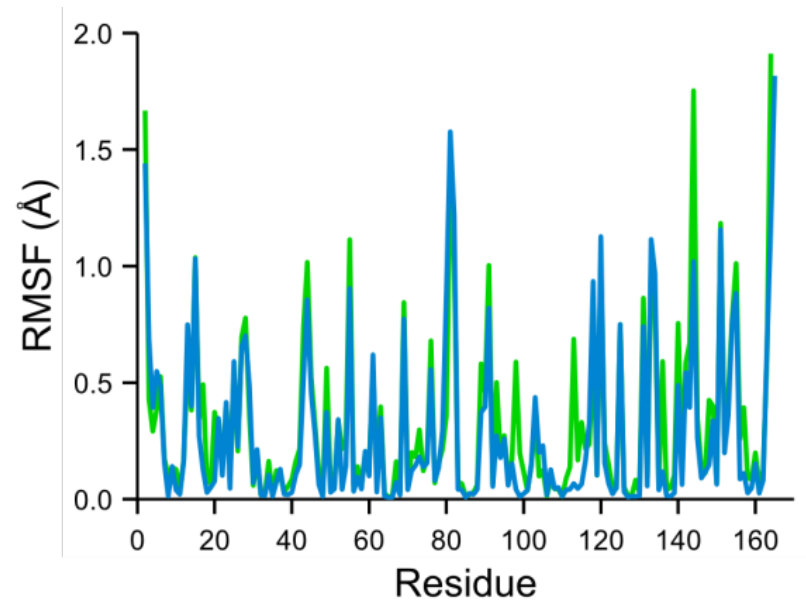
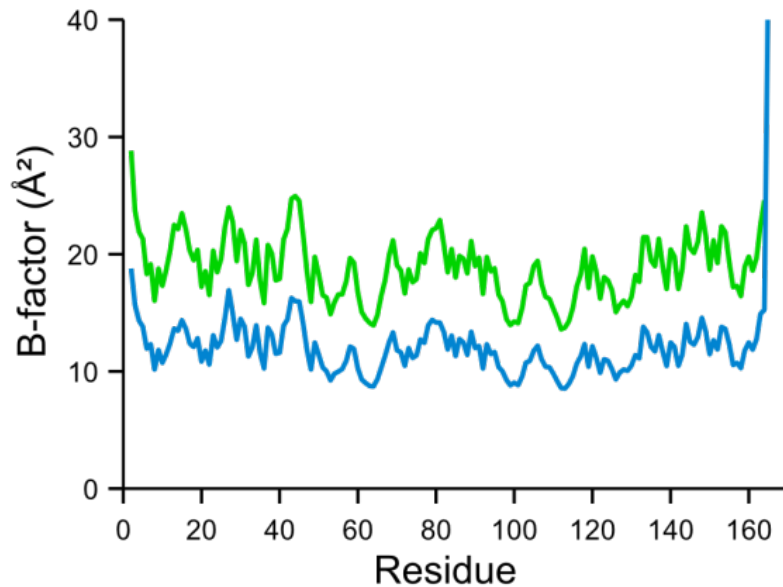


**B**





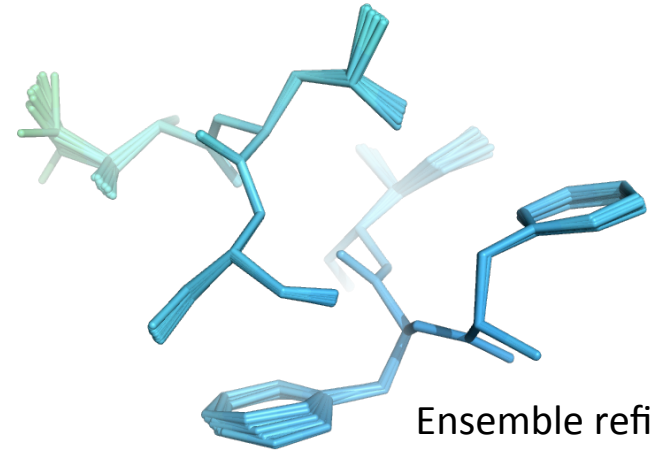
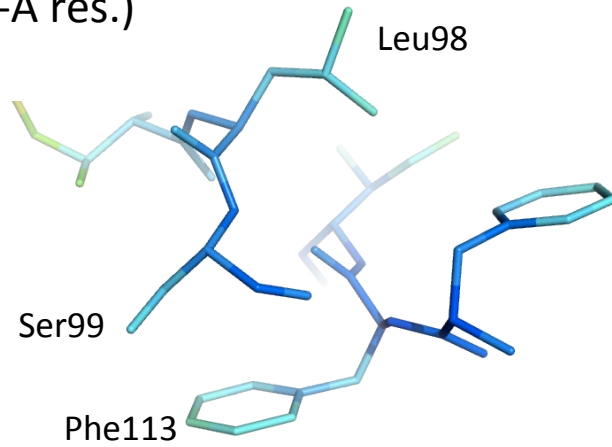
# Isomorphous crystals at 100 and 288 K



3K0N & 3K0M: Proline isomerase (Cyclophilin A)  
Fraser *et al.* (2009)

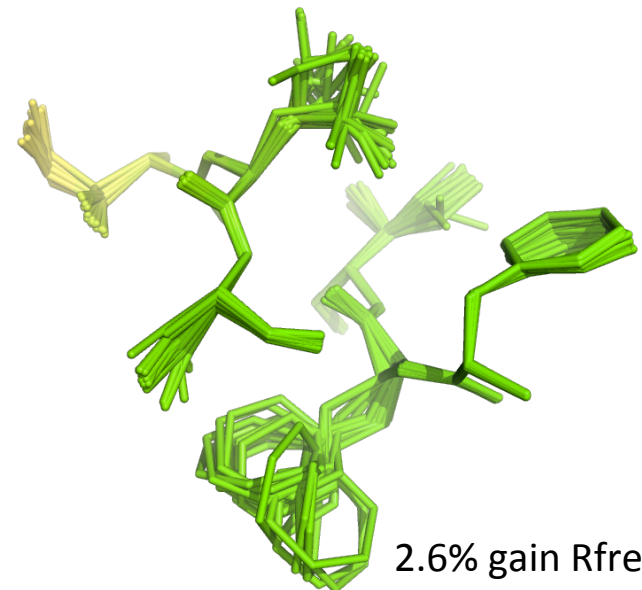
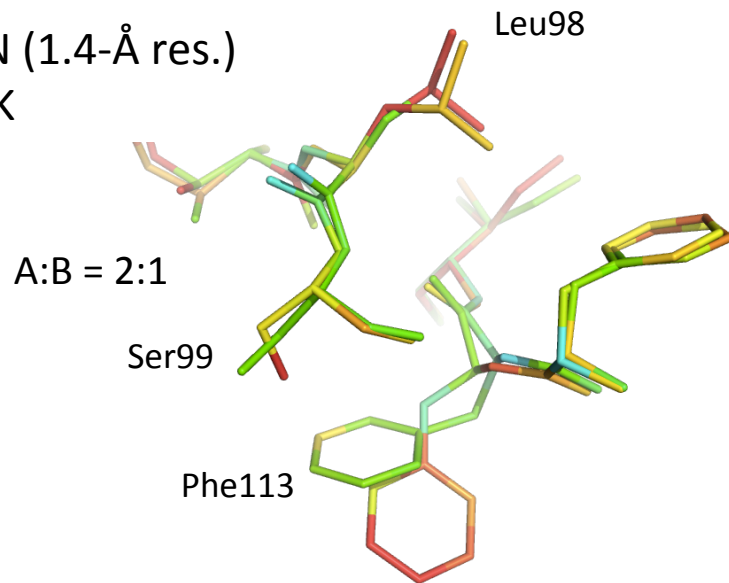
# Multi-conformers in active site

3K0M (1.3-Å res.)  
100 K



Ensemble refinement:  
1.3% gain Rfree

3K0N (1.4-Å res.)  
288 K

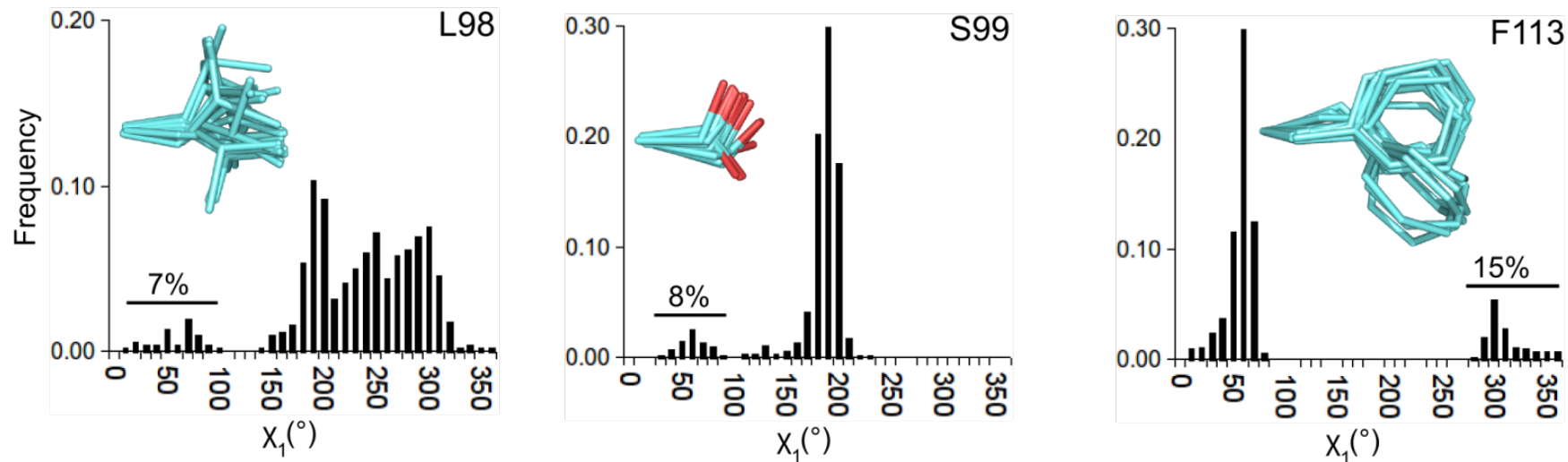


2.6% gain Rfree

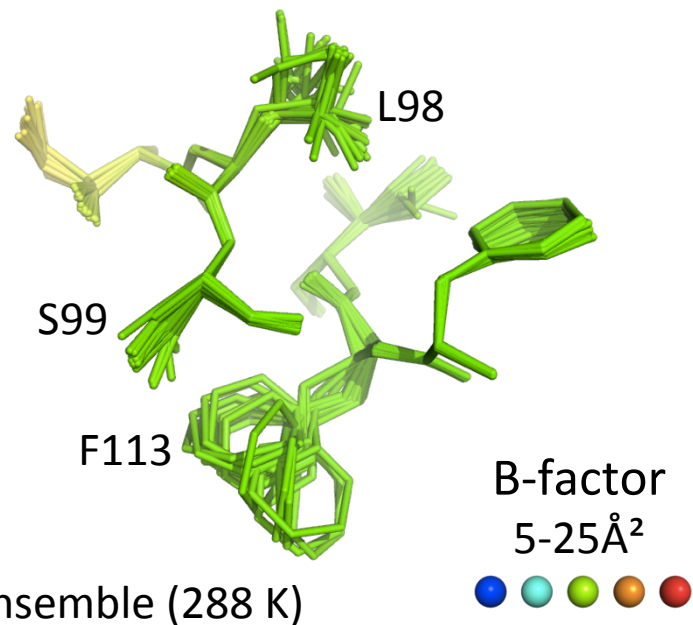
B-factor  
5-25Å<sup>2</sup>



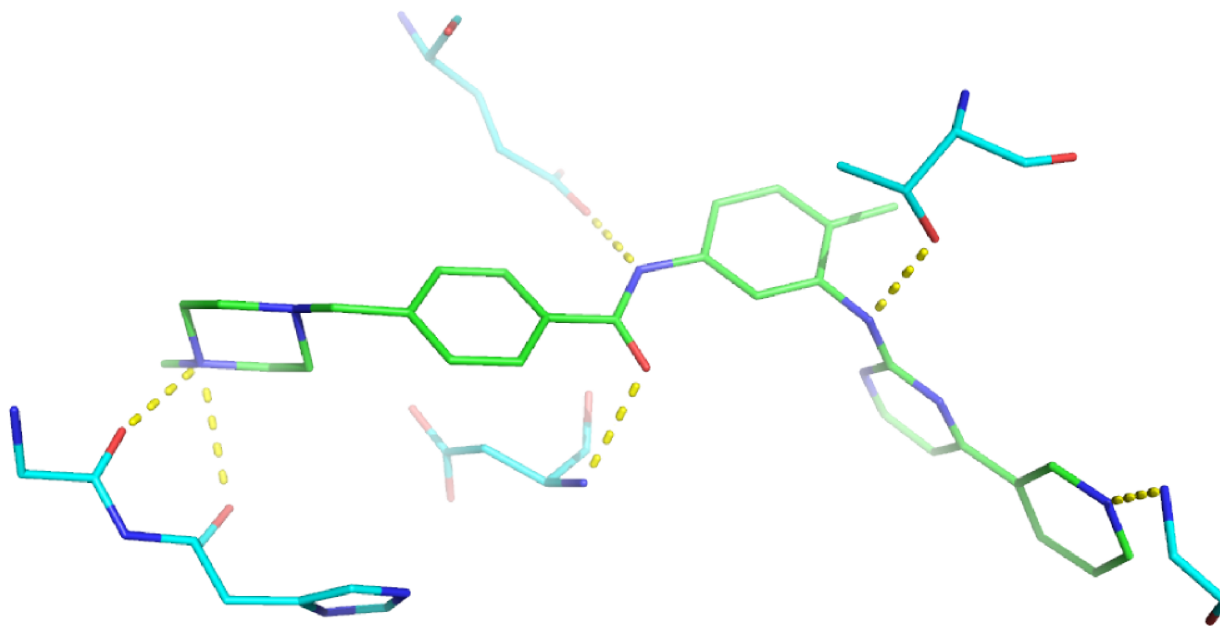
# Occupancies agree with NMR data



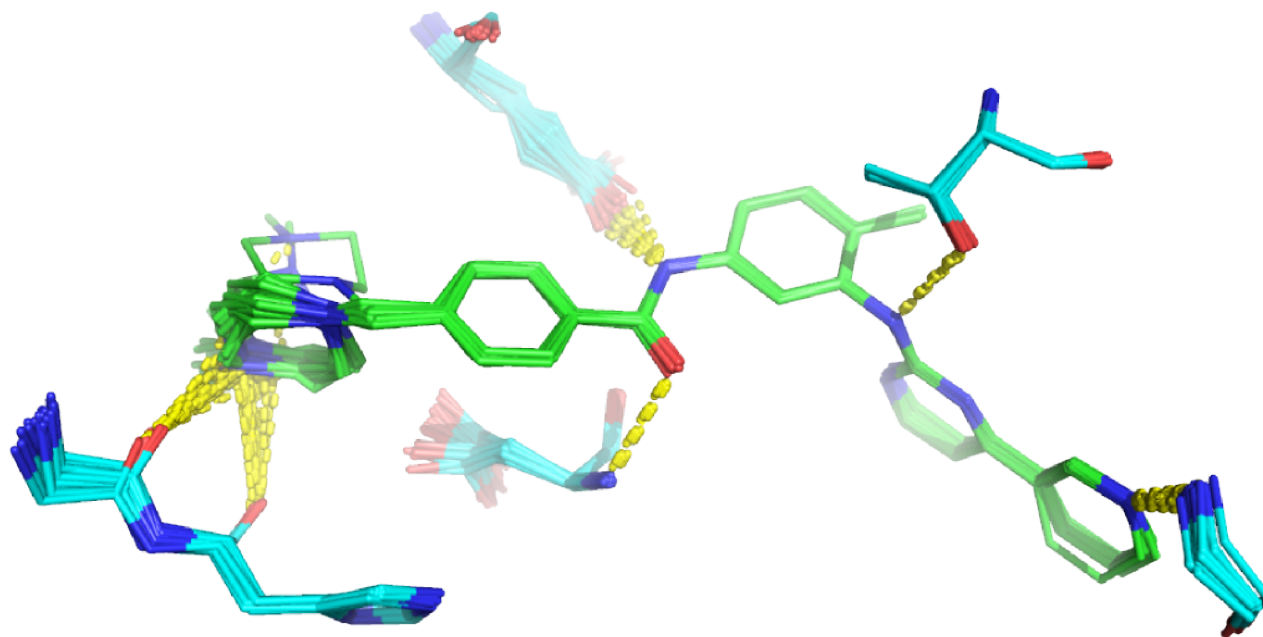
- NMR relaxation dispersion (283 K)  
L98, S99, F113  
Minor population  $\sim 10\%$



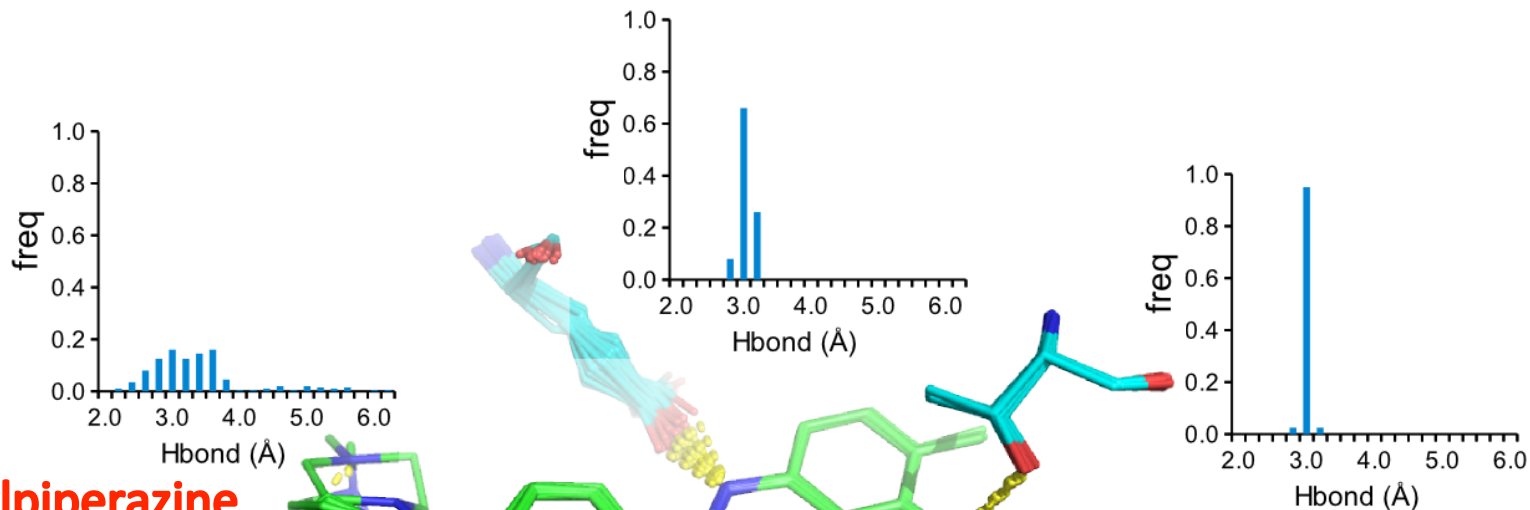
# Imatinib-ABL Tyrosine Kinase (1IEP)



# Imatinib-ABL Tyrosine Kinase (1IEP)



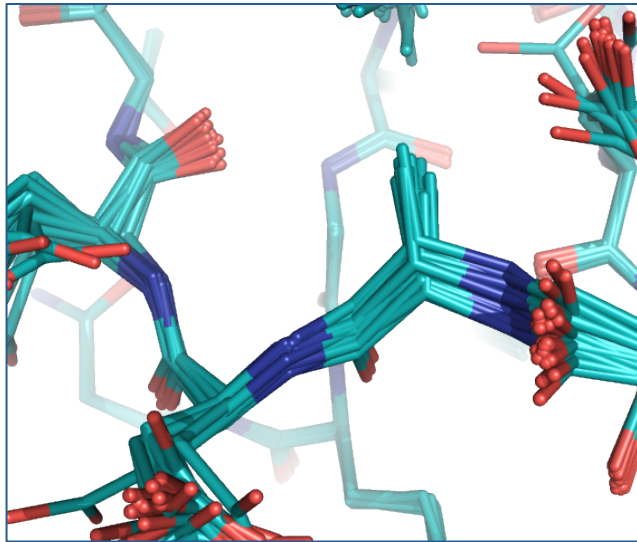
# Imatinib-ABL Tyrosine Kinase (1IEP)



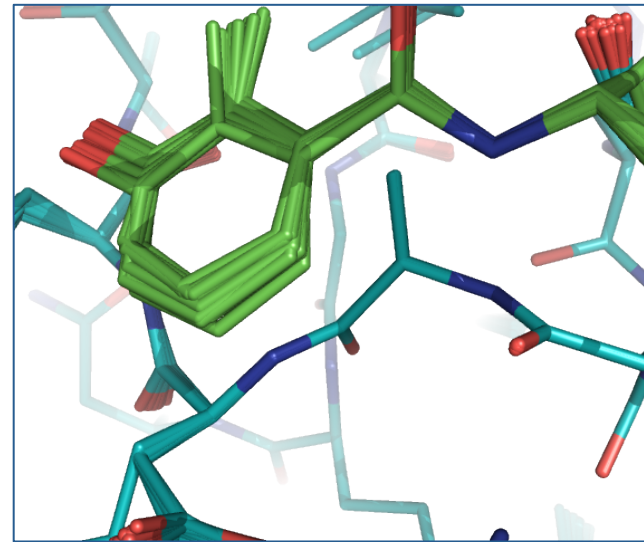
**N-methylpiperazine**

**Thr315 'Gatekeeper'**

# Inhibitor binding to HIV protease



apo-protease (2PC0)



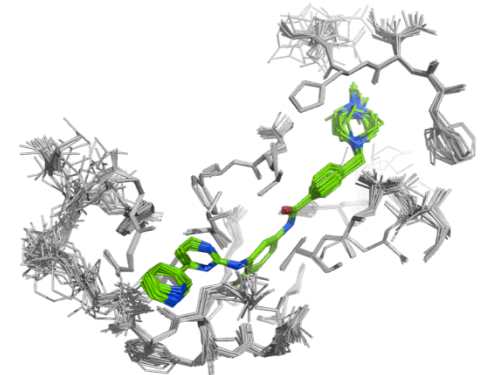
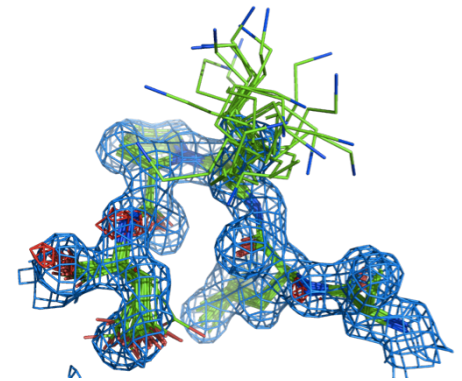
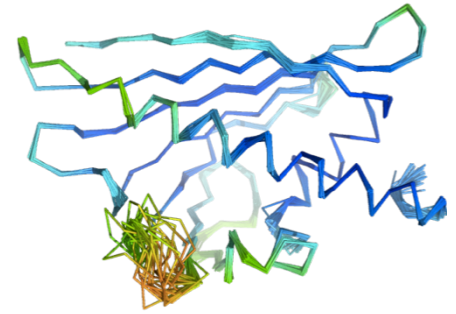
complex with JE-2147 (1KZK)

Comparative analysis of atomic flexibility:

- Isomorphous crystals eliminate differences due to Xtal contacts
- Non-isomorphous crystals allow evaluation of Xtal contact effects

# Conclusions

- Global disorder modeled by TLS and local disorder by MD
- Ensemble refinement improves Rfree and electron density maps
- Suitable for a broad resolution range (1Å – 3Å)
- NCS copies show very similar fluctuations
- Clear representation of local disorder / uncertainty
- Distribution of atom positions allows further structural analysis
- Resolve the finer details of protein structure(s)





# Acknowledgements

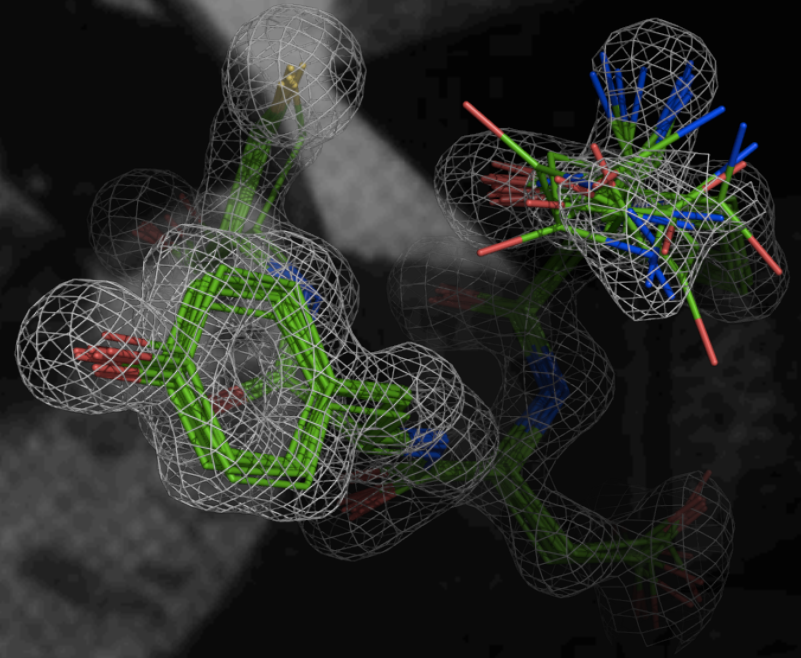
**Piet Gros**

**Gros Laboratory**

**Pavel Afonine, Paul Adams**

**PHENIX Developers**

**Funding: Utrecht University / NWO**



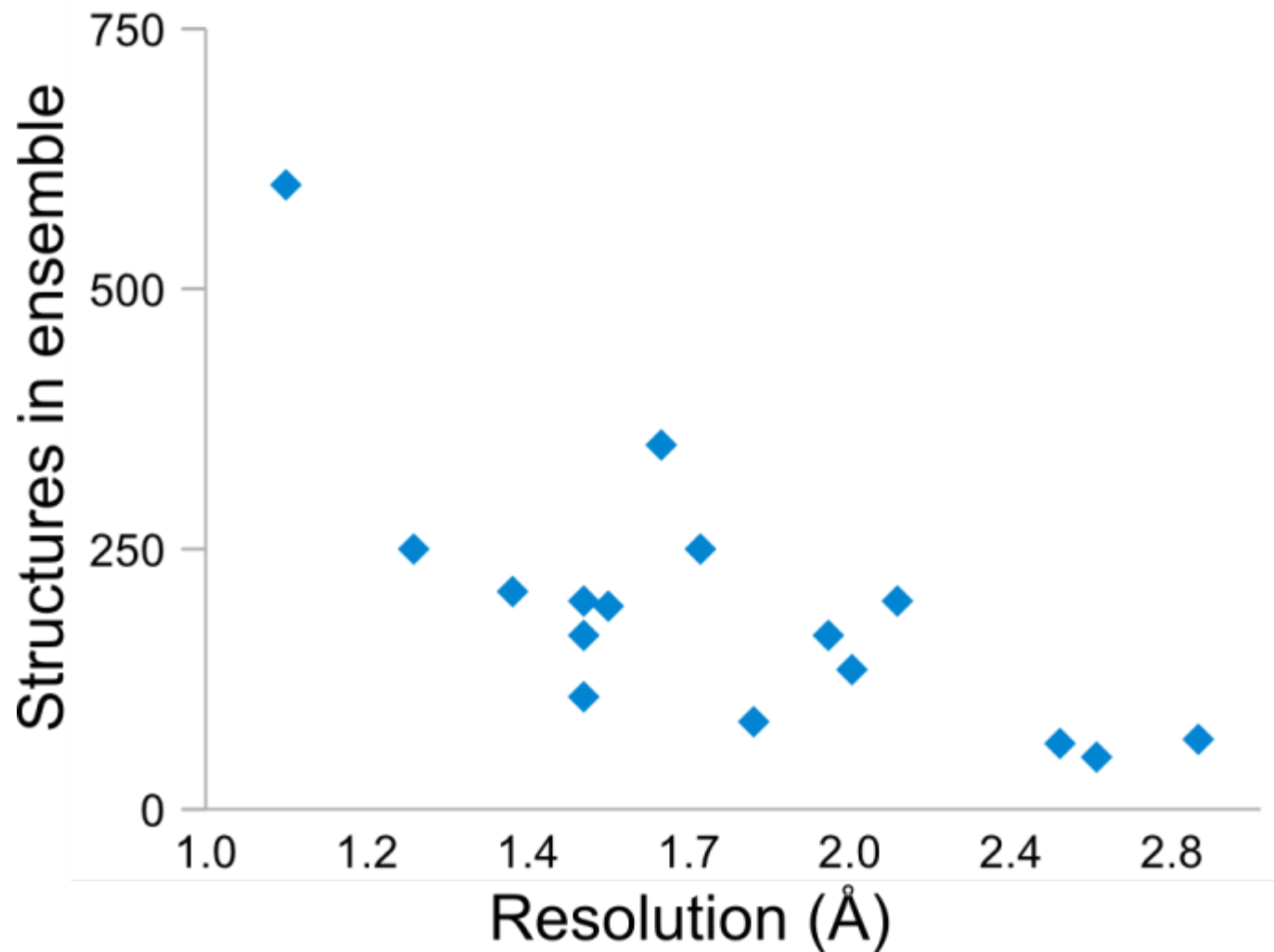
**Universiteit Utrecht**



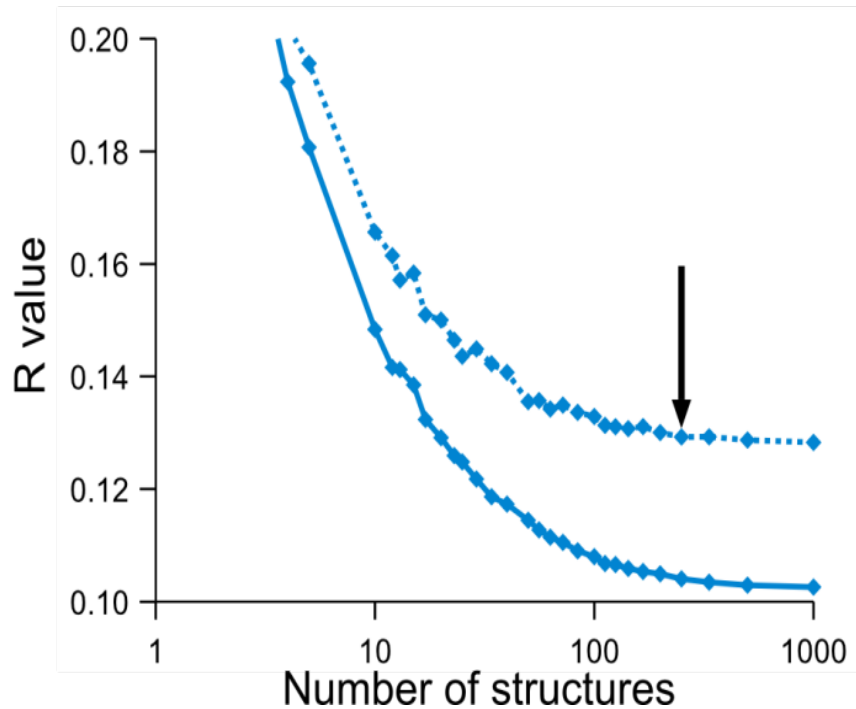
**NWO**

**Phenix**

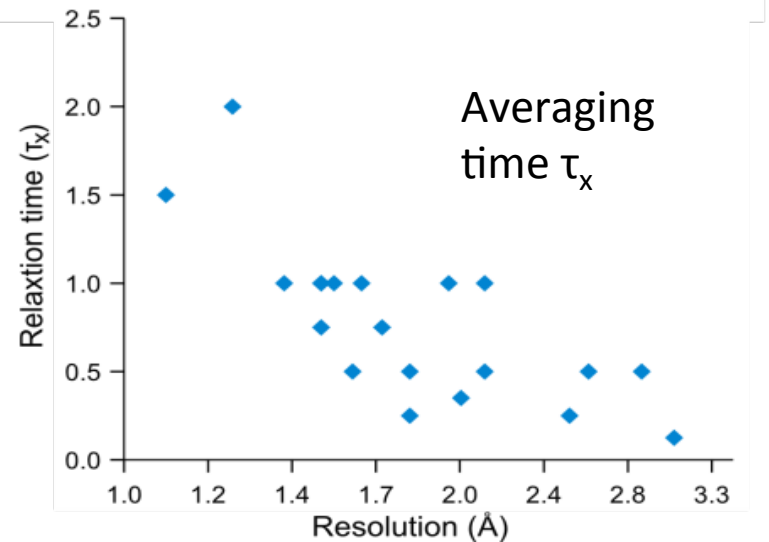
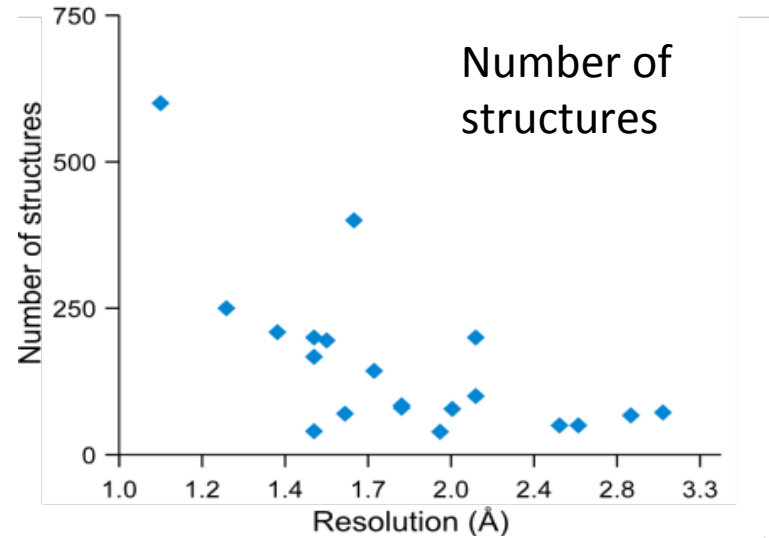
## *Number structures in ensemble vs resolution*



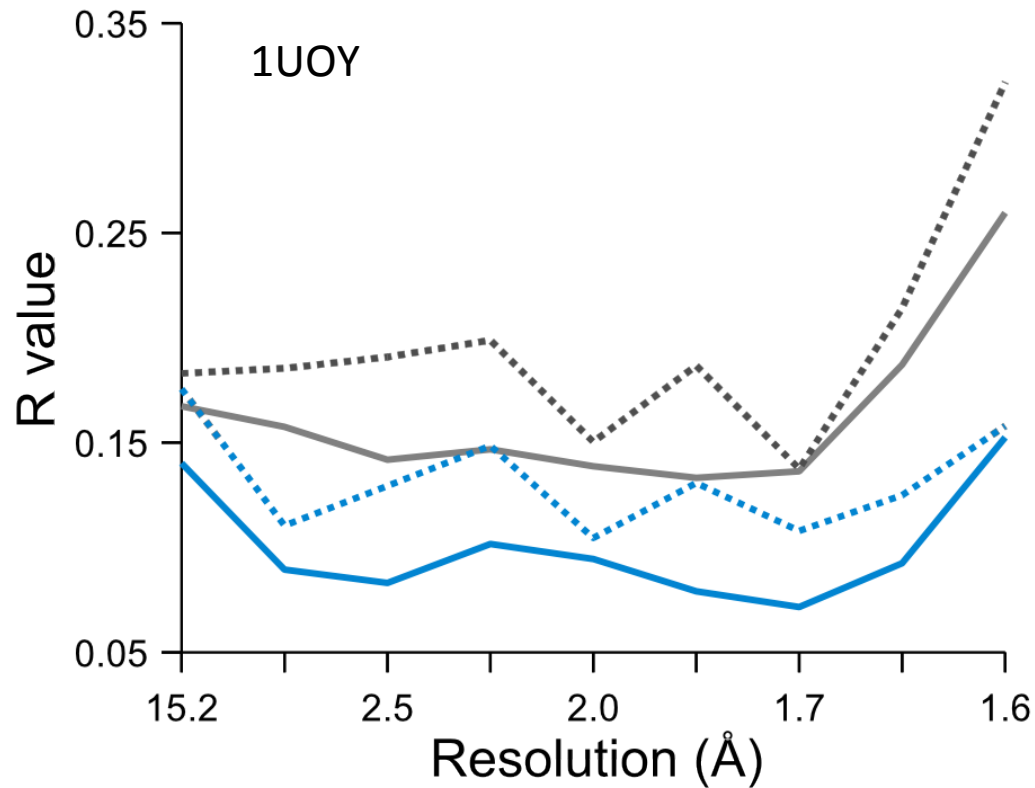
# Number of structures in the ensemble



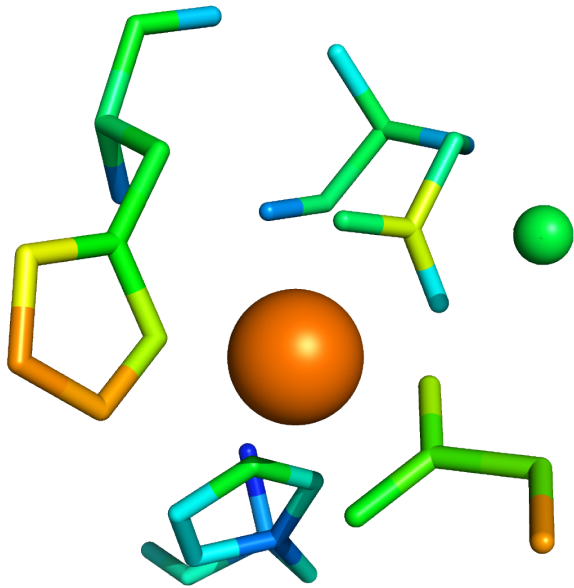
Structures taken equidistant in time to reproduce the Rwork within 0.1%



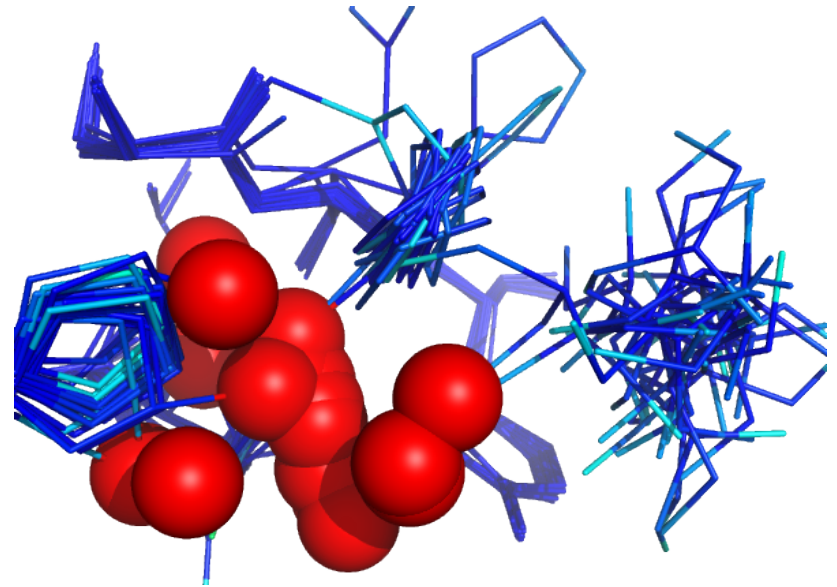
# R-factor by resolution shell



# Partial ligand/ion binding



B-factor of  $\text{Cd}^{2+}$ -ion more than twice its surrounding

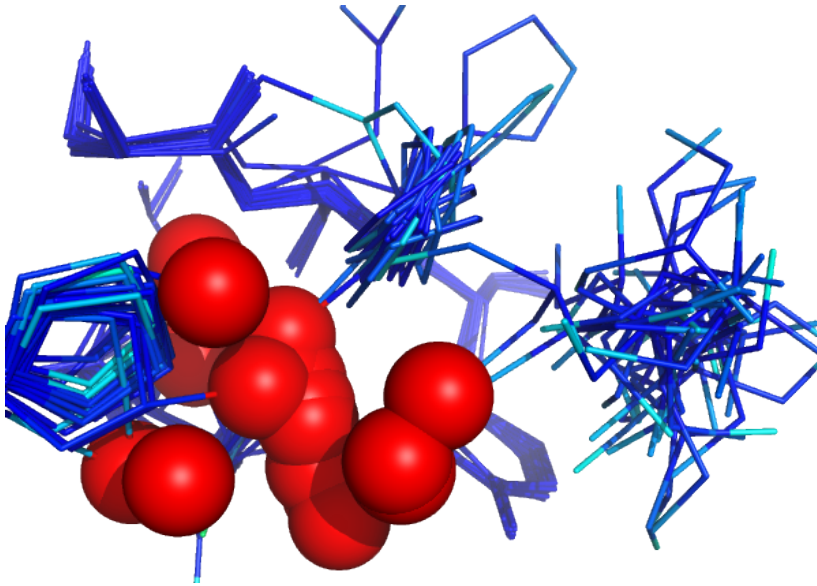


Simulation run at  $Q=1$  for  $\text{Cd}^{2+}$ -ion (atoms coloured by kinetic energy)

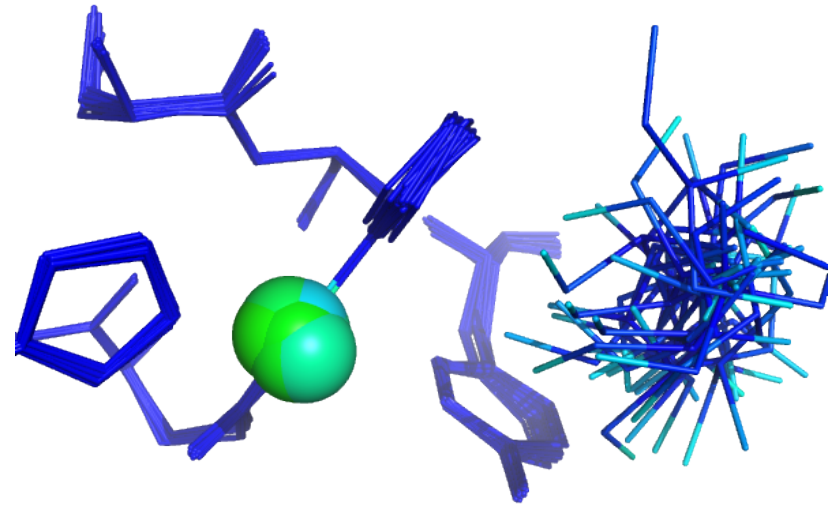
- High B-factor ligand/ion may indicate partial occupancy

# Partial ligand/ion binding

$\text{Cd}^{2+}$ :  $Q=1$



$\text{Cd}^{2+}$ :  $Q=0.7$

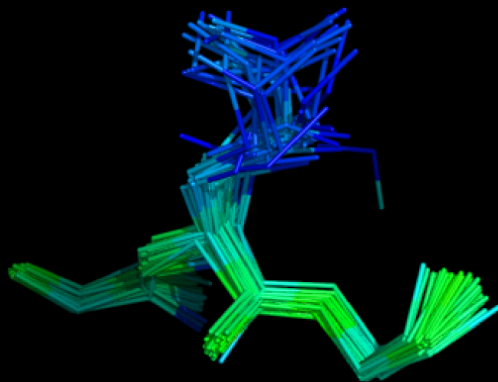


- Heat map (kinetic energy) is validation tool for Ensemble Refinement

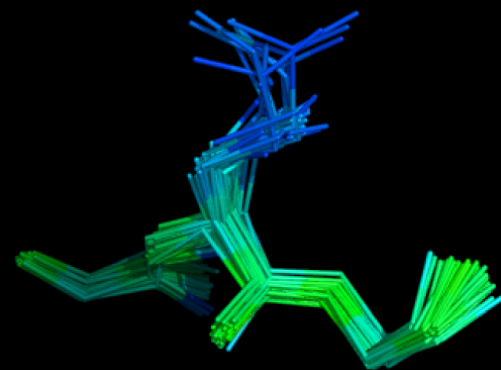
# Ensemble



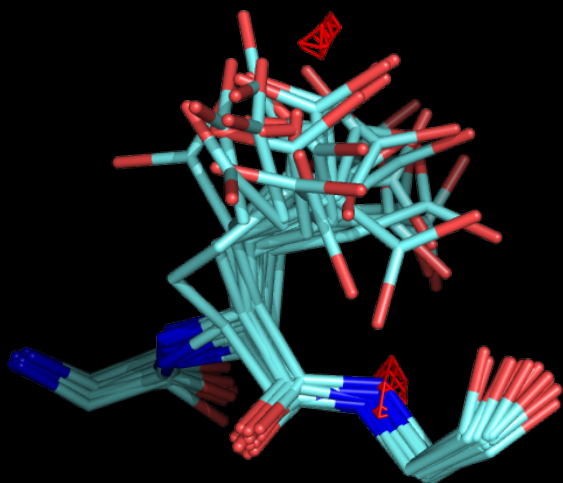
0–10  $\sigma$



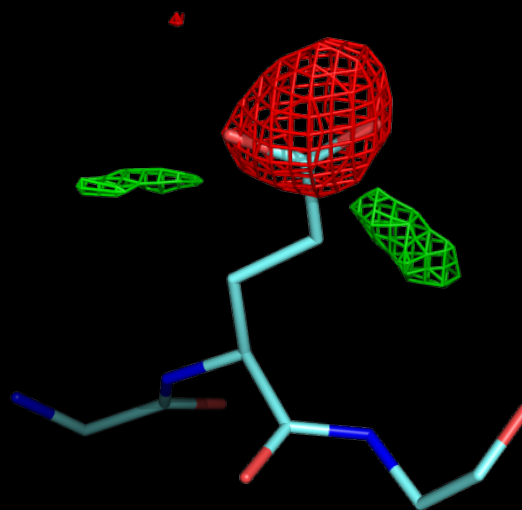
0.5–10  $\sigma$



1–10  $\sigma$



# Ensemble

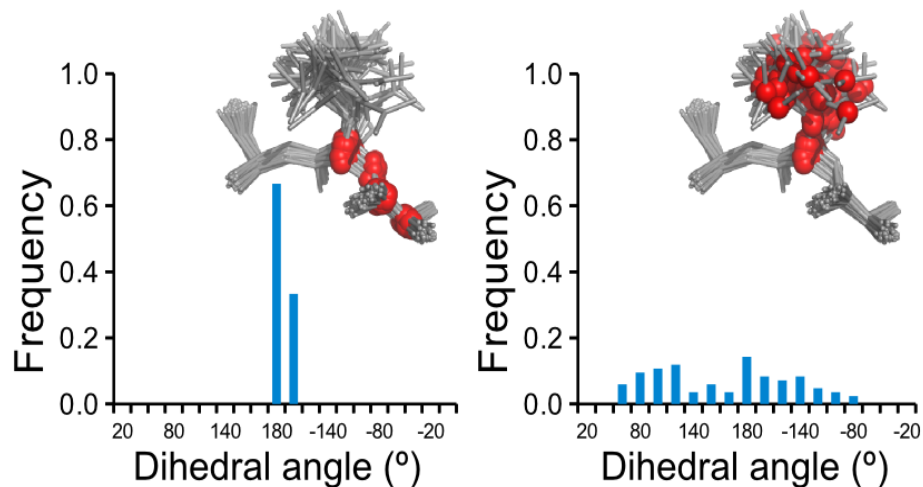


# Single-structure

# Geometric validation

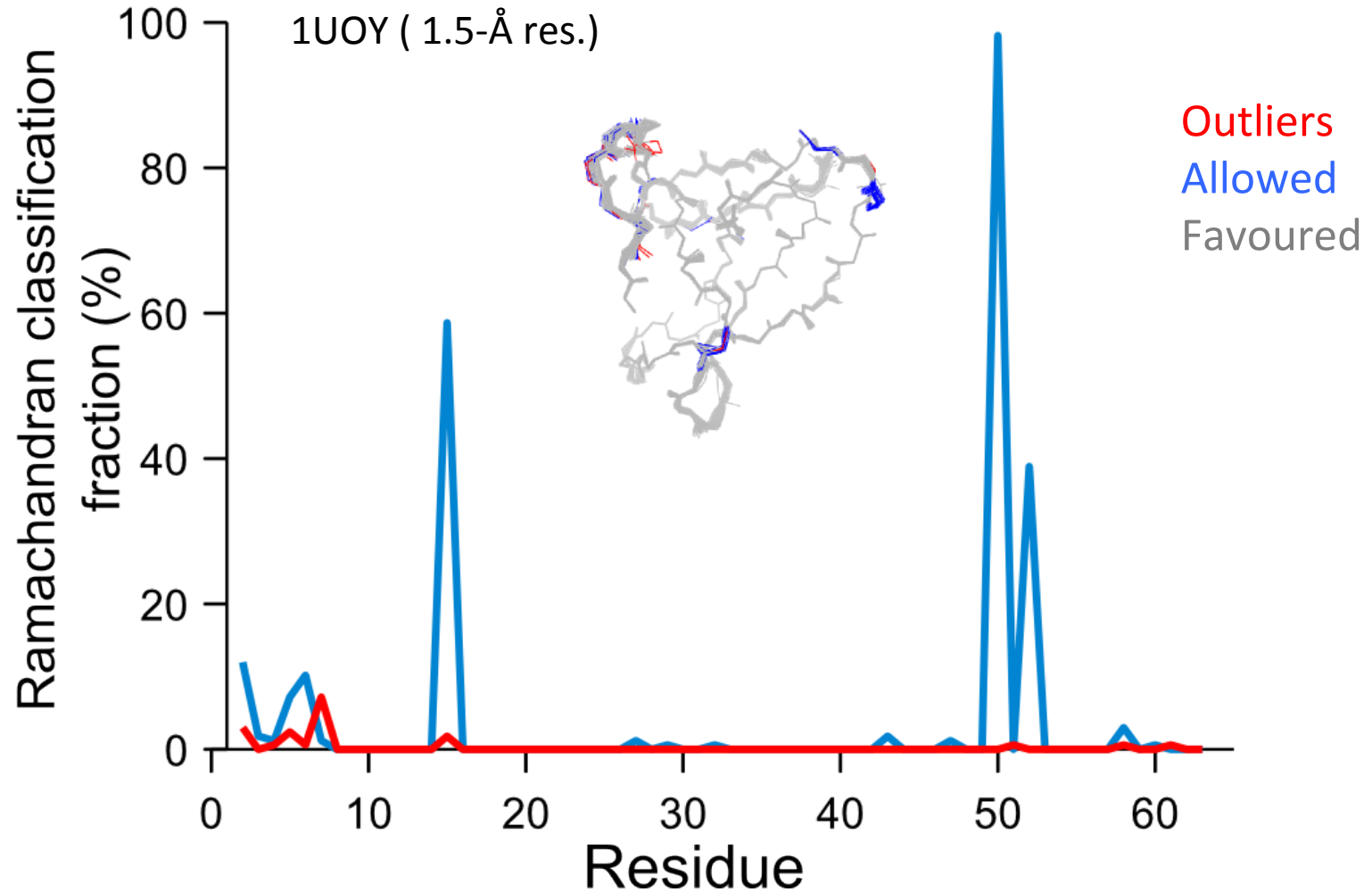
	Bonds	Angles	Dih. angles
Single structure	0.010 Å	1.23°	15.1°
Ens: $\sqrt{\langle(x_{ideal}-\langle x_{model} \rangle)^2\rangle}$	-0.002	-0.28	-6.5
Ens: $\sqrt{\langle(x_{ideal}-x_{model})^2\rangle}$	+0.002	+0.30	+4.0

(statistics for all 20 cases)



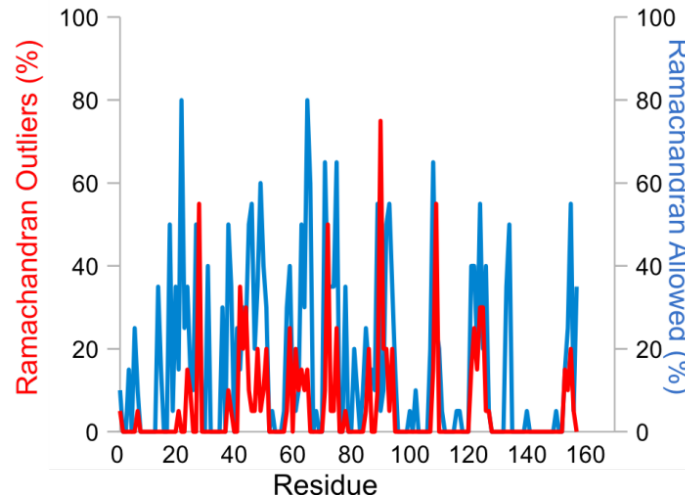
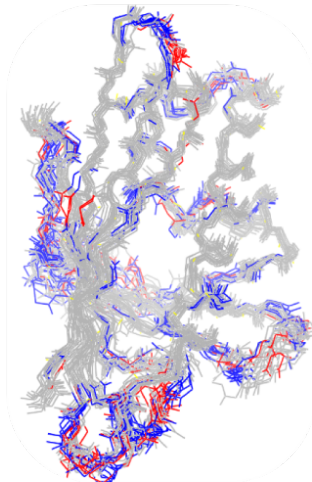
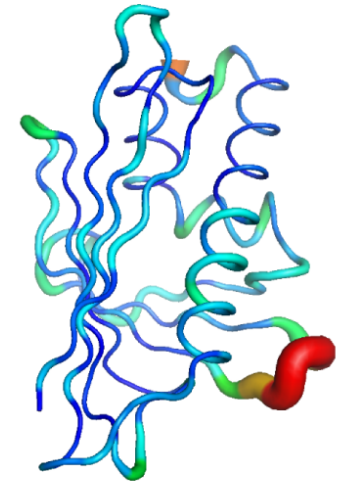
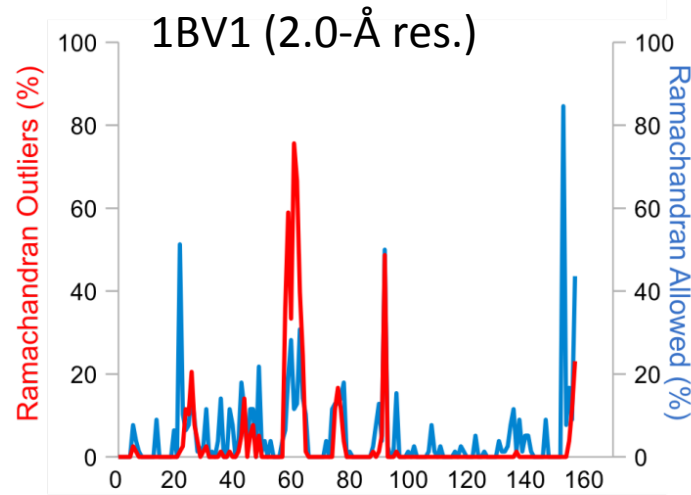
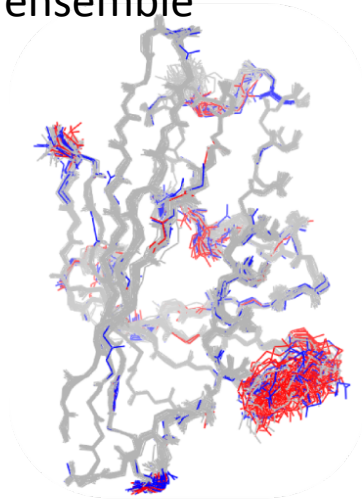


# Ramachandran analysis



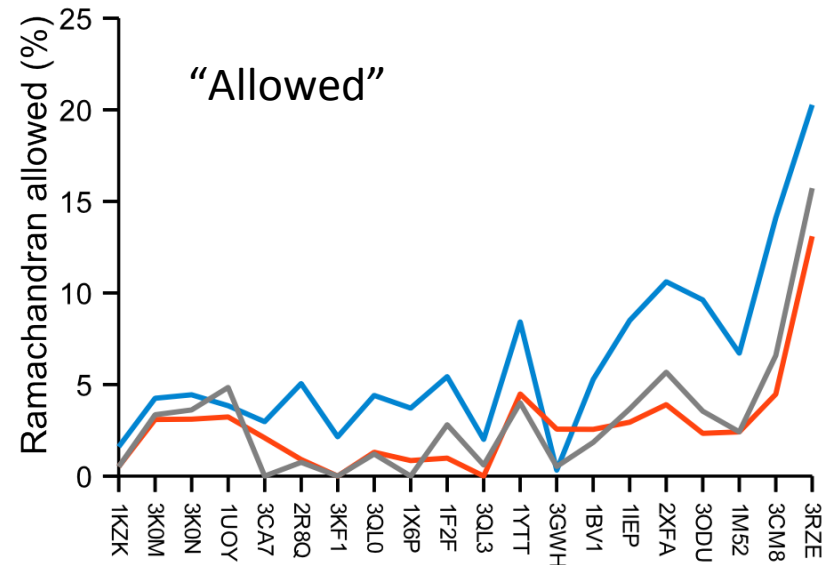
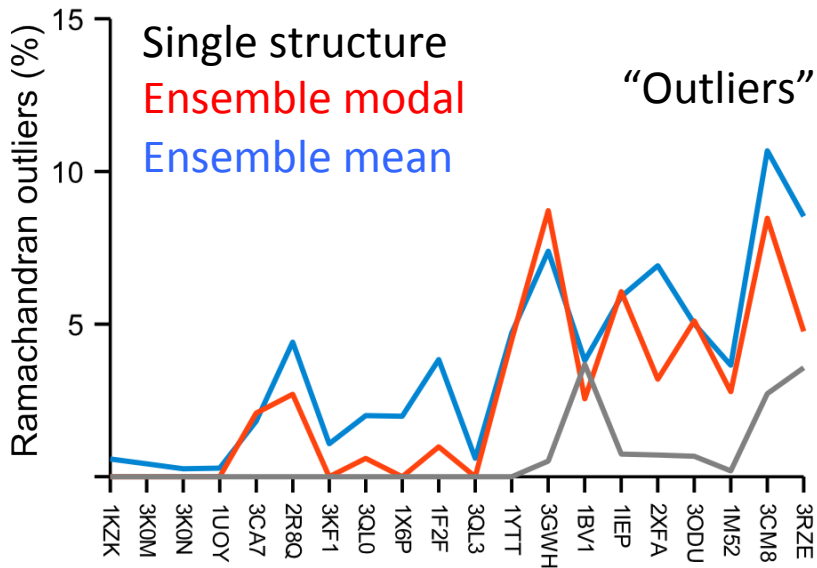
# Outliers occur more in flexible regions

X-ray ensemble



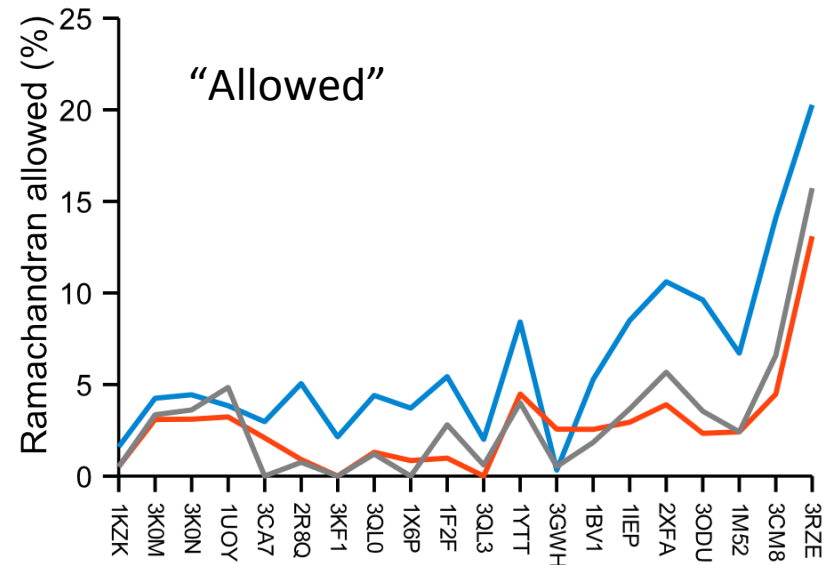
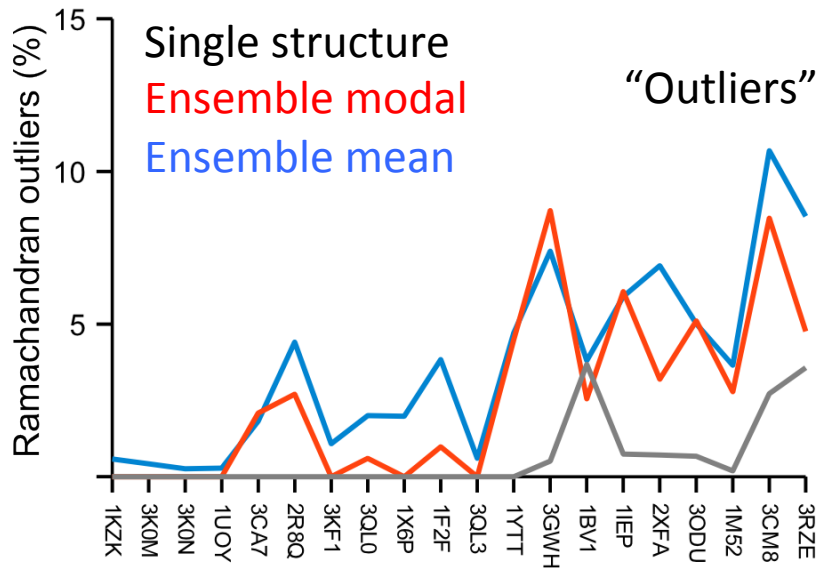
NMR ensemble

# Ramachandran deviations by resolution



- More outliers are observed at lower resolution
- Geometric quality correlates with Rfree
- “Best” run selected by Rfree

# Ramachandran deviations by resolution



- More outliers are observed at lower resolution
- Geometric quality correlates with Rfree
- “Best” run selected by Rfree