



# Acknowledgements

James Fraser Nevan Krogan

Paul Adams John Tainer Greg Hura Nick Sauter

Pavel Afonine Tom Terwilliger Nigel Moriarty, Dorothee Liebschner

Aina Cohen

**UCSF LBNL SLAC**

**ALS 8.3.1: TomAlberTron**

**NIH NIGMS R01 GM124149 (Holton)**

**NIH NIAID P50 AI150476 (Krogan)**

**NIH NIGMS P30 GM124169 (Adams)**

**DOE-BER IDAT (Hura)**

**NIH NIGMS P30 GM133894 (Hodgson)**

# Why is a beamline guy doing refinement things?



Because: the future of crystallography  
has alternate conformations

- ✓ Single-conf models       Dynamics?
- ✓ poor resolution       Mechanism?
- ✓ Low signal/noise       Small signal?
- ✓ Phase Problem       Overlapping states?

# Map Noise: via Parseval's Theorem

$$\frac{\rho}{\sigma(\rho)} \approx \frac{e^-}{7R}$$

$e^-$

number of electrons in peak

7

number of electrons in typical atom

$R$

$R_{\text{iso}}$  or  $R_{\text{free}}$

$\rho$

electron density peak

$\sigma$

rms deviation expected

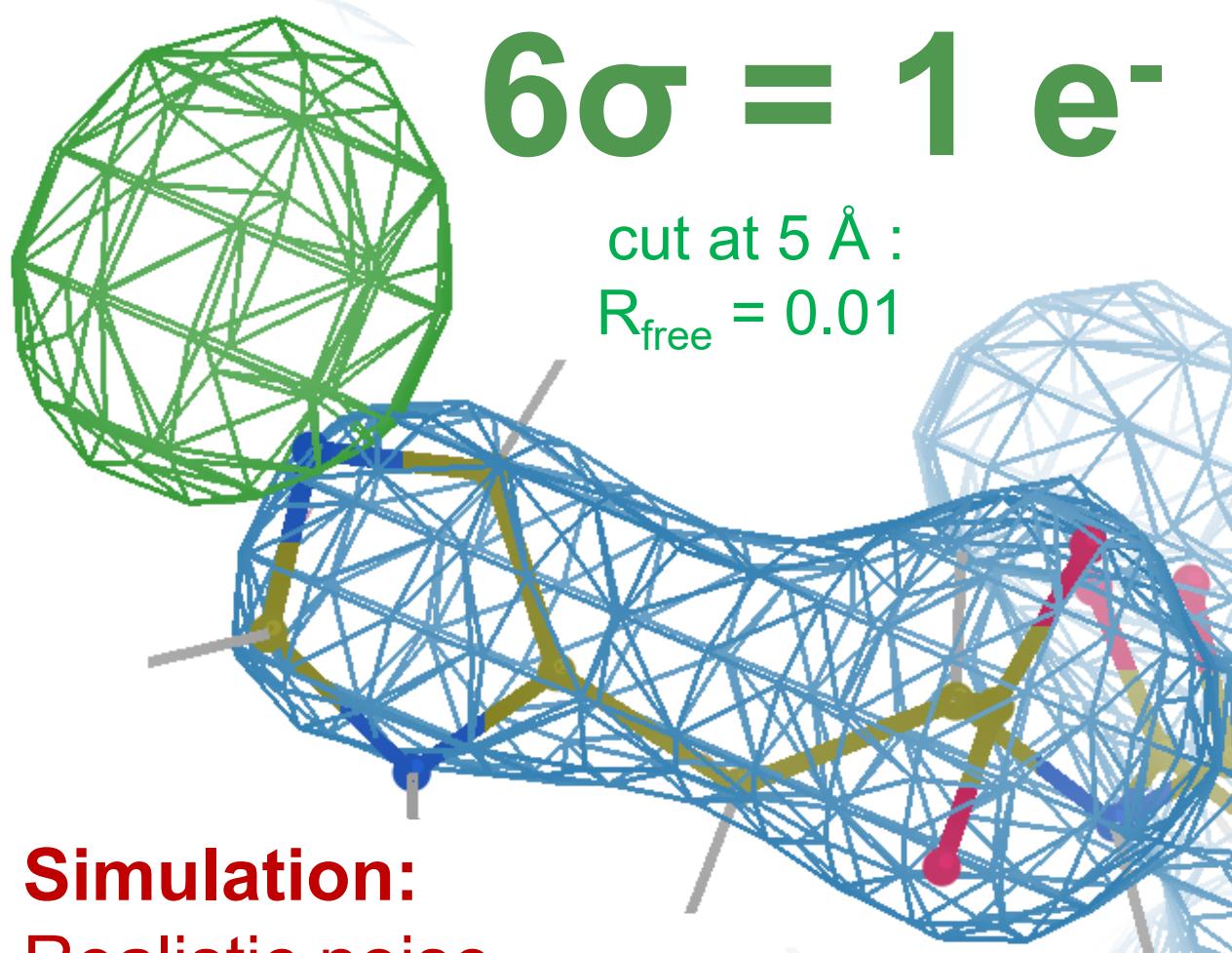
# Hydrogen: visible at 3.5 Å resolution!

$R_{\text{work}} = 0.04$

$R_{\text{free}} = 0.05$

$6\sigma = 1 \text{ e}^-$

cut at 5 Å :  
 $R_{\text{free}} = 0.01$



**Simulation:**  
Realistic noise  
“perfect” model

# Real Data

R<sub>work</sub>: 0.0823  
R<sub>free</sub>: 0.1106  
bond: 0.014 Å  
angle: 1.92°

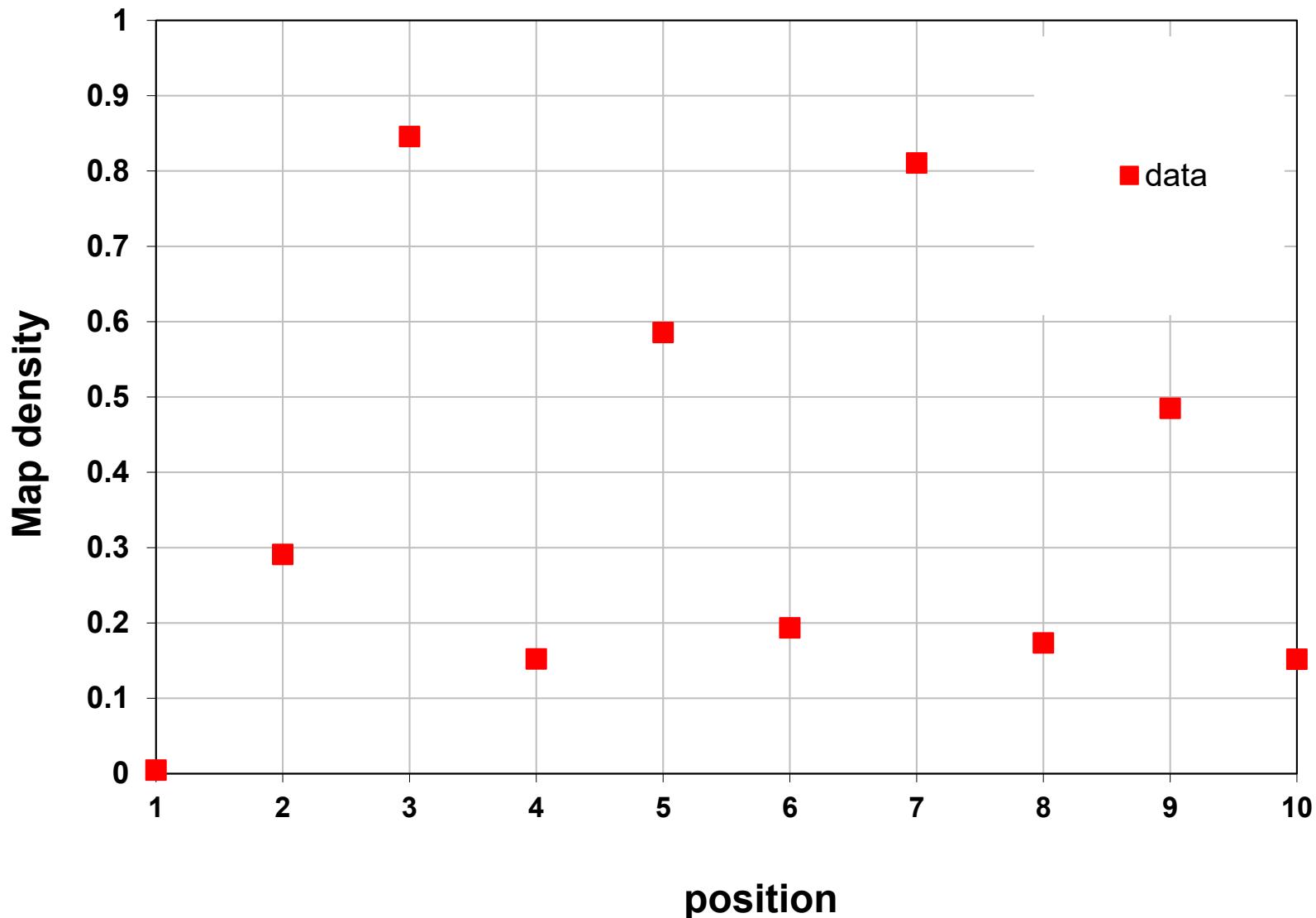
48-copy ensemble  
3 params per observation

MoPro: 1.15  
Clash: 0.49  
C $\beta$ dev: 20  
Rota: 1.7%  
Rama: 0.03%  
reso: 1.0 Å

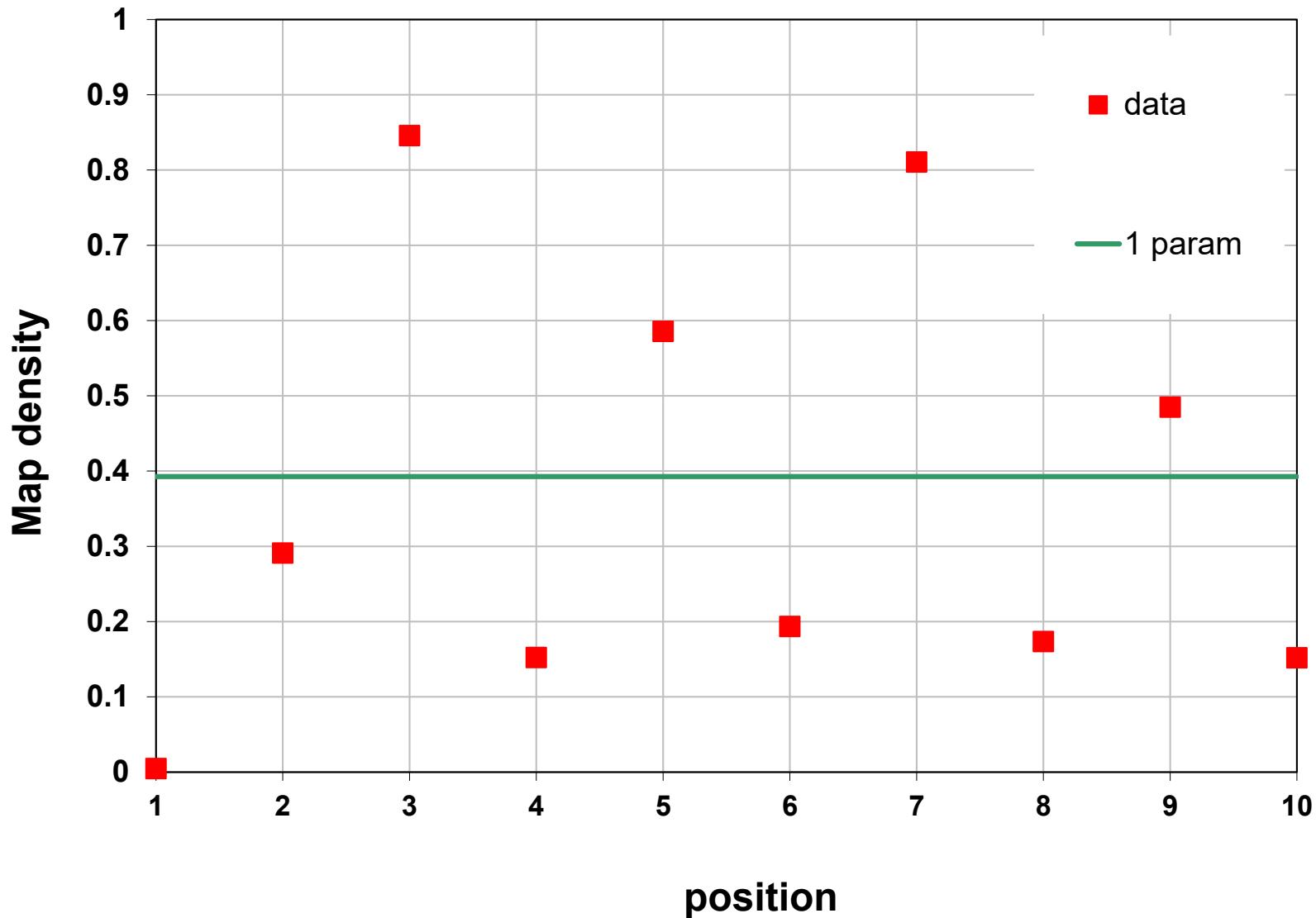
6 $\sigma$

1.2 $\sigma$

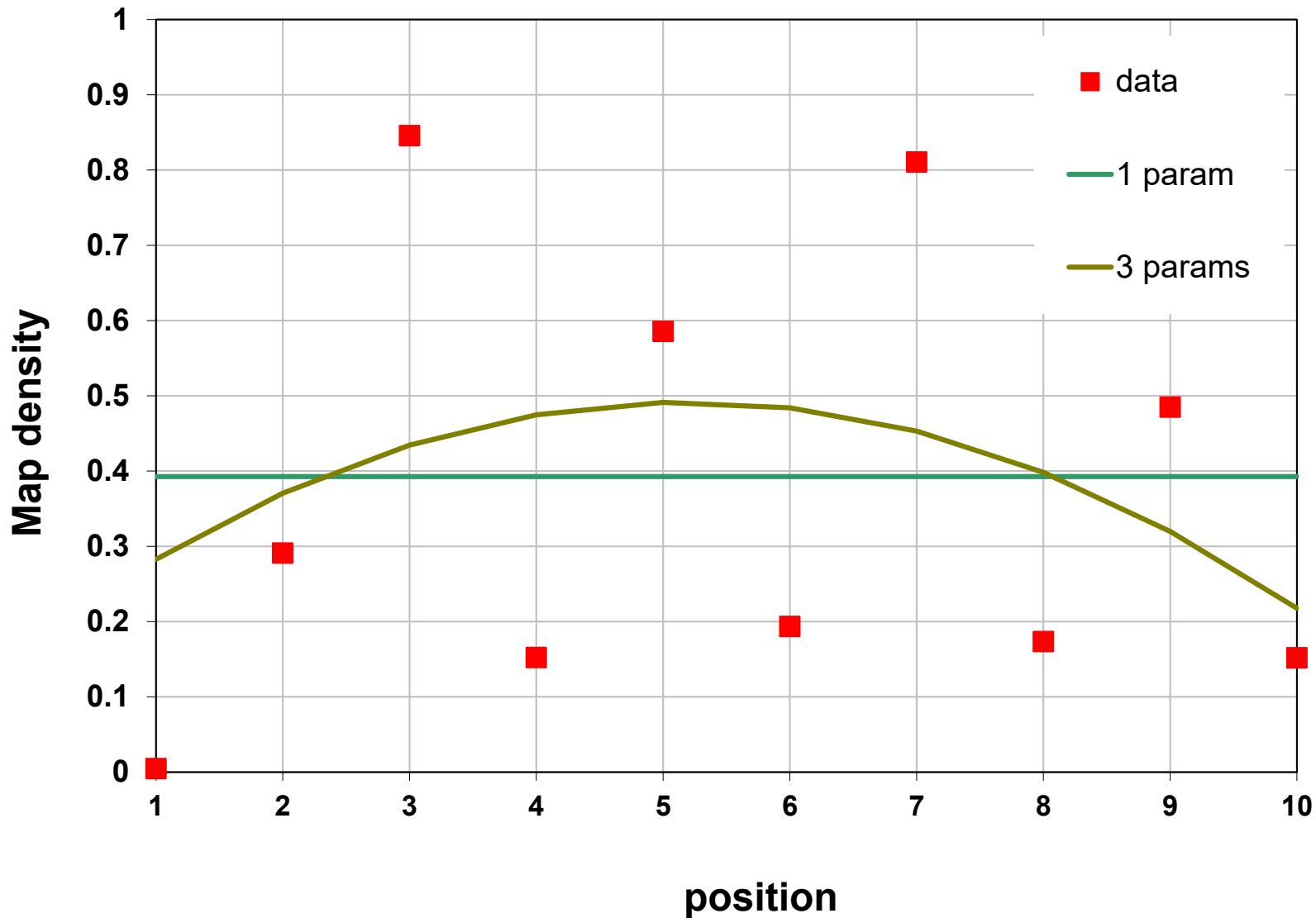
# Over-Fitting data



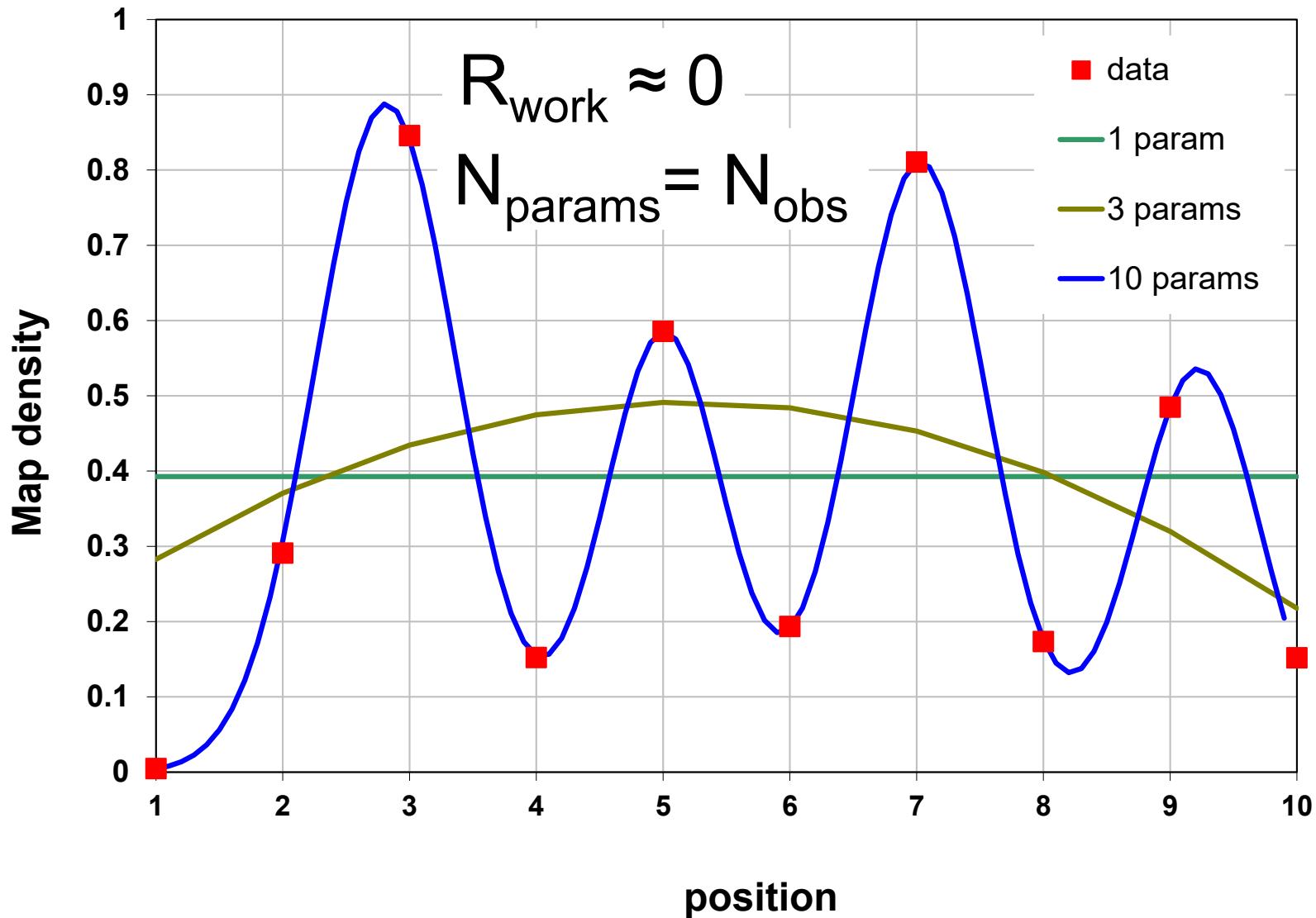
# Over-Fitting data



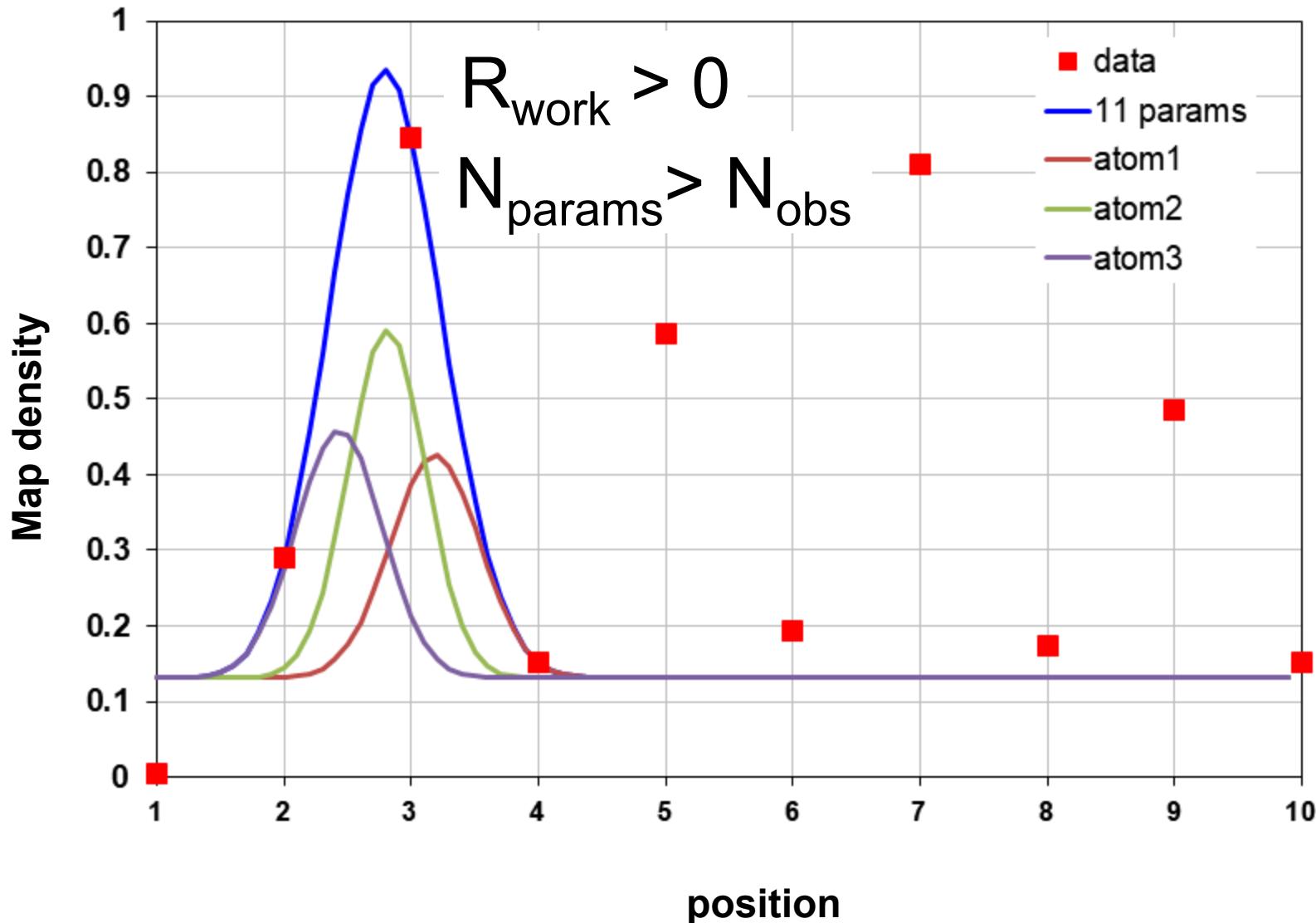
# Over-Fitting data



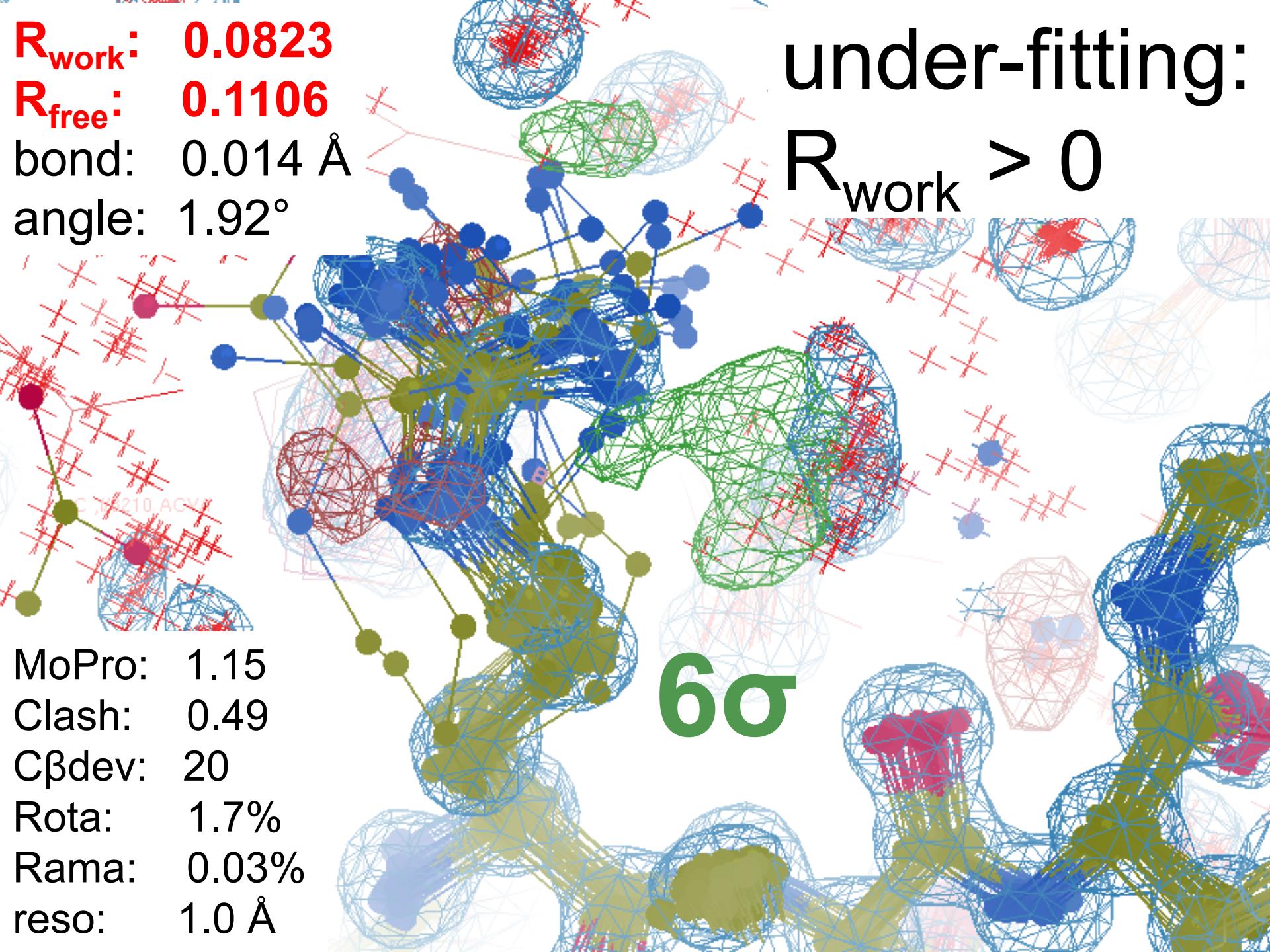
# Over-Fitting data



# Under-Fitting data



$R_{work}$ : 0.0823  
 $R_{free}$ : 0.1106  
bond: 0.014 Å  
angle: 1.92°

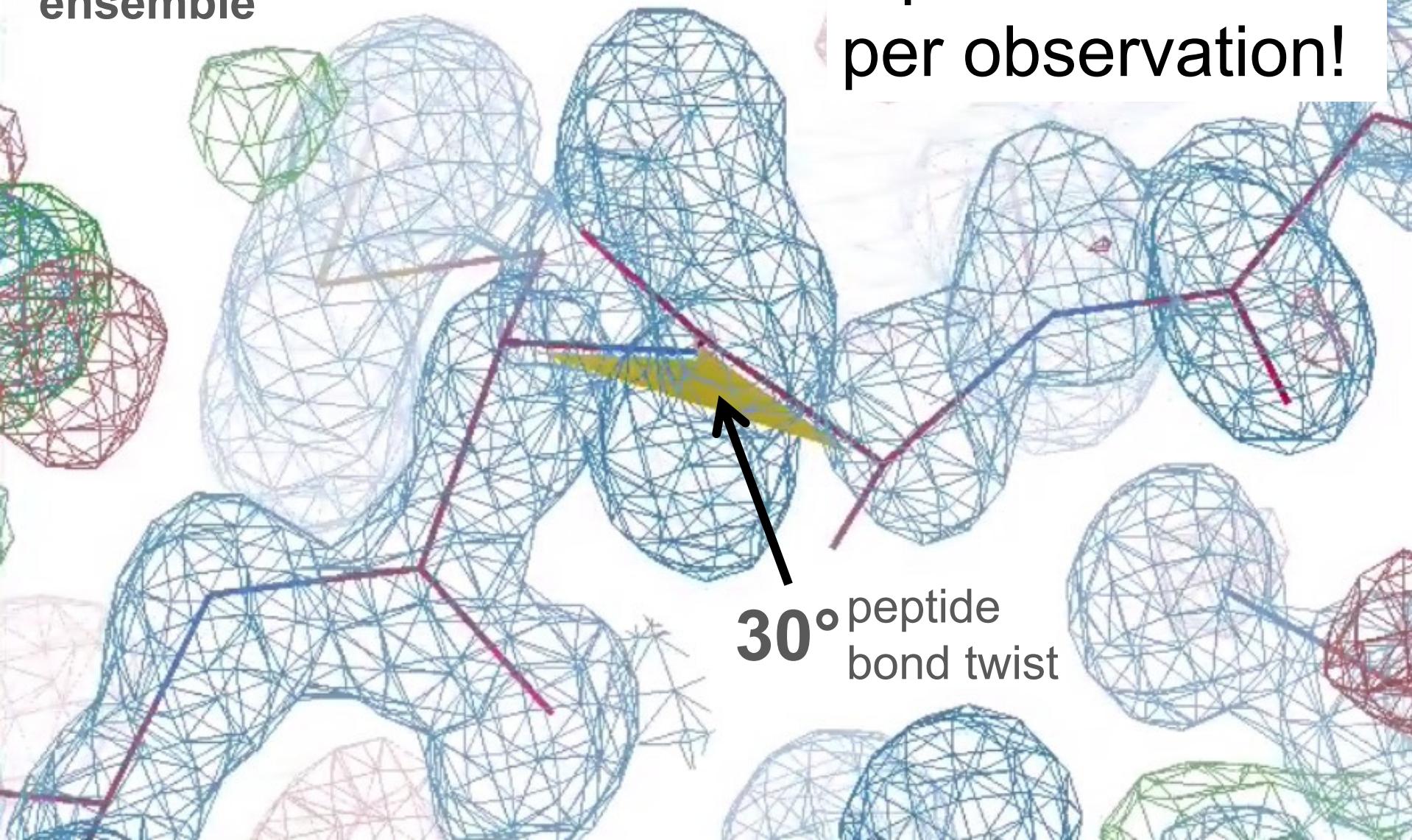


under-fitting:  
 $R_{work} > 0$

# What is holding it back?

1 member of 48-copy ensemble

3 parameters per observation!



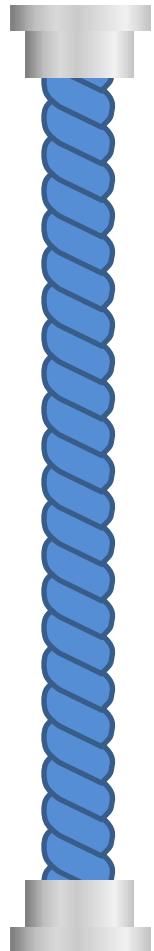
# Why does ensemble refinement **suck?**

(3 parameters/observation and  $R_{\text{work}}$  still high!)

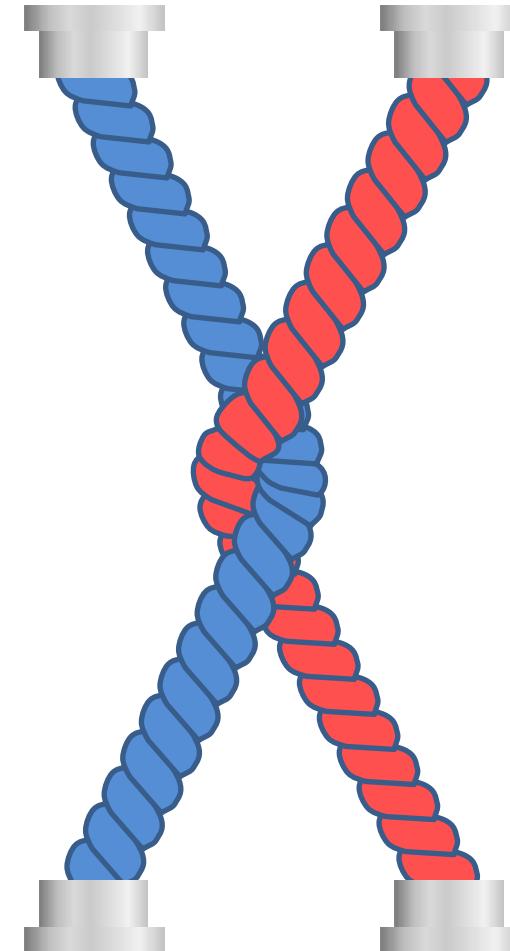
The chains are “tangled”

# What do you mean, “tangled”?

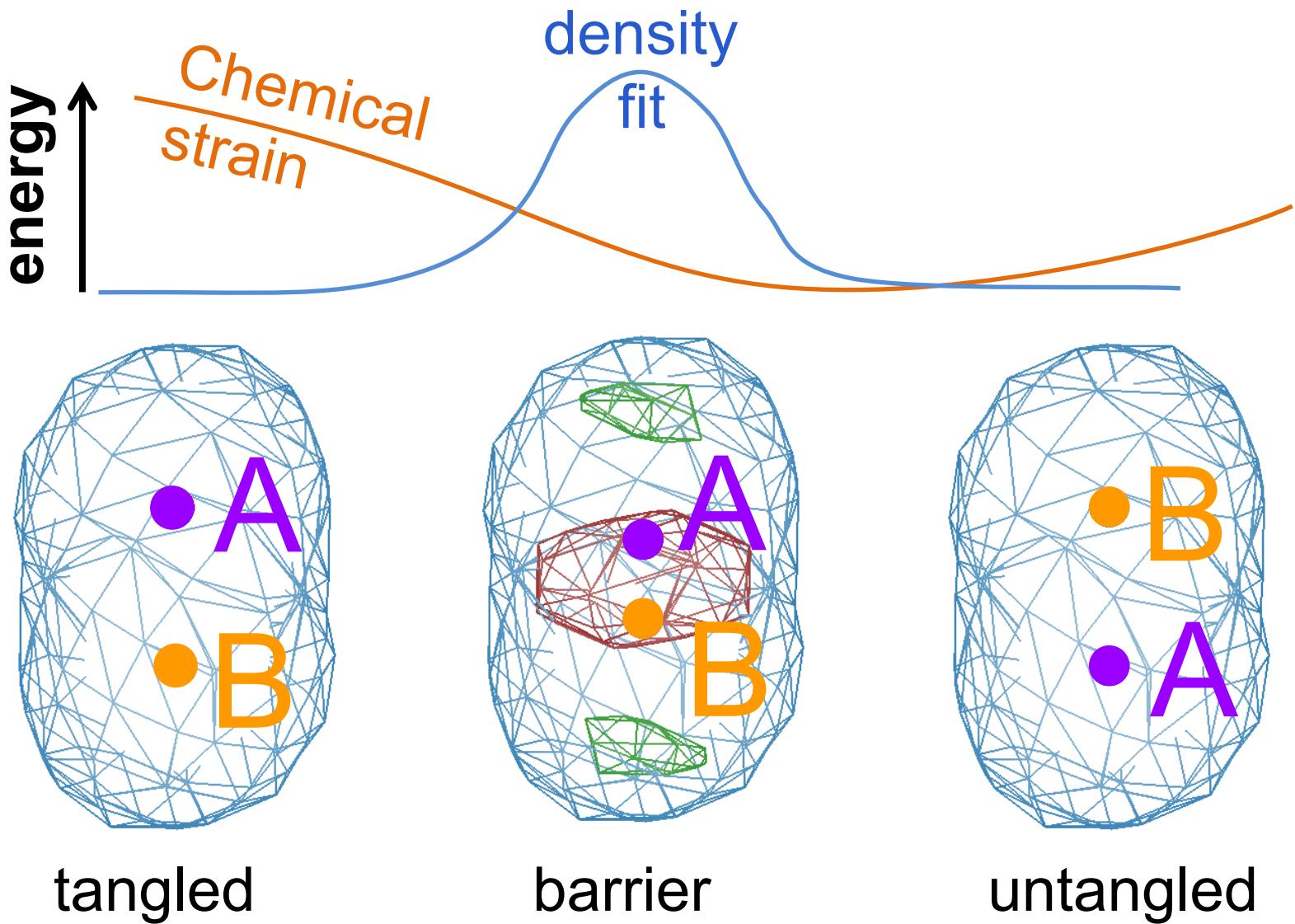
Ground Truth



Local Minimum



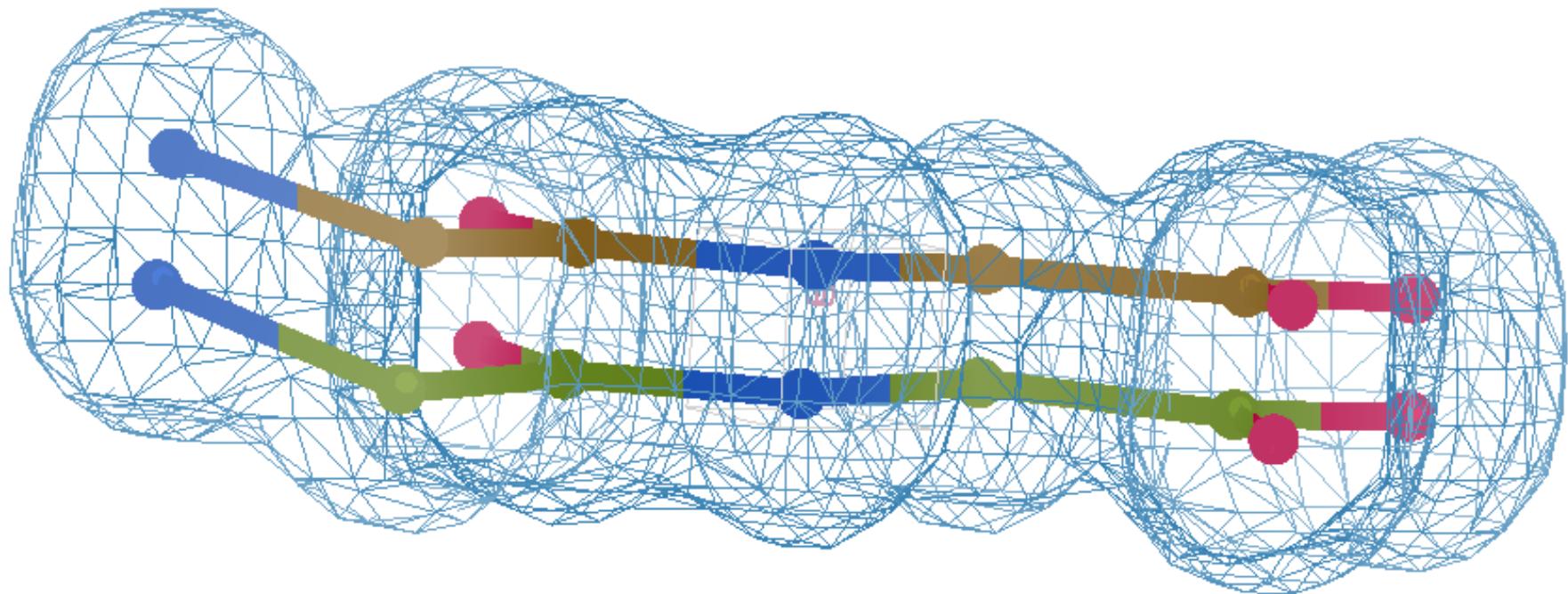
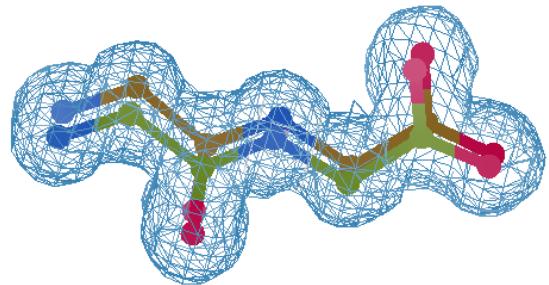
# It takes two to tangle



# Refine from right answer

$R = 0.8\%$

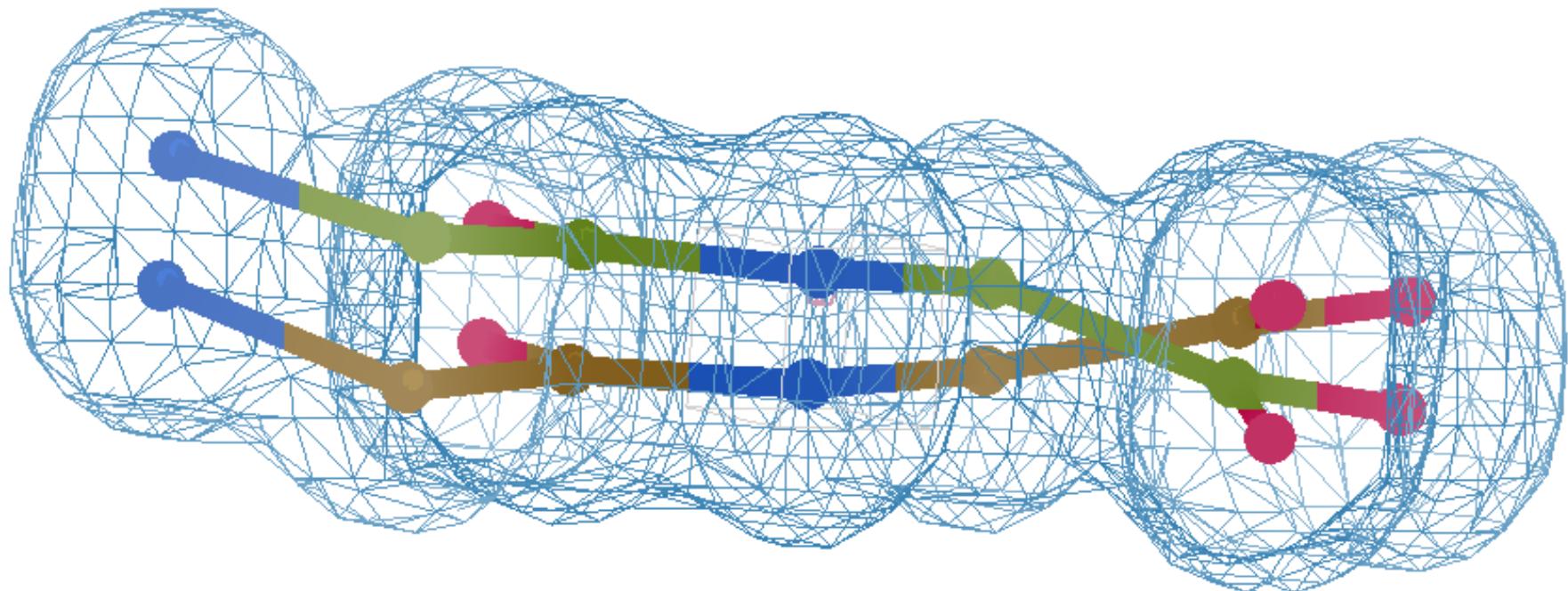
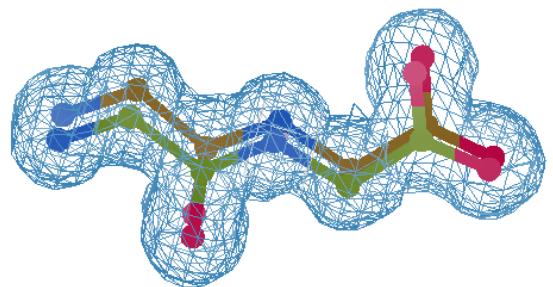
rms bond = 0.0002



# Refine from wrong starting point

$R = 7.07\%$

rms bond = 0.001



# UNTANGLE: sim-data Ground Truth

- As **simple** as possible:
  - 2 conformers
  - 50:50 occupancy
  - flat bulk solvent
  - no false outliers
- **Accessible**
  - 64 residues
  - 78 waters
  - phenix, refmac, shelxl
  - clear score
- Every chance of **success**
  - 1.0 Å reso
  - “one thing wrong”
  - Experimental error only
  - perfect phases

# The UNTANGLE Challenge

**best.pdb**

**refme.mtz**

11) Build better than the best:

as-good or better wE,  $R_{free}$  as **best.pdb**  
non-equivalent ensemble

10) Recover **best.pdb** from **refme.mtz**

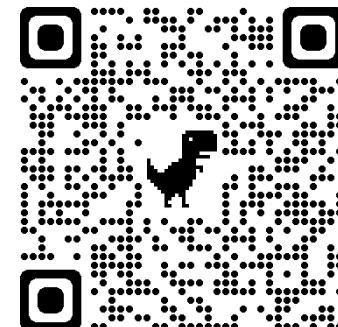
true phases, auto-building allowed

**no cheating** (don't use **best.pdb** as guide)



Prove: true ensemble is unique  
and can be recovered

<https://github.com/jmholton/UnTangle>



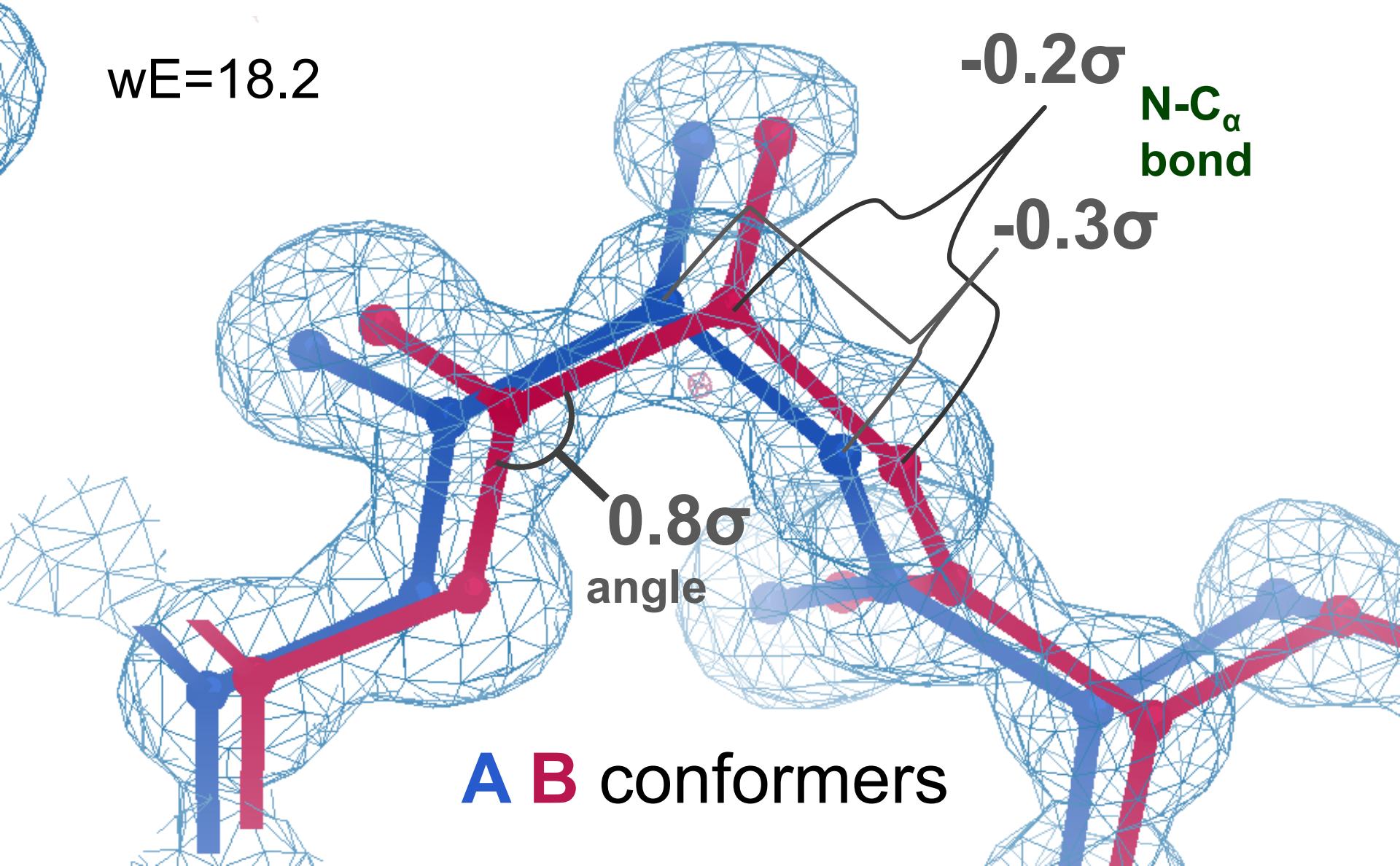
# verified Score sheet

wE	R <sub>work</sub>	R <sub>free</sub>	method
185	20.6	21.4	Alphafold2
109.6	8.87	9.44	Phenix.ensemble_refine
97.3	8.84	9.27	qFit best
91.7	3.64	4.27	manual
80.5	4.31	4.75	lotswrong
65.2	6.35	6.93	Amber24
54.6	5.87	6.57	single conf w aniso B
33.5	2.94	3.32	One thing wrong: Ala39
23.5	2.95	3.31	One thing wrong: Val1
21.4	3.12	3.48	Rectified Simulated annealing
18.1	17.9	18.4	Phenix autobuild
18.2	2.79	3.14	best (ground truth)

# Sim challenge data

Ala39

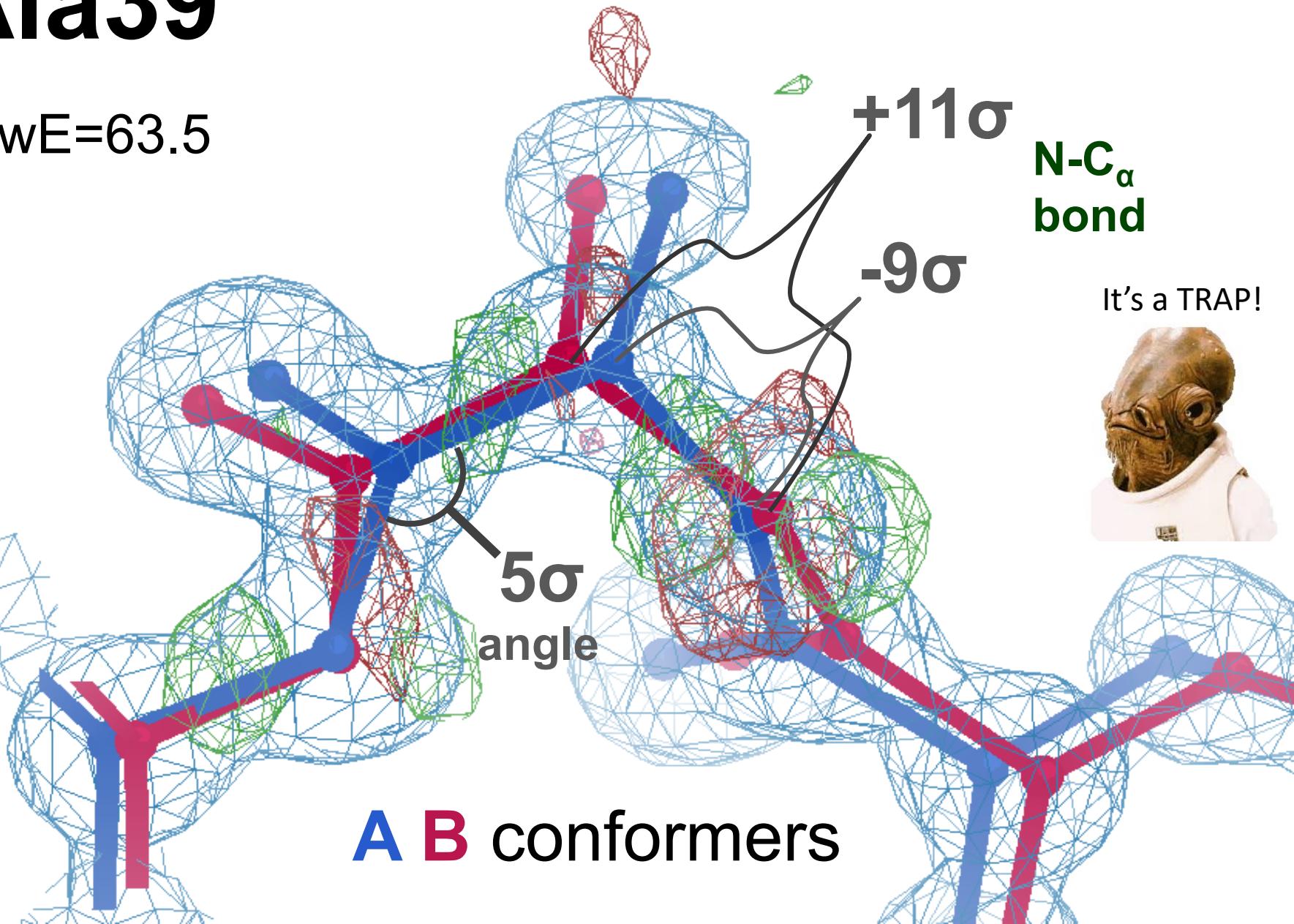
wE=18.2



# Sim challenge data: “otw39”

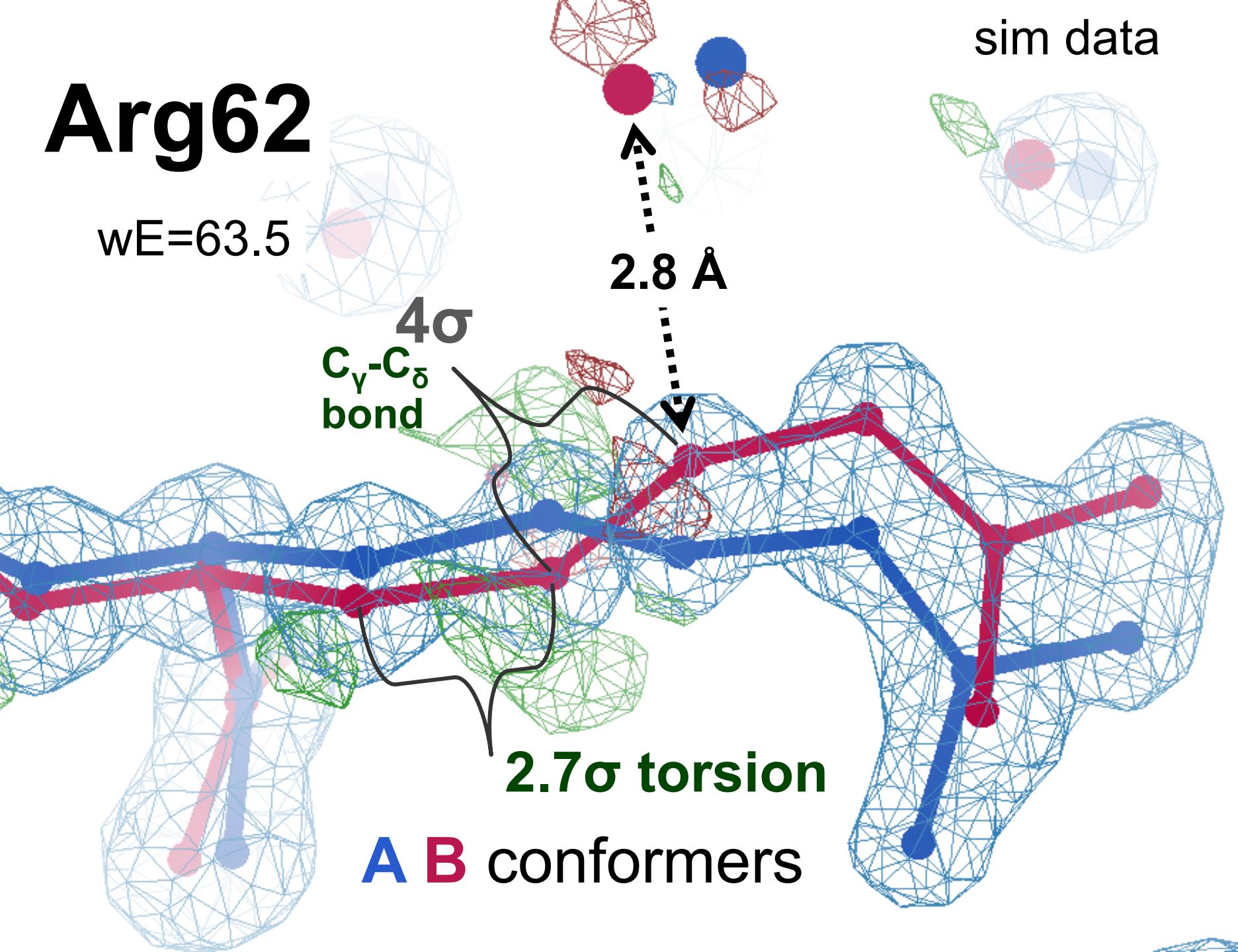
## Ala39

wE=63.5



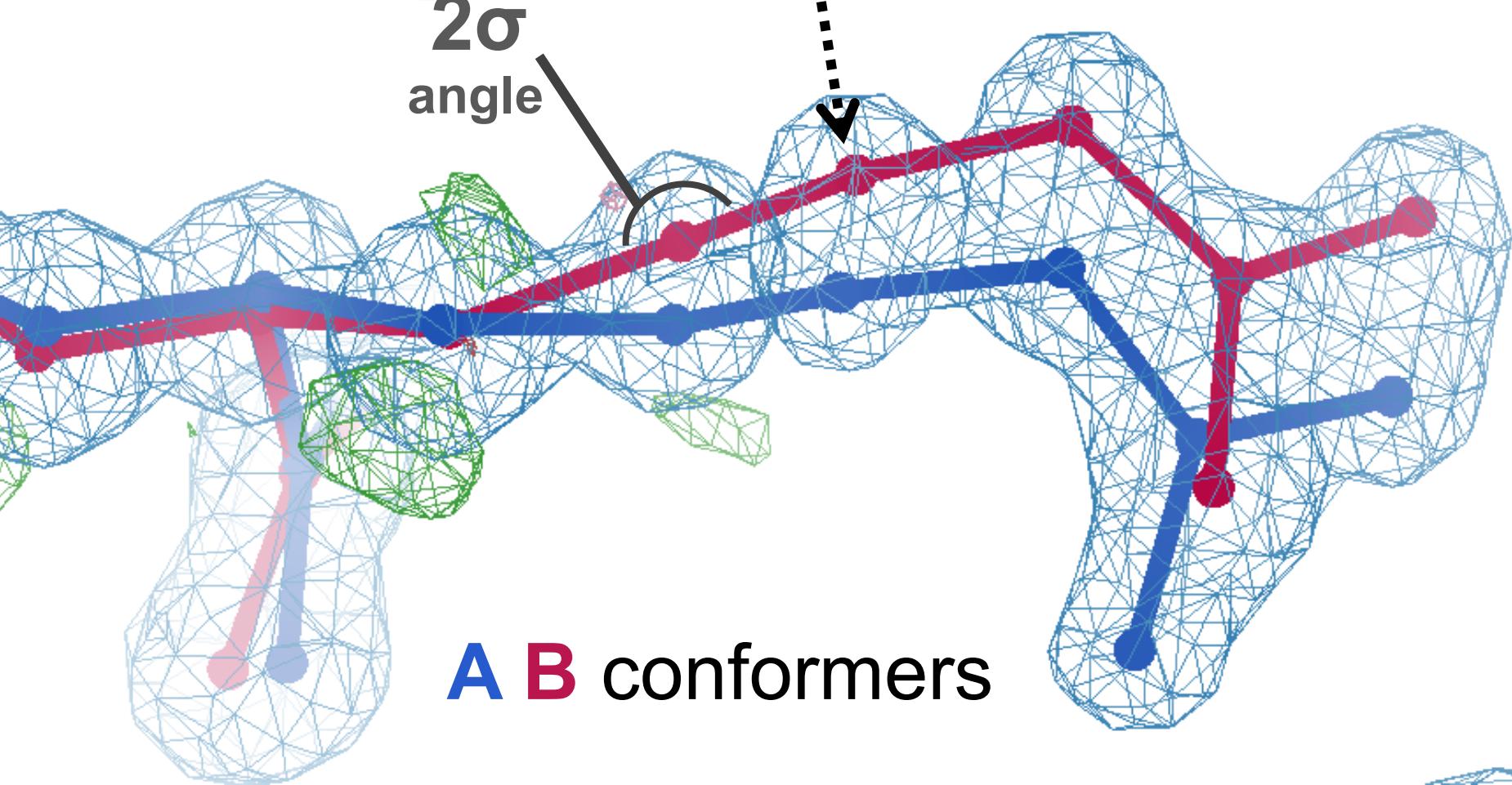
# Arg62

wE=63.5



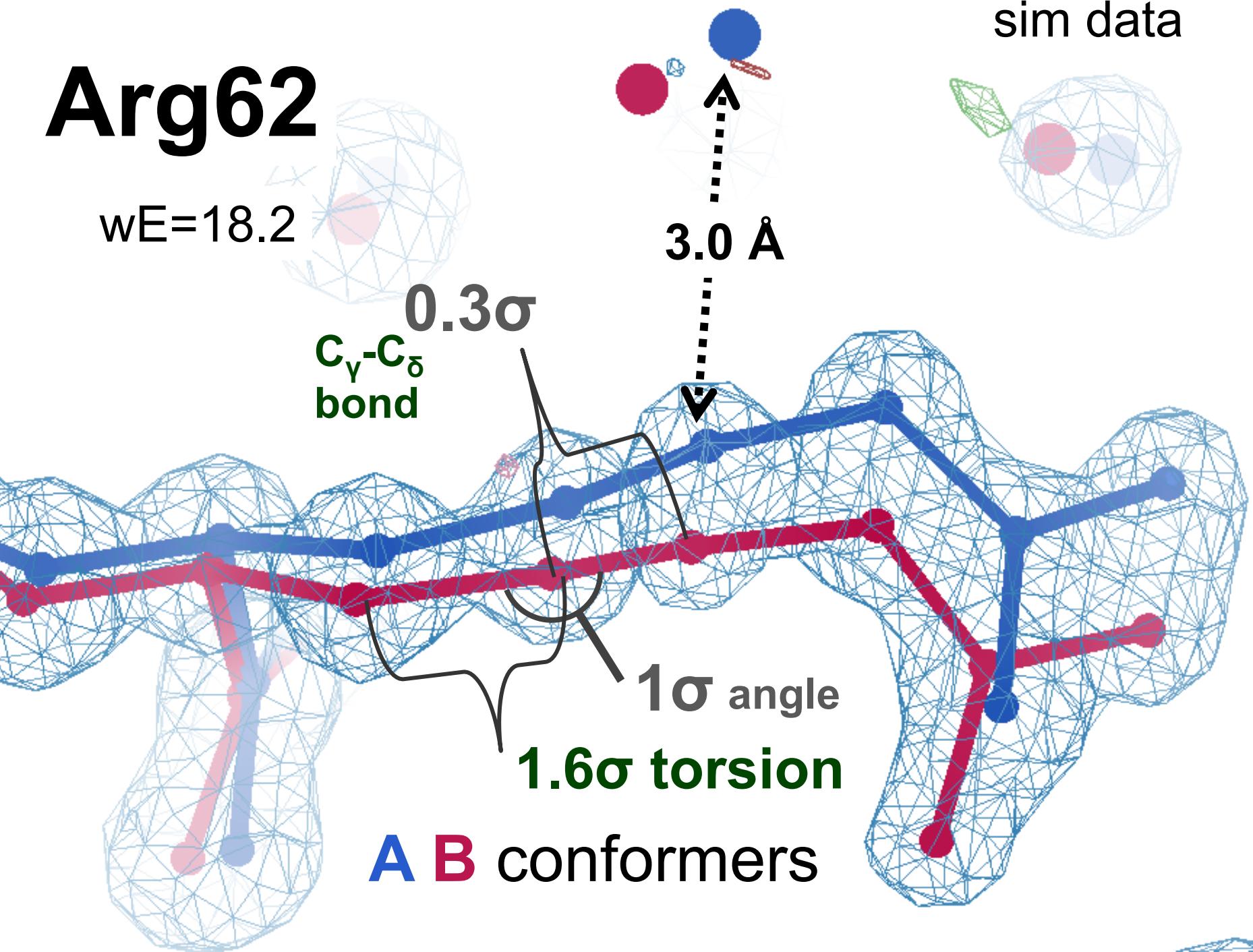
# Arg62

wE=50.7



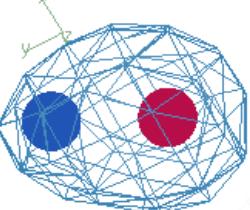
# Arg62

wE=18.2

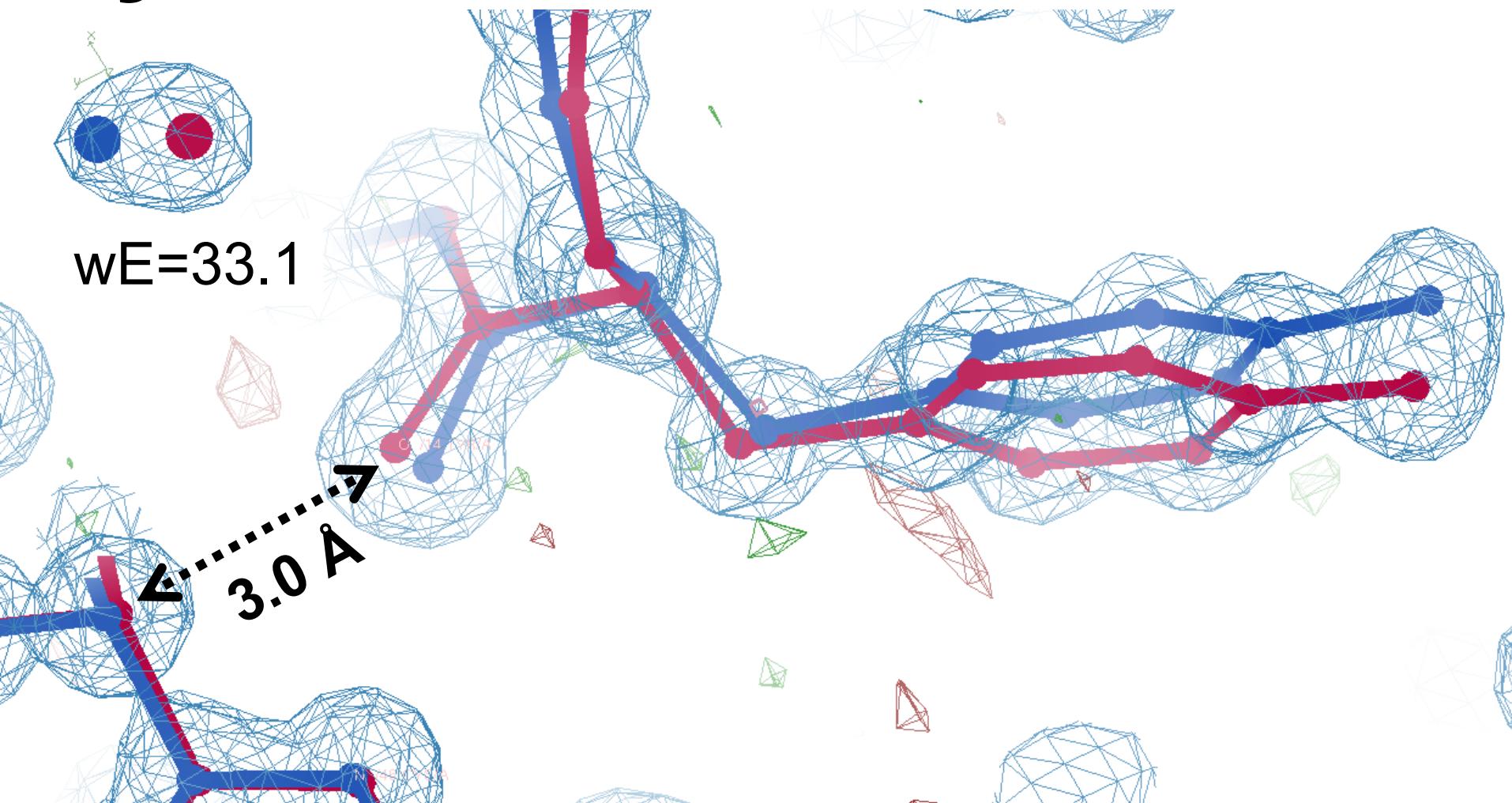


# Tyr14

sim data



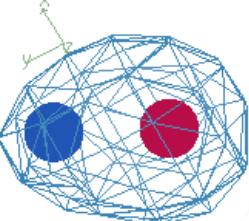
wE=33.1



A B conformers

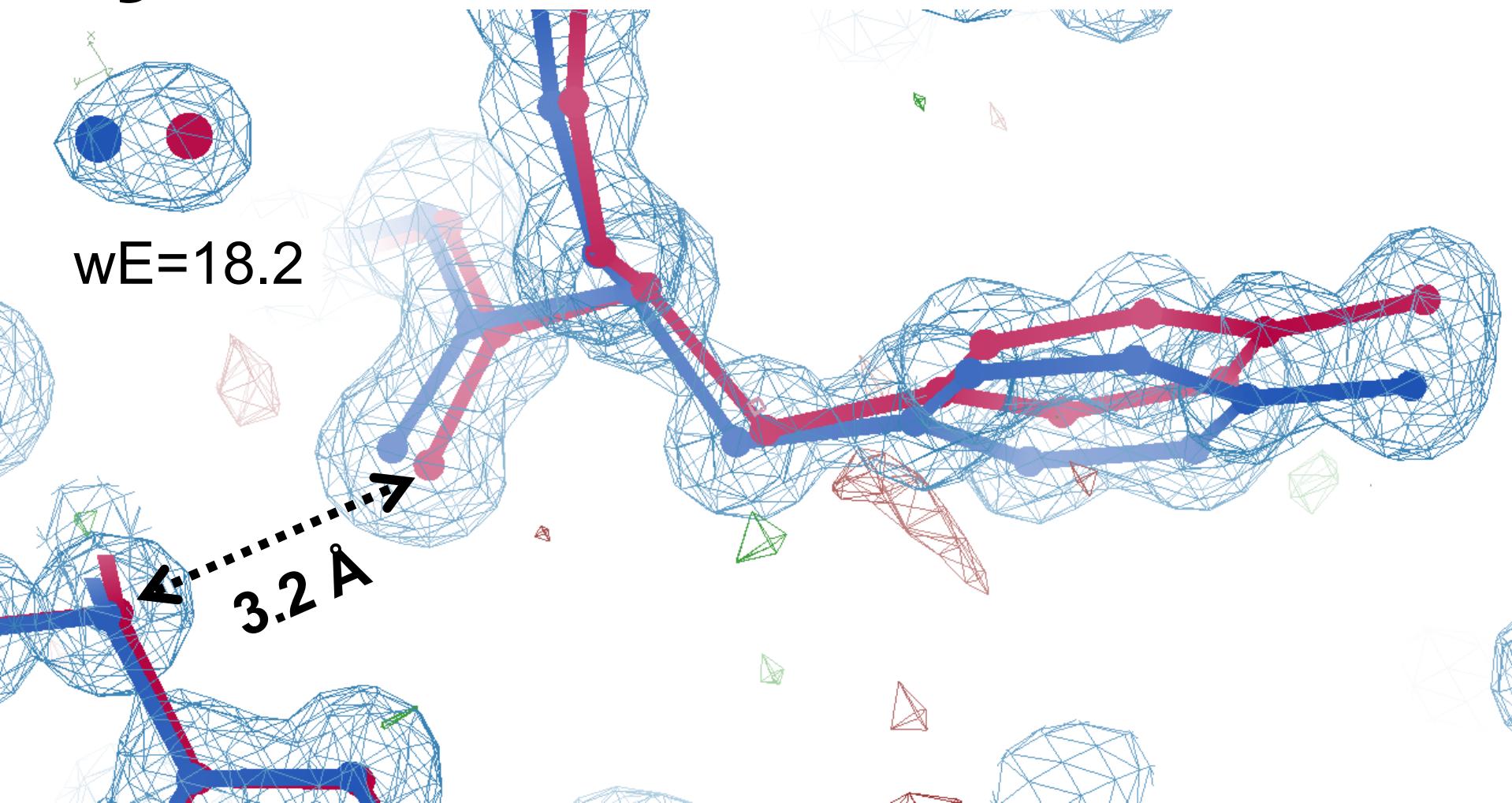
# Tyr14

sim data



3.2 Å

A B conformers



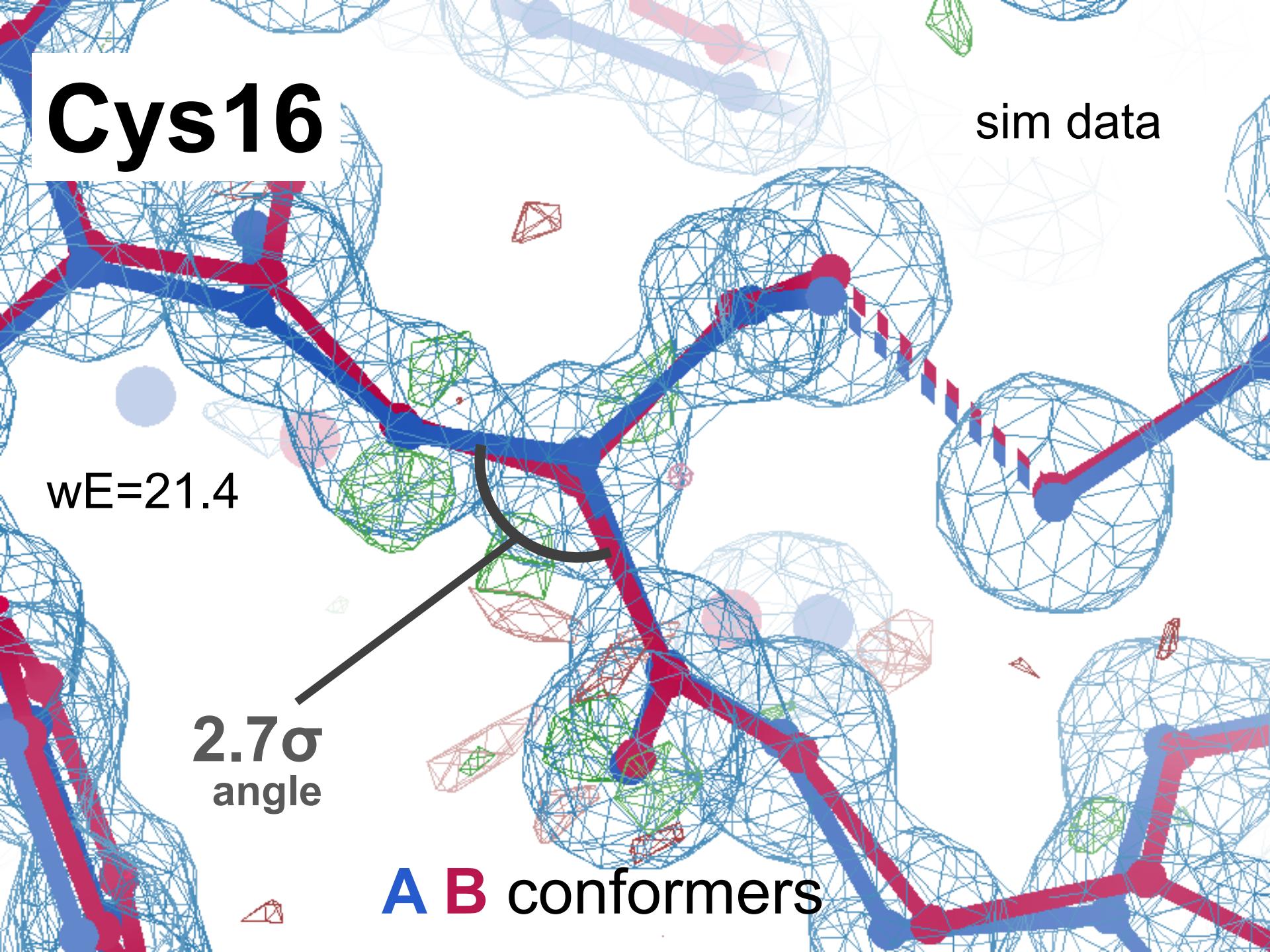
# Cys16

sim data

wE=21.4

$2.7\sigma$   
angle

A B conformers



# Cys16

sim data

$wE=18.2$

$1.5\sigma$   
angle

A B conformers

# wE score: key concepts

- Energy  $\xleftrightarrow{\text{sqrt}()}$  dev/sigma  $\xleftrightarrow{\text{erf}()}$  Probability
- Probability outlier is not noise
  - “clip” outlier energy at “10”
  - Not so sharp transition to outlier
- Tighter Non-bond (Lennard-Jones)

# Statistical Potential

$$E = \left( \frac{v - v_0}{\sigma} \right)^2$$

$E$	“energy”
$v$	model value
$v_0$	ideal value
$\sigma$	rms deviation expected

Outlier problem:  $1 \times 6\sigma \leftarrow 40x 1\sigma$

solution: weight by  
“probability its not noise”

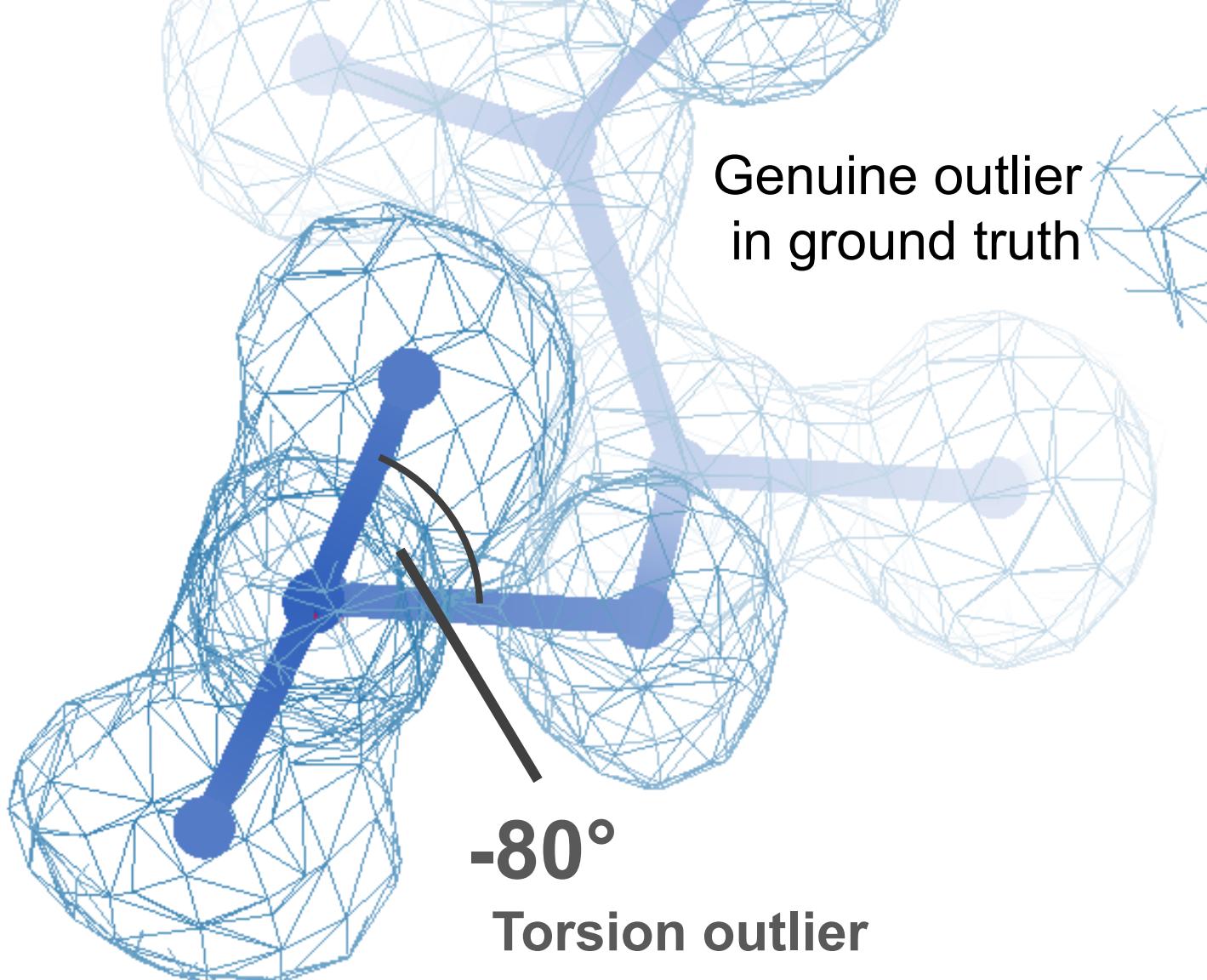
# Glu32



real data



4



# New Tools!

Tom  
Terwilliger



```
phenix.holton_geometry_validation test.pdb
```

```
phenix.create_alt_conf phenix_autobuild.pdb  
refme.mtz \  
nproc=48
```

Pavel  
Afonine

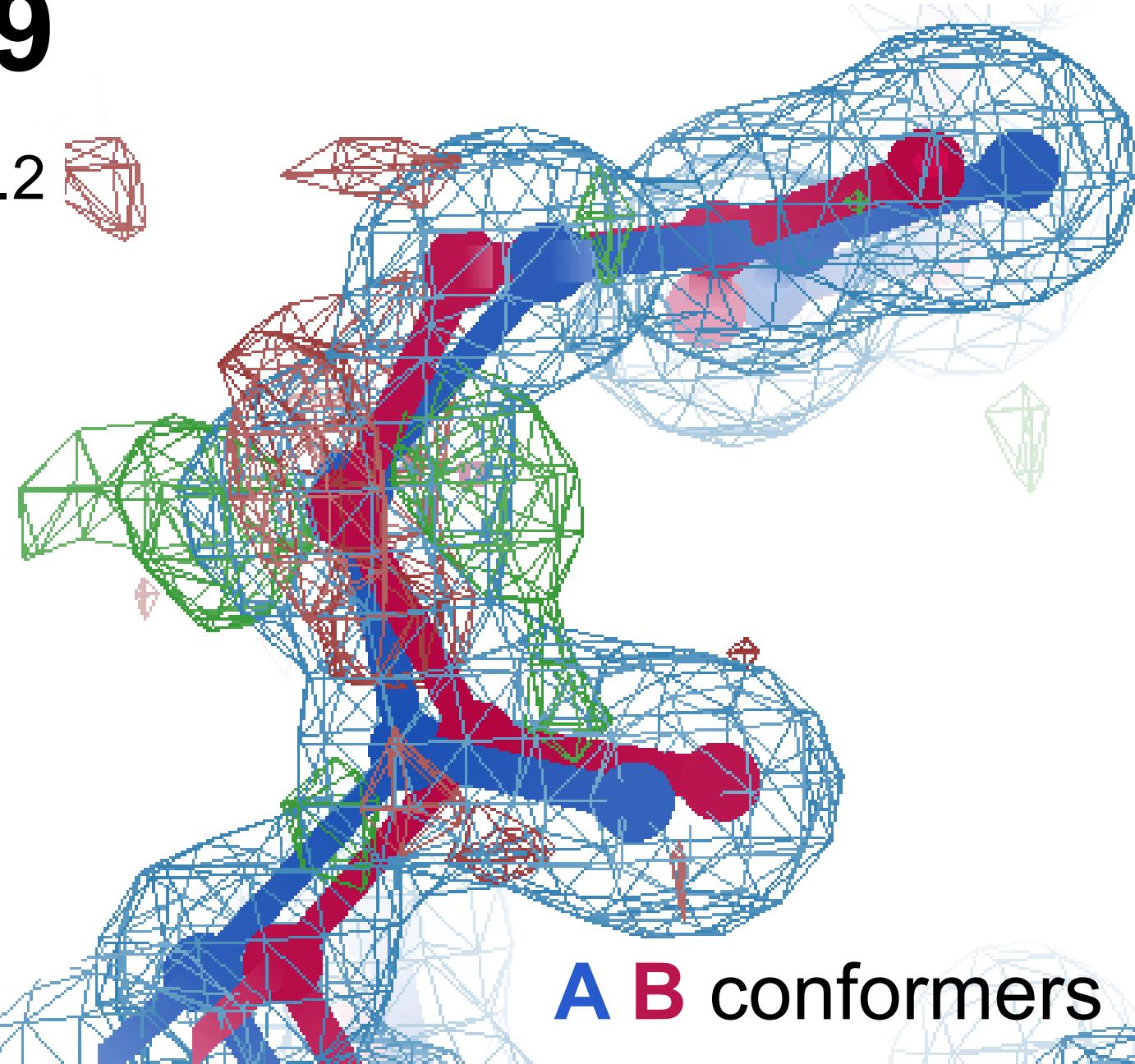
```
phenix.refine \  
model.pdb refme.mtz \  
main.number_of_mac=10 \  
fit_altlocs_method=masking \  
ordered_solvent=true \  
include_altlocs=true \  
ordered_solvent.mode=every_macro_cycle_after_first \  
refine_oat=true
```



# Better than the best?

# Ala39

wE=17.2

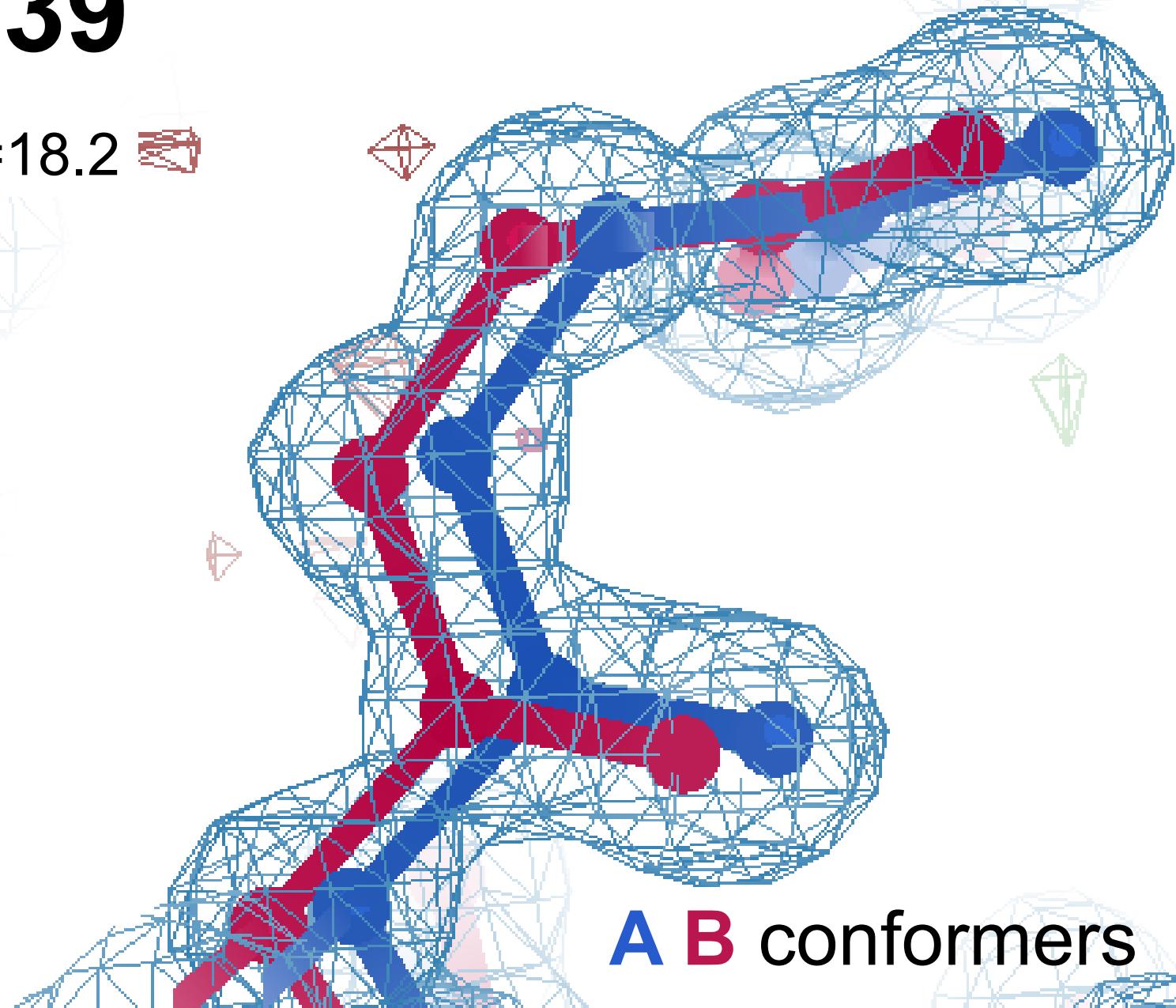


**A B** conformers

best?

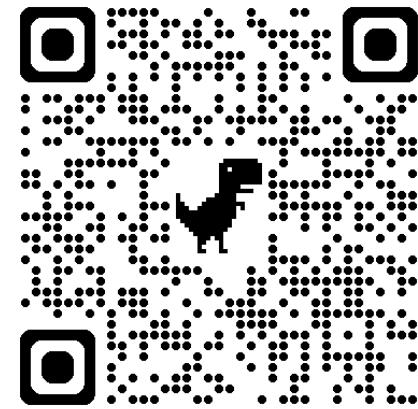
# Ala39

wE=18.2



# Summary

- Rfree < 3% → single e<sup>-</sup>
  - Not over-fitting
- Main barrier: tangled ensemble
- Sensible scoring function
  - Probability its not noise
  - some outliers are real
- Recover underlying ensemble
  - Might be required!
- Cooperative motions



\$1000

\$500

# verified Score sheet

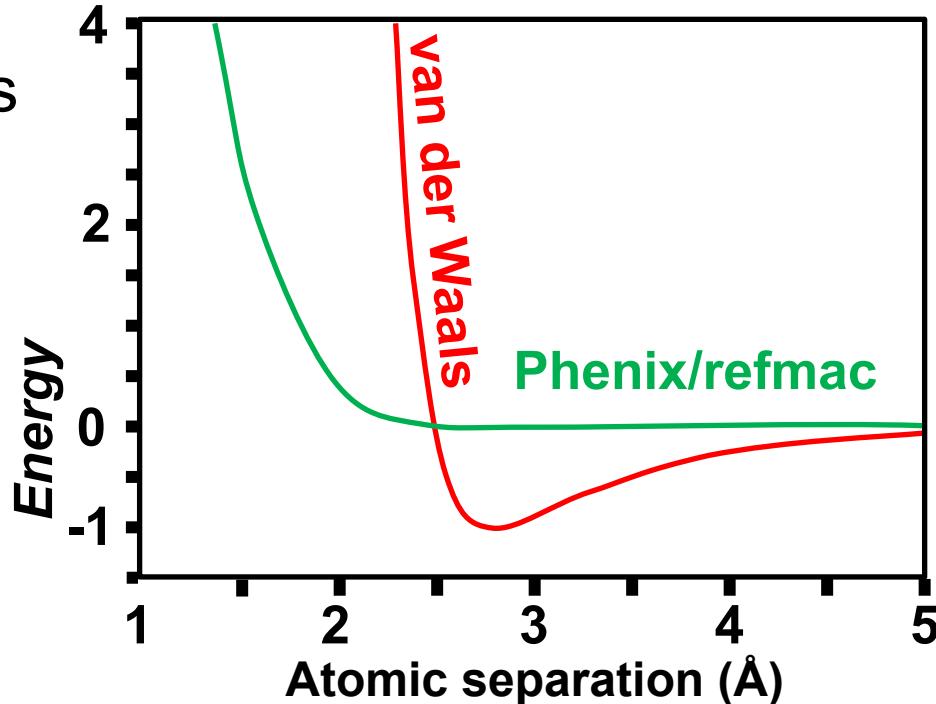
wE	R <sub>work</sub>	R <sub>free</sub>	method
185	20.6	21.4	Alphafold2
109.6	8.87	9.44	Phenix.ensemble_refine
97.3	8.84	9.27	qFit best
91.7	3.64	4.27	manual
80.5	4.31	4.75	lotswrong
65.2	6.35	6.93	Amber24
54.6	5.87	6.57	single conf w aniso B
33.5	2.94	3.32	One thing wrong: Ala39
23.5	2.95	3.31	One thing wrong: Val1
21.4	3.12	3.48	Rectified Simulated annealing
18.1	17.9	18.4	Phenix autobuild
18.2	2.79	3.14	best (ground truth)

# wE score: key concepts

- Probability outlier is not noise
  - “clip” outlier energy at “10”
  - Not so sharp transition to outlier
- Tighter Non-bond (Lennard-Jones)

# Sensible scoring function

- Statistical potential
  - bonds, angles, planes, torsions, chiral
- Lennard-Jones vdW
  - non-bond and clashes
- Probability  $\rightarrow \sigma$  dev
  - Rota, rama
- Peptide bond
  - $\sigma \rightarrow 5^\circ$
- $C_\beta$  dev
  - $\sigma = 0.05 \text{ \AA}$



# Probability that it's not noise

$$P_{nn} = \operatorname{erf}\left(\left|\frac{v - v_0}{\sigma\sqrt{2}}\right|\right)^{N_{things}}$$

$P_{nn}$  probability  
 $v - v_0$  deviation from ideal  
 $\sigma$  rms deviation expected  
 $\operatorname{erf}()$  integral of a Gaussian  
 $N_{things}$  number of trials

$$\operatorname{erf}(1.0/\sqrt{2}) = 68\%$$

$$\operatorname{erf}(1.0/\sqrt{2})^{10} = 2\%$$

$$\operatorname{erf}(2.0/\sqrt{2}) = 95\%$$

$$\operatorname{erf}(2.0/\sqrt{2})^{10} = 63\%$$

$$\operatorname{erf}(4.0/\sqrt{2})^{10,000} = 53\%$$

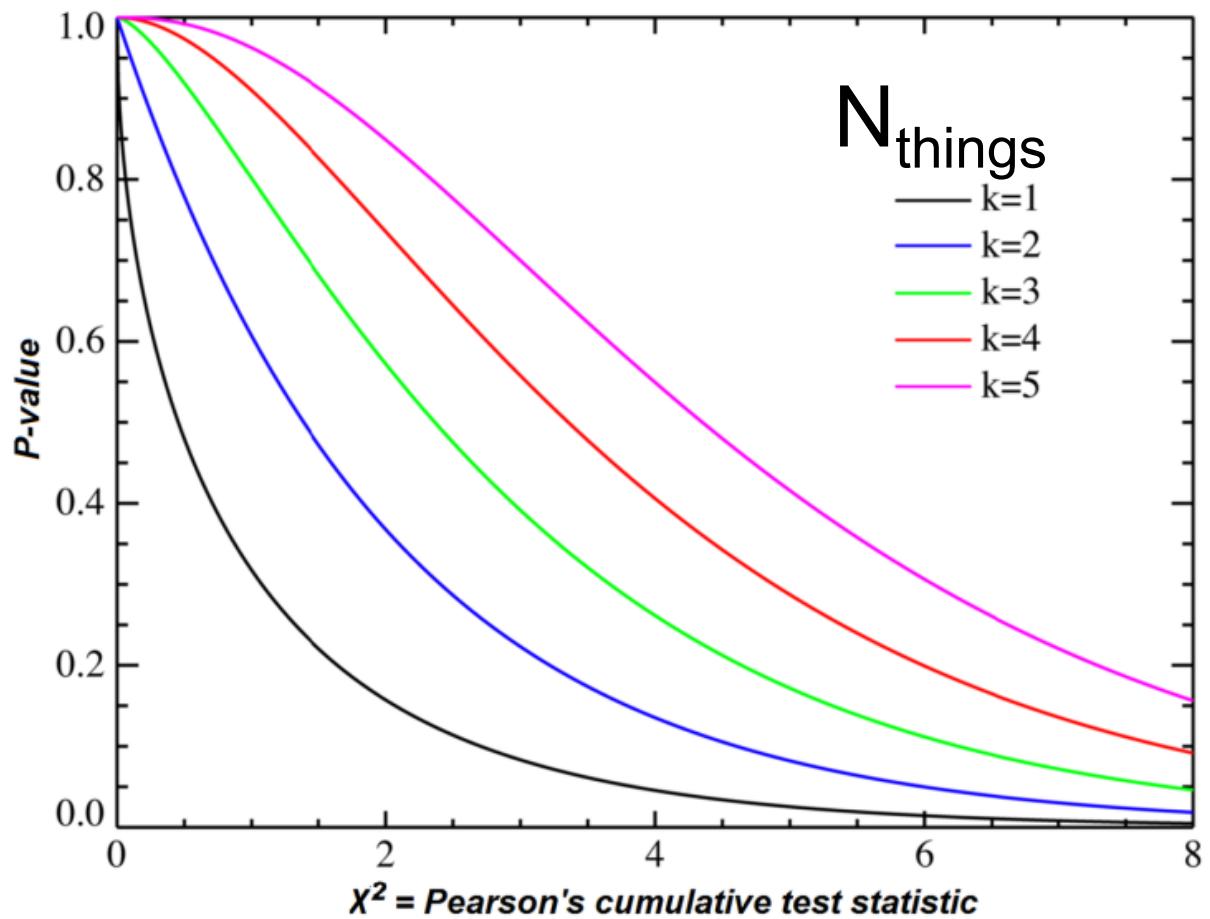
# Probability that it's not noise

average/sum of  $N_{\text{things}}$

$\chi^2$  test

Probability  
That its  
Gaussian

tie breaker



# “clip” for tolerating true outliers

$$\text{clip}(E) = \begin{cases} E, & \text{if } E < 10 \\ 10 + \log(E) - \log(10) & \end{cases}$$

Statistical Energy

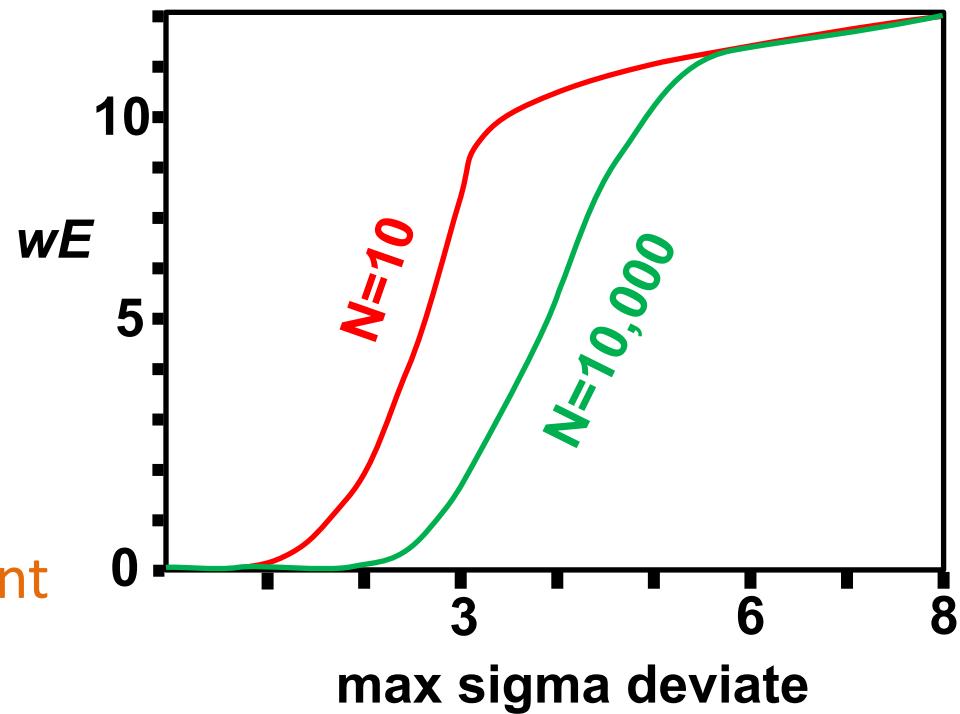
$$E = \left( \frac{\nu - \nu_0}{\sigma} \right)^2$$

Colloquially:

OK! OK! Its an outlier!

Let's not agonize about it

10 $\sigma$  vs 10.1 $\sigma$  **not** more important  
than 3.3 $\sigma$  vs 3.0 $\sigma$



# Probability that it's not noise

softer

$$P_{nn}^{soft} = 1 - 2 \left( \left( \frac{|\nu - \nu_0|}{\sigma_{P50}} \right)^5 \right)$$

$P_{nn}$  probability

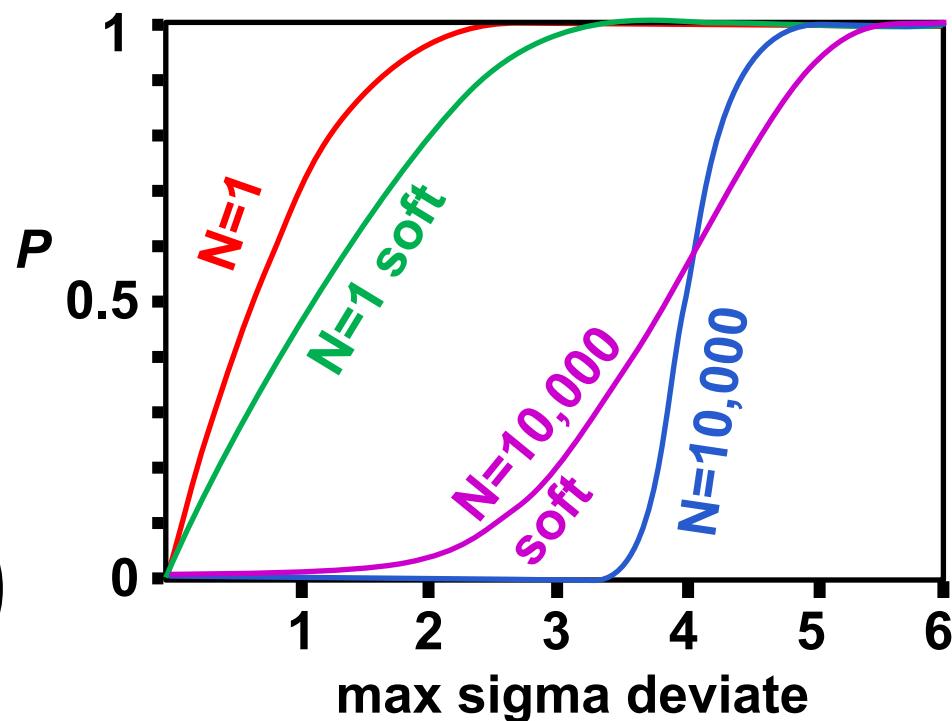
$\nu - \nu_0$  deviation from ideal

$\sigma$  rms deviation expected

inverf() inverse erf()

$N_{things}$  number of trials

$$\sigma_{P50} = \sqrt{2} \text{inverf}\left(0.5^{1/N_{things}}\right)$$



# weighted Energy scoring function

$$wE = \sum_{\text{validation types}} w_{avg} \langle E \rangle + w_{max} \max(\text{clip}(E))$$

`phenix.holton_geometry_validation`

$$w_{max} = P_{nn} \max(\text{clip}(E))$$

$$w_{avg} = \chi^2 \left( \sum_{N_{\text{bonds}}} E \right)$$

$$E = \left( \frac{v - v_0}{\sigma} \right)^2$$

`max()` = worst outlier

$$\text{clip}(E) = \begin{cases} E, & \text{if } E < 10 \\ 10 + \log(E) - \log(10) & \end{cases}$$

