# *Model Refinement*

## Oleg Sobolev

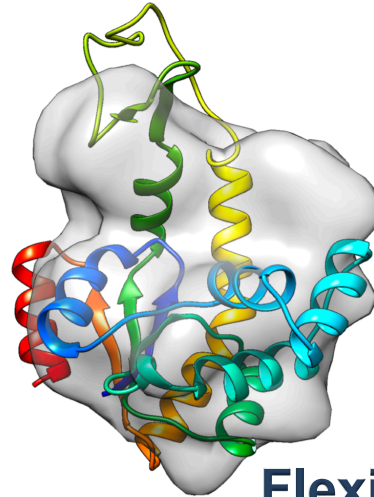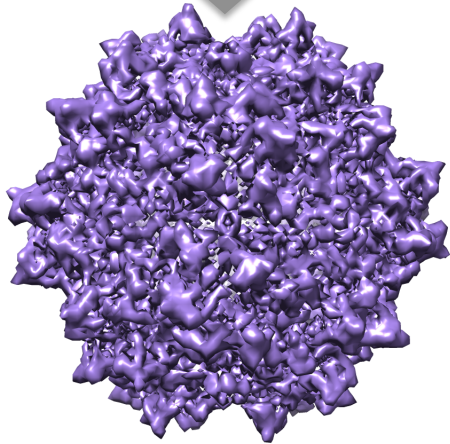**Phenix team**

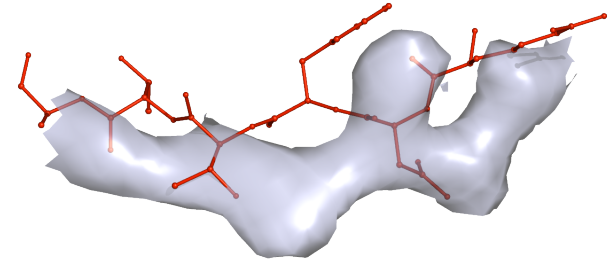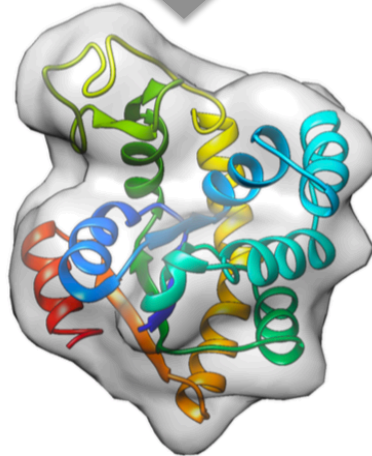**Lawrence Berkeley National Lab, California, USA**

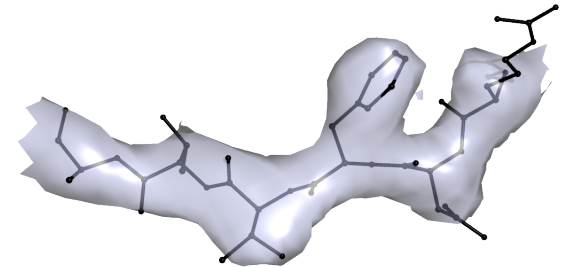# Model refinement vs other model fitting tools
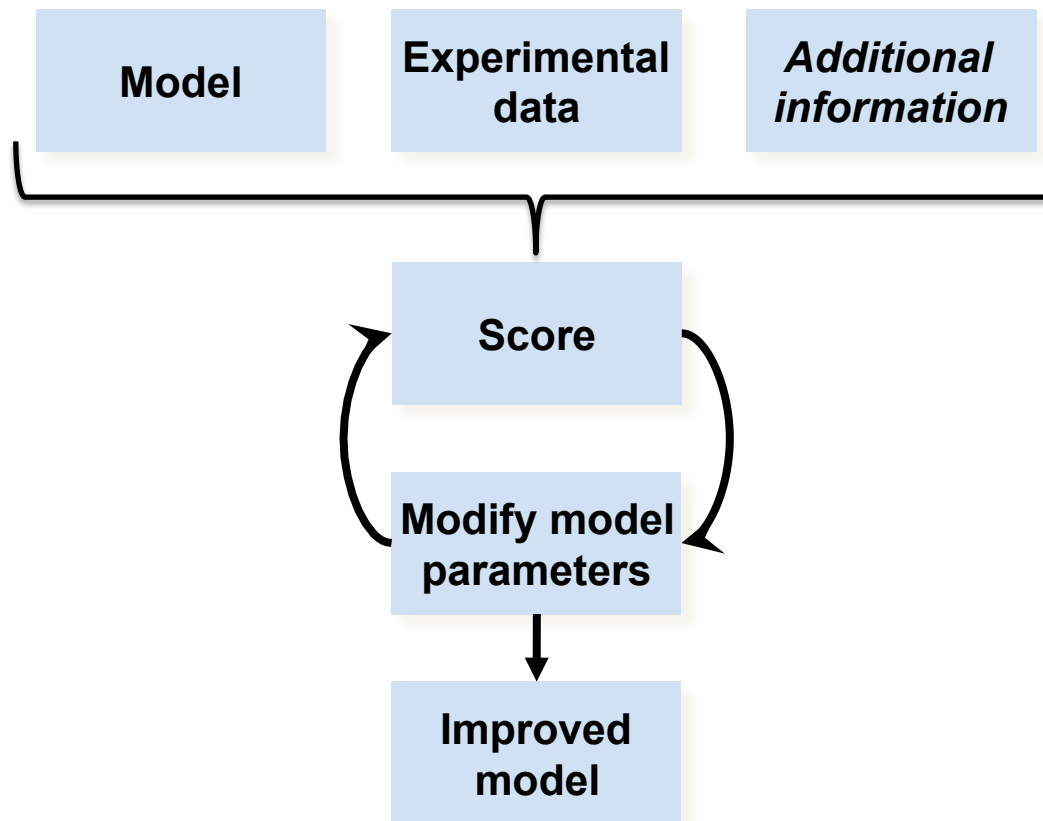


**Docking**

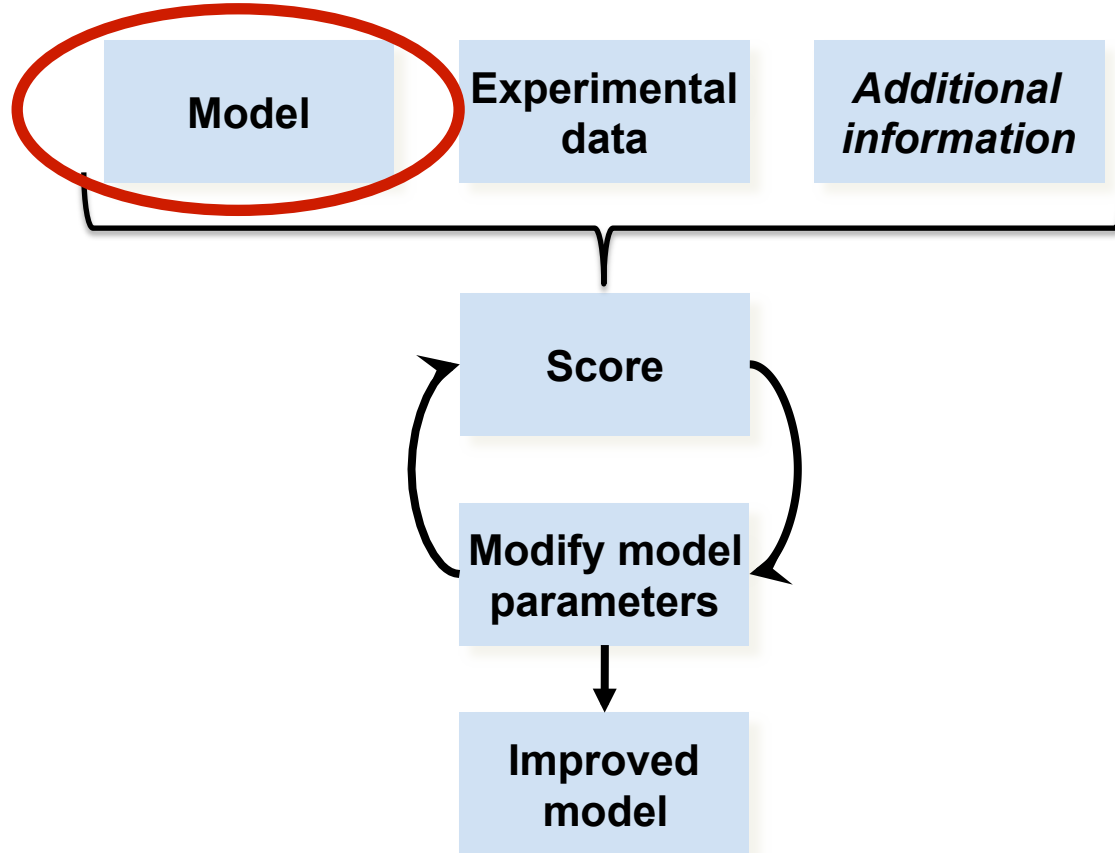**Flexible fitting, morphing**

**Refinement**

• All the above move model to achieve better fit. The difference is: by how much

# Model refinement



**Refinement – optimization process of fitting model parameters to experimental data**

# Refinement: model

# Atomic model parameters



*Position*    *Larger-scale disorder*

```
ATOM     25  CA  PRO A   4        31.309   29.489   26.044   1.00  57.79                    C
ANISOU   25  CA  PRO A   4        8443     7405     6110     2093    -24     -80             C
```

*Local mobility (harmonic vibrations)*

$$\mathbf{F}_{\text{MODEL}} = k_{\text{OVERALL}} \left( \mathbf{F}_{\text{CALC (ATOMS)}} + \mathbf{F}_{\text{BULK}} \right)$$
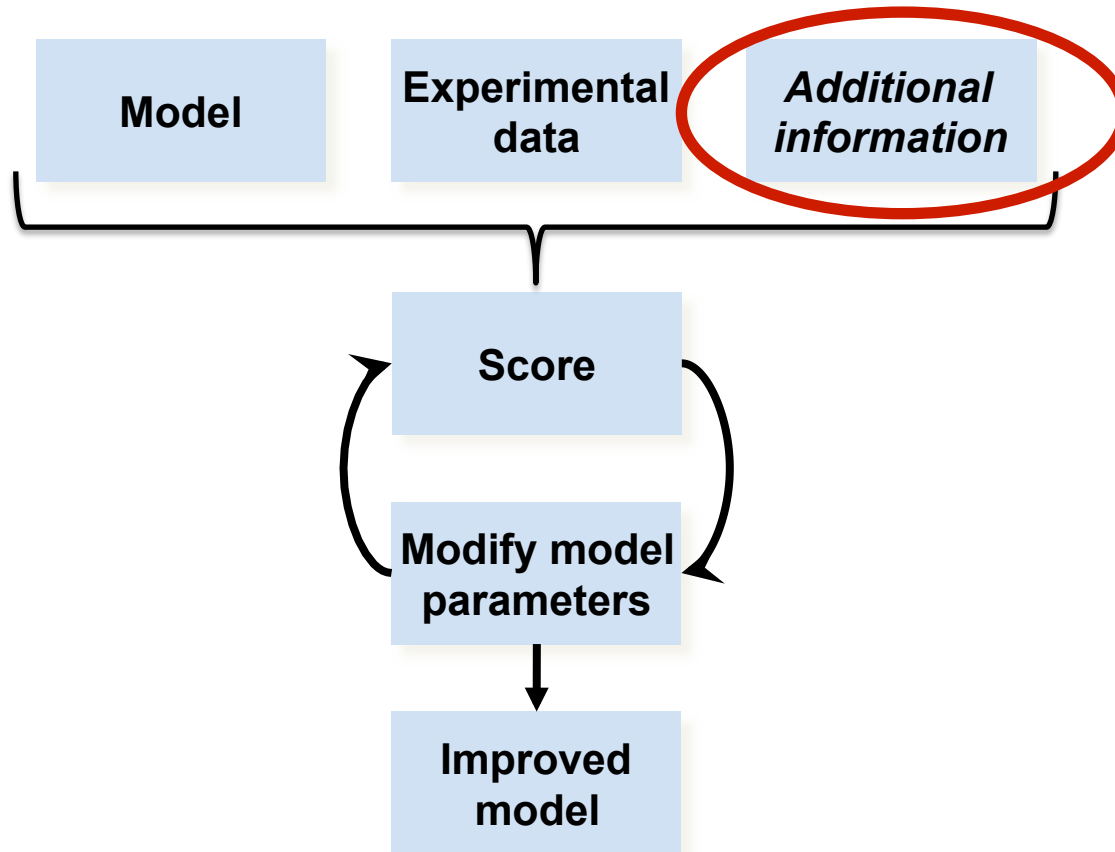
Occupancy **1.00**          **57.79** ADP (B-factor)

$$\mathbf{F}_{\text{CALC (ATOMS)}}(h,k,l) = \sum_{n=1}^{Natoms} q_n f_n(s) \exp\left( -\frac{B_n s^2}{4} \right) \exp\left( 2i\pi \mathbf{r}_n \mathbf{s} \right)$$

Atom type **C**          **31.309 29.489 26.044**

Atomic coordinates

# Additional information (restraints, constraints)

# Restraints for coordinate refinement

- The weight *w* balances data and restraints

$$T = T_{\text{DATA}}(F_{\text{OBS}}, F_{\text{MODEL}}) + wT_{\text{RESTRAINTS}}$$
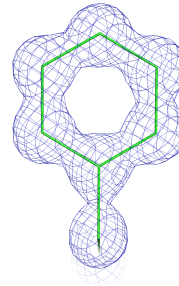
- Too much restraints: model may not adequately describe the data

- Too much data: model may not obey prior knowledge about model geometry

- Using optimal weight is very important
  - Programs know how to calculate it optimally
  - Sometimes programs fail to calculate it optimally
    - You need to be able to recognize this situation
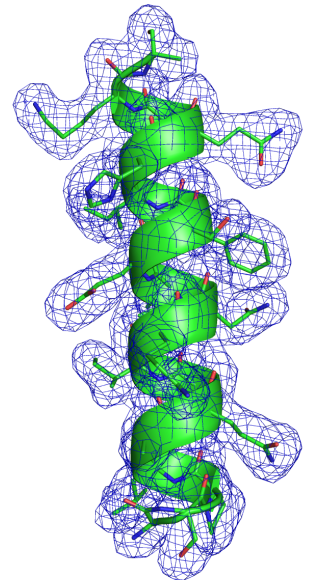
# Restraints in structure refinement

- Refinement target is usually a weighted sum of experimental data and *a priori* chemical knowledge terms

$$T = T_{\text{DATA}}(F_{\text{OBS}}, F_{\text{MODEL}}) + wT_{\text{RESTRAINTS}}$$
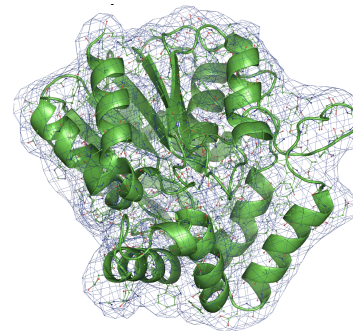
- At ultra-high resolution (<1Å) an unrestrained refinement sometimes may be possible.

- At 'typical' resolutions (1-3Å) *standard* restraints are necessary:
  covalent bond, angles, etc

- At lower resolution (lower than 3Å) more restraints needed:
  NCS, Secondary Structure, Ramachandran, …

# Restraints for coordinate refinement

$$T = T_{\text{DATA}}(F_{\text{OBS}}, F_{\text{MODEL}}) + wT_{\text{RESTRAINTS}}$$

$$T_{\text{RESTRAINTS}} = T_{\text{BOND}} + T_{\text{ANGLE}} + T_{\text{DIHEDRAL}} + T_{\text{PLANE}} + T_{\text{REPULSION}} + T_{\text{CHIRALITY}} + \ldots$$
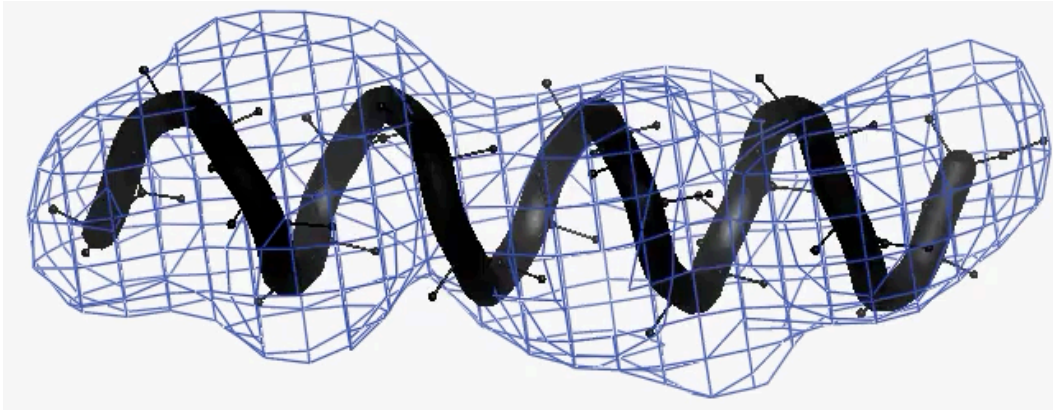
$$T_{\text{BOND}} = \Sigma_{\text{all bonded pairs}} \, w \, (d_{\text{ideal}} - d_{\text{model}})^2$$

From libraries (CCP4 Monomer Library or GeoStd in Phenix)
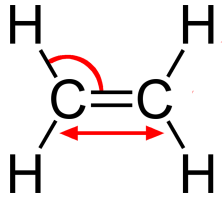
Calculated from actual atomic model

# Importance of additional restraints

- Toy example: refinement of a perfect α-helix into low-res map
  - Standard restraints on covalent geometry isn't sufficient
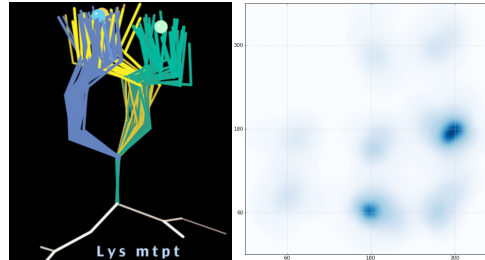    - Model geometry deteriorates as result of refinement
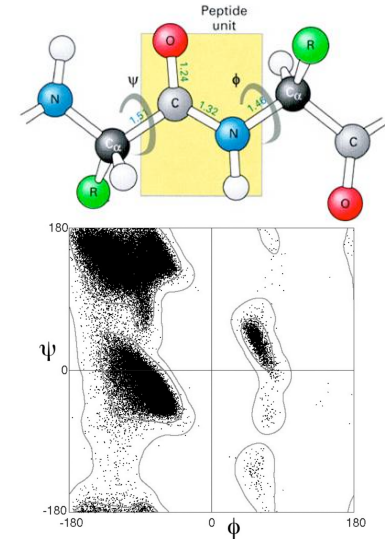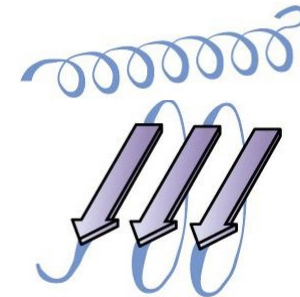
# Restraints for low resolution

Side chain distributions

Main chain distributions
(Ramachandran)

Covalent geometry

Internal
symmetry
(NCS)

Similar (homologous) structures
(reference model restraints)

Secondary structure

# Secondary structure restraints

## Helices



n+4

n

## Sheets



## Restraints used

H-bond lengths
H-bond angles

## Proteins

## Base pairs



## Stacking pairs



## Nucleic Acids

H-bond lengths
H-bond angles
Planarity
Parallelity

# Ramachandran plot restraints



- Ramachandran plot restraints
  - Use to stop outliers from occurring

Before refinement

After refinement





Good idea to use Ramachandran plot restraints!

# NCS (internal symmetry): constraints vs restraints



*Source: Internet*

- **Constraints**: molecules 1, 2 and 3 are required to be identical

- **Restraints**: molecules 1, 2 and 3 are required to be similar but not necessarily identical

# Refinement target function (score)

# Refinement target function (score)

$$T = T_{\text{DATA}}(F_{\text{OBS}}, F_{\text{MODEL}}) + wT_{\text{RESTRAINTS}}$$
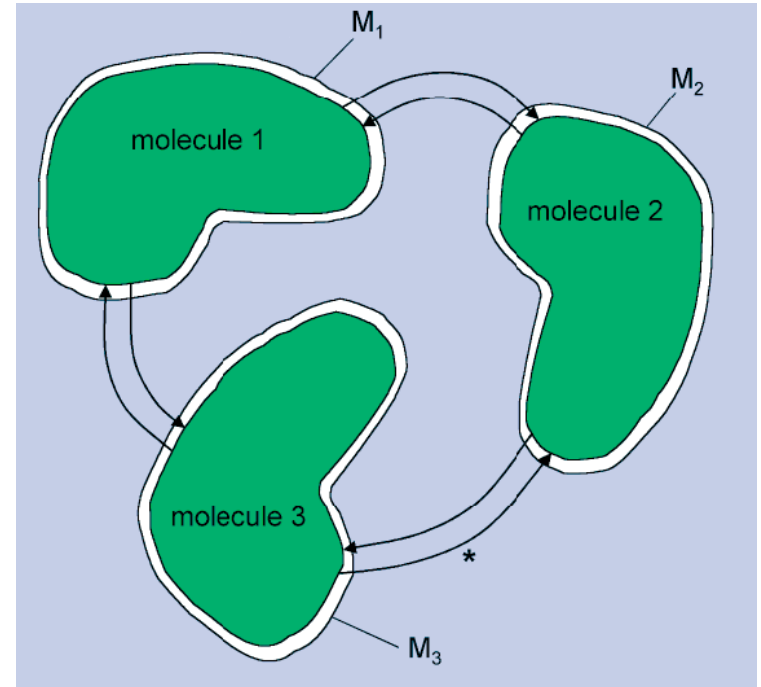
**Model**



$\mathbf{F}_{\text{MODEL}}$

**Data**



$\mathbf{F}_{\text{OBS}}$

$$\boldsymbol{T_{DATA}} = \sum_{hkl} (F_{obs} - F_{model})^2$$

$$\boldsymbol{T_{DATA}} = \sum_{hkl} \frac{||F_{obs}| - |F_{model}||}{|F_{obs}|}$$

$\boldsymbol{T_{DATA}} =$
Maximum-Likelihood score

# Refinement

## Crystallography



| Initial model | Experimental data | *A priori* knowledge |
|---|---|---|

Score

Modify model parameters

Improved model

phenix.refine
Available since 2005

## Cryo-EM



| Initial model | Experimental data | *A priori* knowledge |
|---|---|---|

Score

Modify model parameters

Improved model

phenix.real_space_refine
Available since 2013

# Atomic model refinement: crystallography vs cryo-EM

## Crystallographic refinement

- Improving model improves map
  - (2mFo-DFc, Model phase), (mFo-DFc, Model phase)
  - Better model leads to better map
  - Better map leads to more model built
  - Improving model in one place lets build more model elsewhere in the unit cell
  - Refine all model parameters (XYZ, B) from start to end of structure solution
  - Build solvent (ordered water) early
- Experimental data never changed
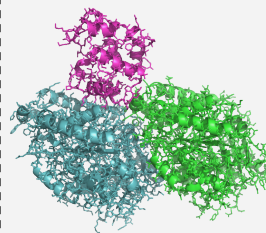- Data / restraints weight is global and time expensive to find best value
- Whole model needs to be refined

## Cryo-EM refinement

- Changing model does not change map
  - Build solvent (water) last
  - Get as complete and accurate model as possible before refining B factors and occupancies
- Experimental data changes a lot during the process (filtering, boxing, using maps with implied symmetry or not, etc.)
  - What map to use in refinement?
  - Refined B factors depend on map used
- Data / restraints weight can be local and is always optimal
- Boxed parts of the model can be refined

# Refinement: command line

- ## Real-space (cryo-EM)

  `phenix.real_space_refine model.pdb map.mrc resolution=3.4`

  `phenix.real_space_refine model.pdb map_coeffs.mtz`

  `phenix.real_space_refine model.pdb map_coeffs.mtz ligans.cif`

  `phenix.real_space_refine parameters.eff`

- ## Reciprocal-space (crystallography)

  `phenix.refine model.pdb data.mtz`

  `phenix.refine model.pdb data.mtz ligans.cif`

  `phenix.refine parameters.eff`

# Refinement tools in *Phenix*

# Refinement protocol



Data
Model
Parameters

Inputs

Rigid body

Rotamer fitting

Simulated Annealing

ADP

Refinement macro-cycle

Morphing

Occupancy

Weight calculation

XYZ minimization

Macro-cycle

(PDB or mmCIF)

Refined model

Trajectory

Log file

# Understanding <u>inputs</u> and outputs

- Real-space inputs

    - Atomic model (PDB, mmCIF)
    - Map (real map: MRC or Fourier map: MTZ)
    - Ligand restraints ("ligand CIF")
    - Parameter files (as command line arguments or a file)

- Reciprocal-space inputs

    - Atomic model (PDB, mmCIF)
    - Reflection data (typically MTZ but most other formats are OK)
    - Ligand restraints ("ligand CIF")
    - Parameter files (as command line arguments or a file)

# Understanding inputs and <u>outputs</u>

- Real-space outputs

  - Atomic model (PDB, mmCIF)
  - .log file
  - .eff file – summary of all input parameters
  - .geo file (optionally)

- Reciprocal-space outputs

  - Atomic model (PDB, mmCIF)
  - .log file
  - .eff file – summary of all input parameters
  - MTZ file with copy of input data and 2Fo-Fc and Fo-Fc maps
  - .geo file (optionally)

.geo file contains description of all the geometry restraints used in refinement

## Understanding inputs and <u>outputs</u>

- MTZ outputted by `phenix.refine` contains

  1. Verbatim copy of input data considered for use

  2. Data that was actually used in refinement

  3. Total model structure factors $\mathbf{F}_{model}$

  4. Fourier maps
     - $2mF_{obs} - DF_{model}$ 'filled'
     - $2mF_{obs} - DF_{model}$
     - $mF_{obs} - DF_{model}$
     - Anomalous difference map (if anomalous data)

# Refinement: practical considerations

# Aggressive optimization methods

- Simulated annealing (SA)

- Model morphing

  - Only use if model has gross errors (correction requires large movements)

  - Do not use if model is relatively good and only needs small corrections

# Use Hydrogen atoms

- Half of the atoms in a protein molecule
- Make most interatomic contacts
- Add to model towards the end, data resolution does not matter
- Once added, do not remove before the PDB deposition
- H do contribute to R-factors (expect 0.1-2% drop in R)



**A structure without (left) and with (right) hydrogen atoms**

# Know when to stop refinement



Colored bars are histograms showing distribution of values for structures at similar resolution

The black polygon shows where the statistics for the user's structure fall in each histogram

**Crystallographic model quality at a glance.**
L.Urzhumtseva, P.V.Afonine, P.D.Adams & A.Urzhumtsev. *Acta Cryst.* **D**65, 297-300 (2009)

# Know when to stop refinement

## Likely overall good model



| Average B | RMSD(angles) |
|---|---|
| 8.9 | 0.88 |
| 15.6 | 1.90 |
| 26.4 | 2.75 |

| RMSD(bonds) | R-free |
|---|---|
| 0.004 | 0.116 |
| 0.018 | 0.189 |
| 0.027 | 0.260 |

| R-work |
|---|
| 0.107 |
| 0.156 |
| 0.218 |

## Clearly there are problems



| Average B | RMSD(angles) |
|---|---|
| 8.9 | 0.88 |
| 35.0 | 0.29 |
| 26.4 | 2.75 |

| RMSD(bonds) | R-free |
|---|---|
| 0.004 | 0.116 |
| 0.001 | 0.387 |
| 0.027 | 0.260 |

| R-work |
|---|
| 0.107 |
| 0.385 |
| 0.218 |

# Low resolution (3Å or worse)

- Use:

  - Ramachandran plot restraints

  - Secondary structure restraints

  - Reference model restraints (if quality homology model is available)

  - NCS (restraints or constraints)

# NCS (Non-crystallographic symmetry)

- Constraints vs restraints
  - Constraints:
    - 4-5 Å or worse
    - Highly symmetric molecules
  - Restraints:
    - 2-4 Å

- Torsion vs Cartesian NCS
  - Torsion is preferable in most cases

- Symmetry related copies:
  - Can be found automatically as part of refinement
  - Can be specified manually
  - Automatic determination relies on model quality
    - Always check automatically detected NCS copies

# Secondary structure (SS) restraints

- Always use at 3Å and worse
- Better than 3Å: use if needed
- Require SS annotation
- SS annotation must be accurate
  - Errors in SS annotation may propagate into refined model

- Secondary structure (SS) annotation
  - SS information
    - HELIX/SHEET records in PDB file or equivalent in mmCIF
    - *Phenix* generated parameter files
  - Tools to create SS annotation
    - Command line (*phenix.secondary_structure_restraints*)
    - *Phenix* GUI
  - Quality of SS annotation:
    - Depends on quality of input model (GIGO)
    - No software can annotate SS fully reliably and correctly
    - Manual validation and editing almost always required
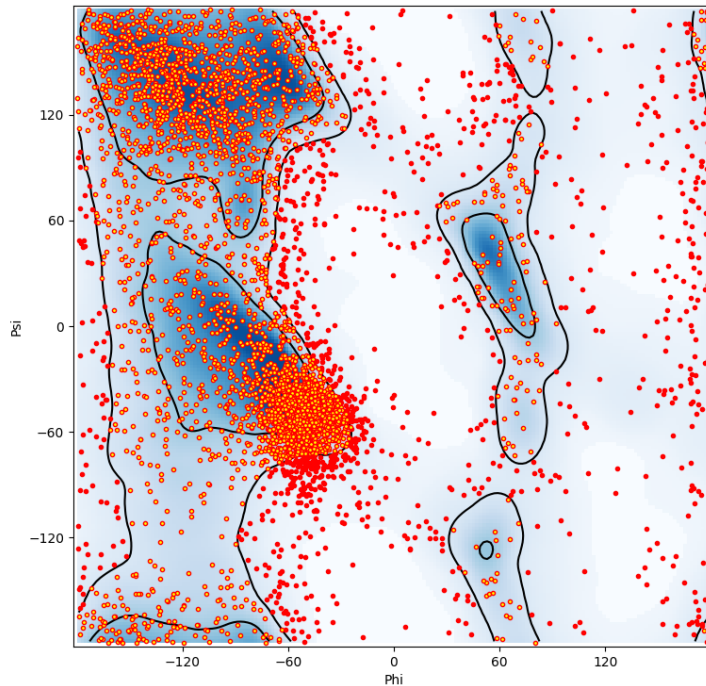
# Ramachandran plot restraints

- Likely need at about 3Å and worse

- Better than 3Å: use if needed (preserve good initial model from deterioration)

- Check Ramachandran plot regularly

- Don't use to fix outliers. Fix outliers first (manually), then use Ramachandran plot restraints to stop re-occurring outliers.

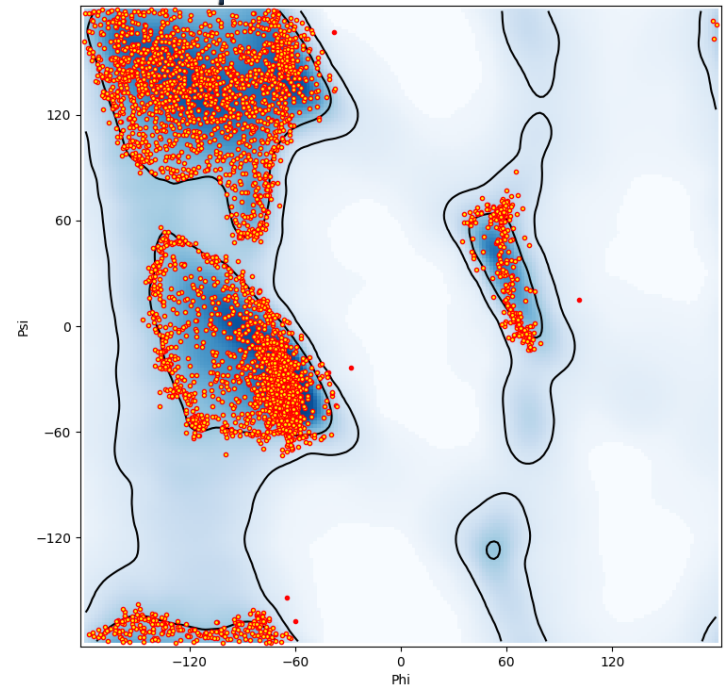# Ramachandran plot restraints

- Ramachandran plot restraints
  - Don't use to fix outliers. Fix outliers first, then use Ramachandran plot restraints to prevent re-occurring outliers.

**PDB code: 5a9z**

Original
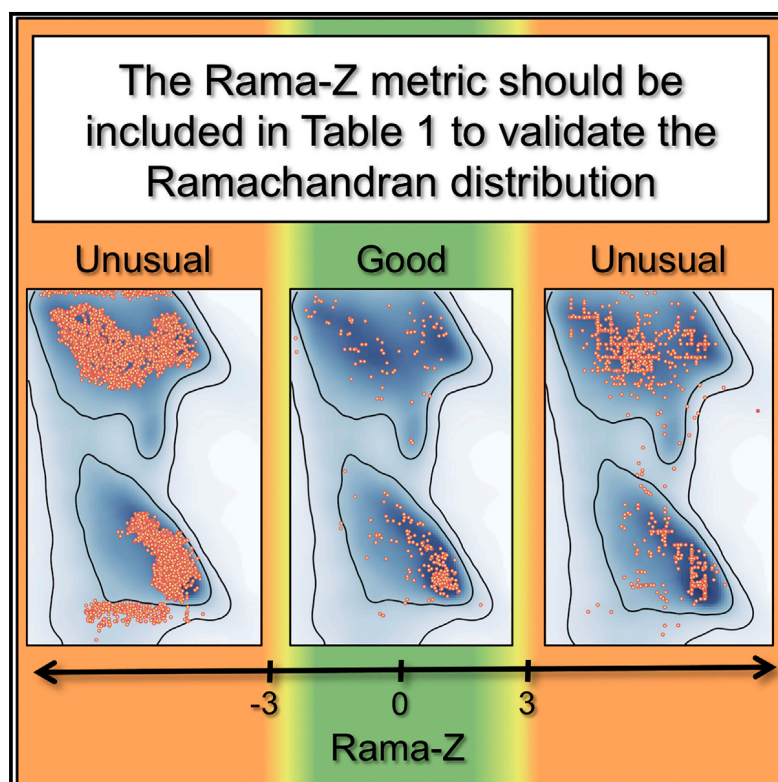
Refined with Ramachandran plot restraints



Bad idea to use Ramachandran plot restraints in this case. Fix outliers first!

# Rama-Z score

## Structure

## A Global Ramachandran Score Identifies Protein Structures with Unlikely Stereochemistry

**Graphical Abstract**

**Authors**

Oleg V. Sobolev, Pavel V. Afonine,
Nigel W. Moriarty,
Maarten L. Hekkelman,
Robbie P. Joosten,
Anastassis Perrakis, Paul D. Adams

**Correspondence**

osobolev@lbl.gov (O.V.S.),
r.joosten@nki.nl (R.P.J.)

**In Brief**

Counting the number of Ramachandran outliers is not sufficient for protein backbone validation. Sobolev et al. revisited the underutilized Ramachandran $Z$ score. The authors describe its reimplementation in Phenix and PDB-REDO and showcase its utility. They advocate including it in the validation reports provided by the Protein Data Bank.

Talk tomorrow (Aug 23) 1:30 pm, room 209 Session A 20

# Refinement: practical considerations

- Final stages

  - Make the model as complete as possible

  - Build alternative conformations

  - Use Hydrogen atoms (and keep them in the final model!)

  - Add ordered solvent components

- Remember: the better the model, the better the map

  - You may see and model your ligands better!

# Reading

# Phenix resources



Phenix paper

Video tutorials

Documentation

Relevant papers

Bi-annual newsletters

Slides from workshops

# User support

- **Feedback, questions, help**

  Mailing list (anyone signed up):        phenixbb@phenix-online.org
  Bug reports (developers only):        bugs@phenix-online.org
  Ask for help (developers only):        help@phenix-online.org

- **Reporting a bug or asking for help:**

  - We can't help you if you don't help us to understand your problem

  - Make sure the problem still exist using the latest *Phenix* version

  - Send us all inputs (files, non-default parameters) and tell us steps that lead to the problem

  - All data sent to us is kept confidentially

# The Phenix Project

**Lawrence Berkeley Laboratory**

Paul Adams, Pavel Afonine, Dorothee Liebschner, Nigel Moriarty, Billy Poon, Christopher Schlicksup, Oleg Sobolev

**BERKELEY LAB**

**Los Alamos National Laboratory**
**New Mexico Consortium**

Tom Terwilliger, Li-Wei Hung

New Mexico CONSORTIUM

Los Alamos NATIONAL LABORATORY

**UTHealth**

Matt Baker, Corey Hyrc

UTHealth®
The University of Texas
Health Science Center at Houston

**University of Cambridge**

Randy Read, Airlie McCoy, Tristan Croll, Claudia Millán Nebot, Rob Oeffner

UNIVERSITY OF CAMBRIDGE

**Duke University**

Jane & David Richardson, Christopher Williams, Vincent Chen