

Refinement

DLS-CCP4 Data Collection and Structure Solution Workshop

December 13-20 2016

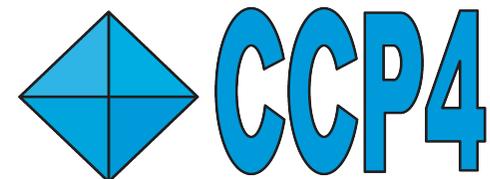
Diamond Light Source, Oxfordshire, UK

Oleg Kovalevskiy

okovalev@mrc-lmb.cam.ac.uk

Rob Nicholls

nicholls@mrc-lmb.cam.ac.uk

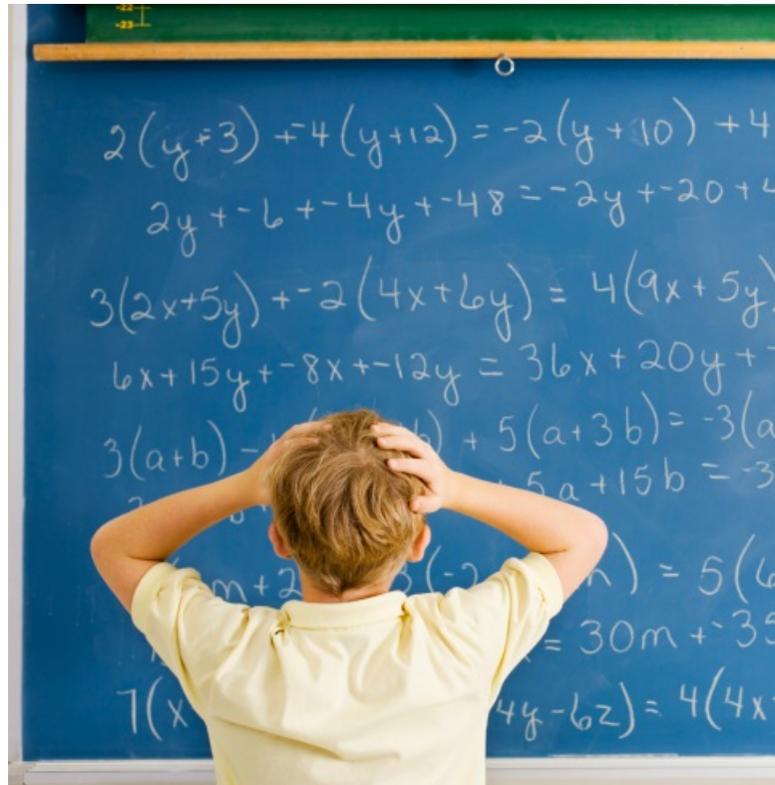


Contents

- **Refinement**
 - Purpose of Refinement
 - Crystallographic Data
 - Model Parameterisation
 - Restraints
- **Low-resolution Refinement**
 - Jelly-body restraints
 - ProSMART External Restraints
 - ProSMART Comparative Structural Analysis
 - LIBG Restraints
 - LORESTR Automated Pipeline
 - Map Sharpening and Blurring

Contents

**(almost) No math due to time constraints,
focus on practical usage**



Purpose of Refinement

Crystallographic refinement has one major purpose:

to fit atomic model into observed X-ray crystallographic data

Model should agree with the observed data

Model must be chemically and structurally sensible

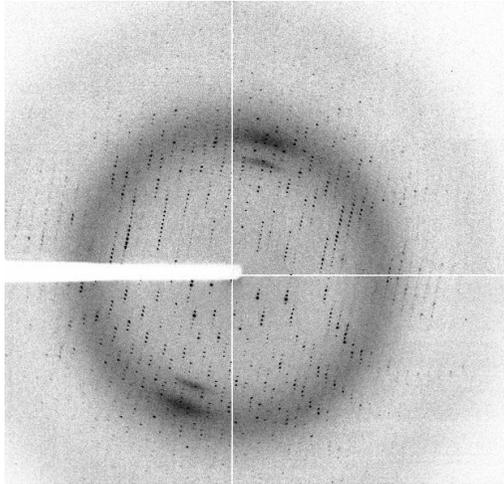
and one immediate most valuable consequence:

to calculate best possible electron density map

Allowing the atomic model to be visualised, criticised and analysed

Direct relation between model quality and phase quality (corresponds to the quality of the electron density maps)

Fourier Transform

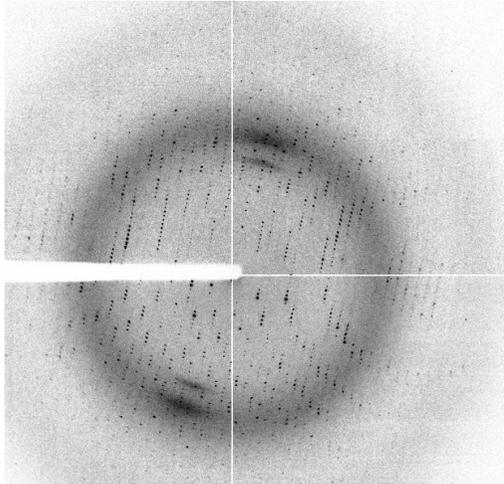


H, K, L	$ F_{\text{obs}} $	ϕ
...		
5, 5, 5	348	-
5, 5, 6	392	-
5, 5, 7	157	-
5, 5, 8	312	-
...		

We have observed amplitudes: $|F_{\text{obs}}|$

But we don't have phases: ϕ

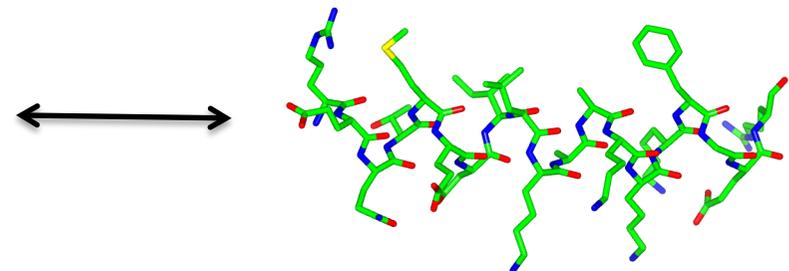
Fourier Transform



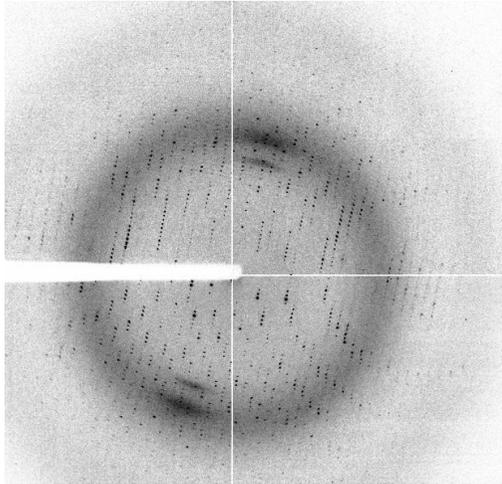
H, K, L	$ F_{\text{obs}} $	ϕ
...		
5, 5, 5	348	-
5, 5, 6	392	-
5, 5, 7	157	-
5, 5, 8	312	-
...		

Suppose we have a starting model:

H, K, L	$ F_{\text{calc}} $	ϕ_{calc}
...		
5, 5, 5	355	27°
5, 5, 6	387	8°
5, 5, 7	146	75°
5, 5, 8	340	31°
...		



Model Refinement



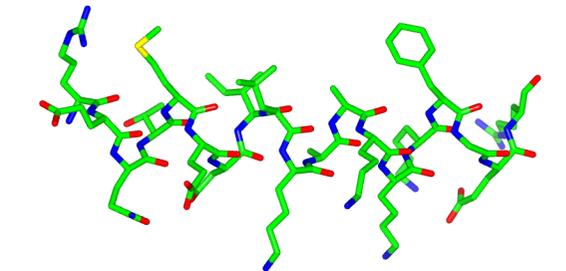
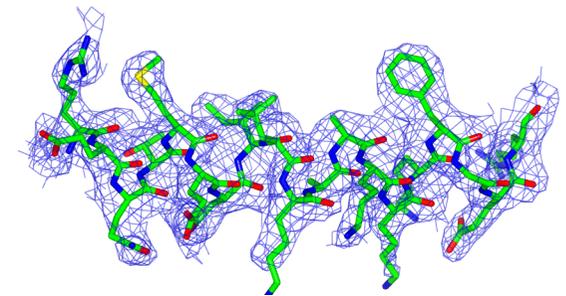
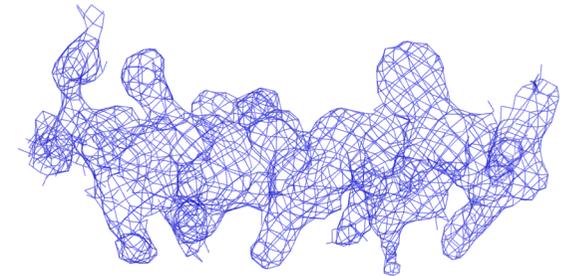
Idea:

Iteratively improve the model, optimising the agreement between

$|F_{obs}|$ and $|F_{calc}|$

Purpose: improve phase estimates: ϕ_{calc}

H, K, L	$ F_{obs} $	ϕ
...		
5, 5, 5	348	-
5, 5, 6	392	-
5, 5, 7	157	-
5, 5, 8	312	-
...		



H, K, L	$ F_{calc} $	ϕ_{calc}
...		
5, 5, 5	355	27°
5, 5, 6	387	8°
5, 5, 7	146	75°
5, 5, 8	340	31°
...		

Model Refinement

Idea:

Iteratively improve the model to optimise the agreement between $|F_{\text{obs}}|$ and $|F_{\text{calc}}|$

Note – we are not refining against a density map

We are optimising the agreement between $|F_{\text{obs}}|$ and $|F_{\text{calc}}|$

How to assess correspondence between the model and experimental observations?

R-factor:
$$R = \frac{\sum ||F_{\text{obs}}| - |F_{\text{calc}}||}{\sum |F_{\text{obs}}|}$$

Model Refinement

Refinement essentially tries to minimise the R-factor

How do we know that the model is reliable?

What if we improve the amplitudes $|F_{calc}|$ but worsen the phases φ_{calc} ?

Such overfitting can happen if there are too many parameters

How to validate?

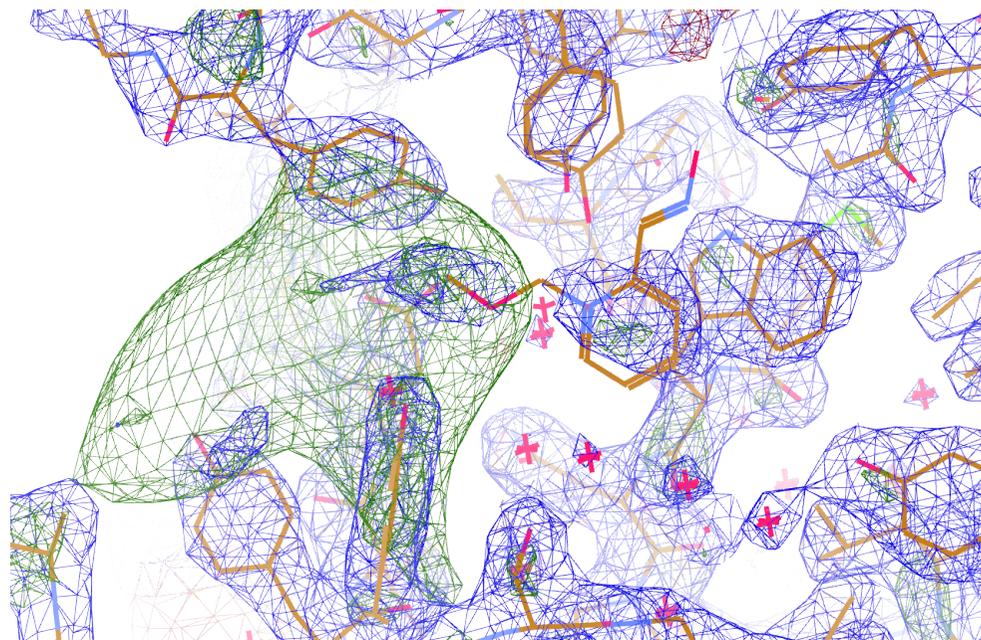
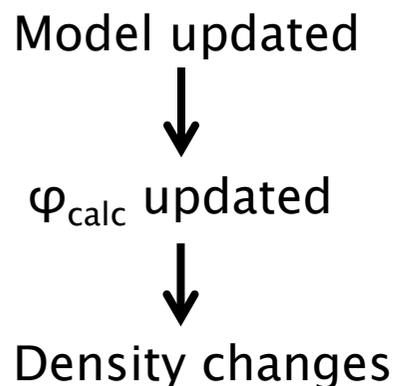
- **R_{free}** – reserve a portion of data for cross-validation (usually 5%)
- **Chemical & structural validation** – ensure that the model is physically sensible
- **Inspect electron density map** – manual intervention

Map Calculation

Two types of maps:

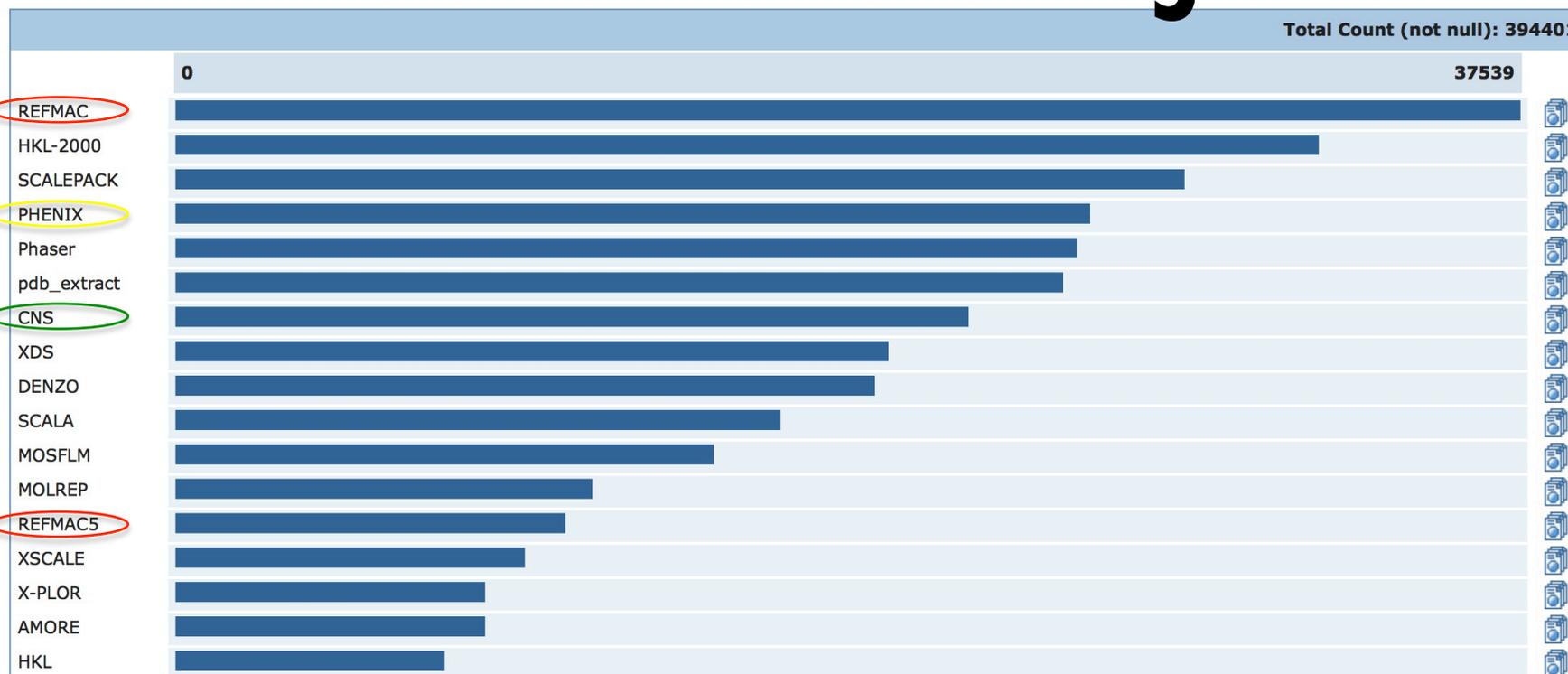
- $2F_{\text{obs}} - F_{\text{calc}}$: “standard” electron density – represents crystal contents
- $F_{\text{obs}} - F_{\text{calc}}$: *difference density* – represents differences

Maps are calculated using phase estimates from the current model: φ_{calc}



Note – contrast with real space refinement

Available Refinement Programs



- SHELXL
- CNS
- REFMAC5
- TNT
- BUSTER/TNT
- phenix.refine
- RESTRAINT
- MAIN
- MOPRO
- XD

Your Crystal Peculiarities

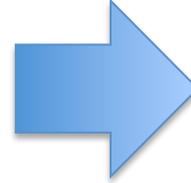
Why there is no single universal refinement protocol?

Why do you need to adjust refinement parameters?

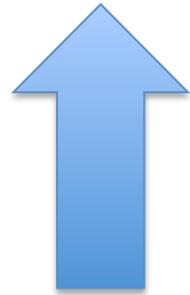
Three important things to consider:

1. You describe your data with the model. Data of different quality shall be described differently (mathematically speaking, different number parameters of the model could be estimated given particular data).
2. The model is defined not only by the PDB file, but also in the parameters of the refinement program.
3. Different quality of the model at different stages of refinement

Model (PDB file)
Atomic coordinates,
B-factors



Refined
model



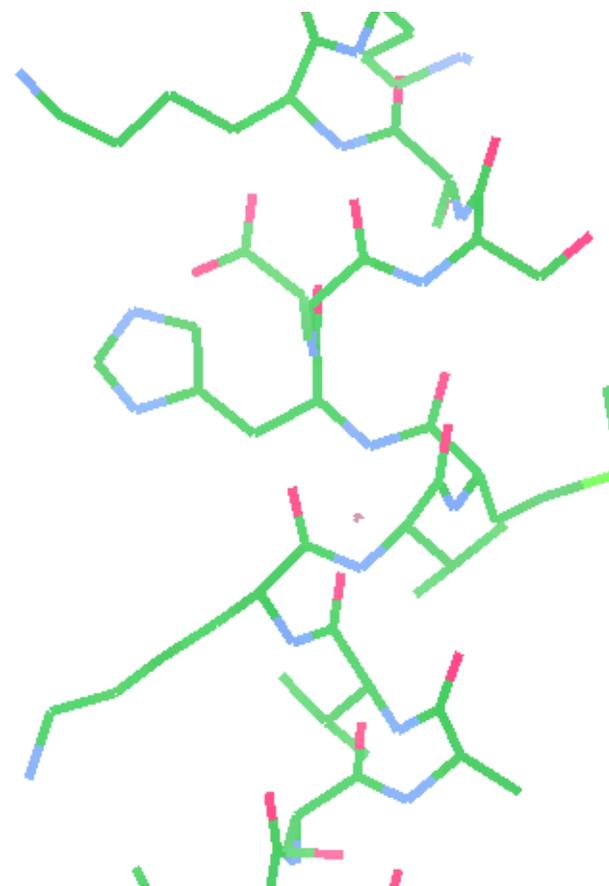
X-ray data
(experiment)

Model Parameterisation

Standard refinable parameters

Atomic model:

- Position – (x,y,z) coordinates
- Uncertainty – B-factors
- (Occupancies)



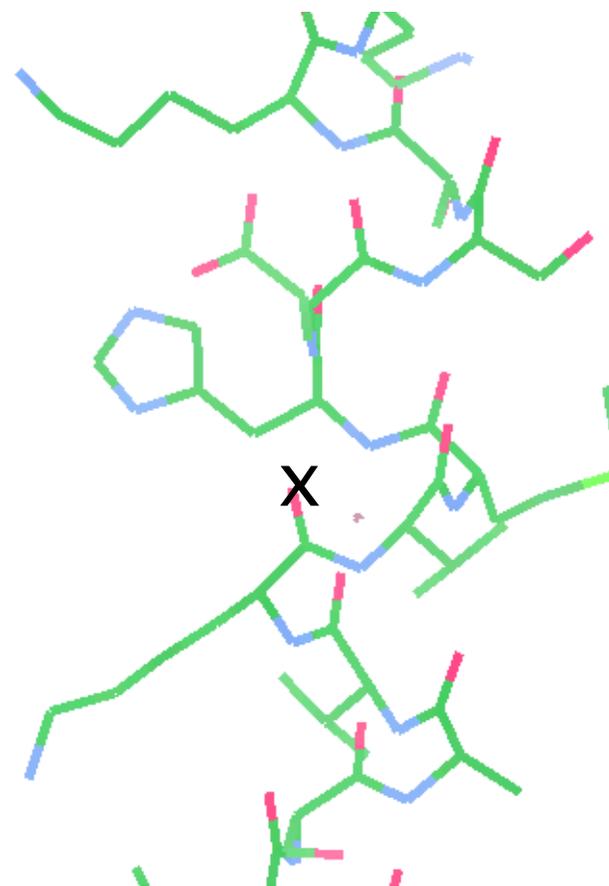
ATOM	5	CB	ASP	A	8	-30.909	9.723	18.264	1.00	33.70	C
ATOM	6	CG	ASP	A	8	-31.252	9.345	16.825	1.00	41.96	C
ATOM	7	OD1	ASP	A	8	-31.072	10.248	15.981	1.00	46.18	O

Model Parameterisation

Standard refinable parameters

Atomic model:

- Position – (x,y,z) coordinates
- Uncertainty – B-factors
- (Occupancies)



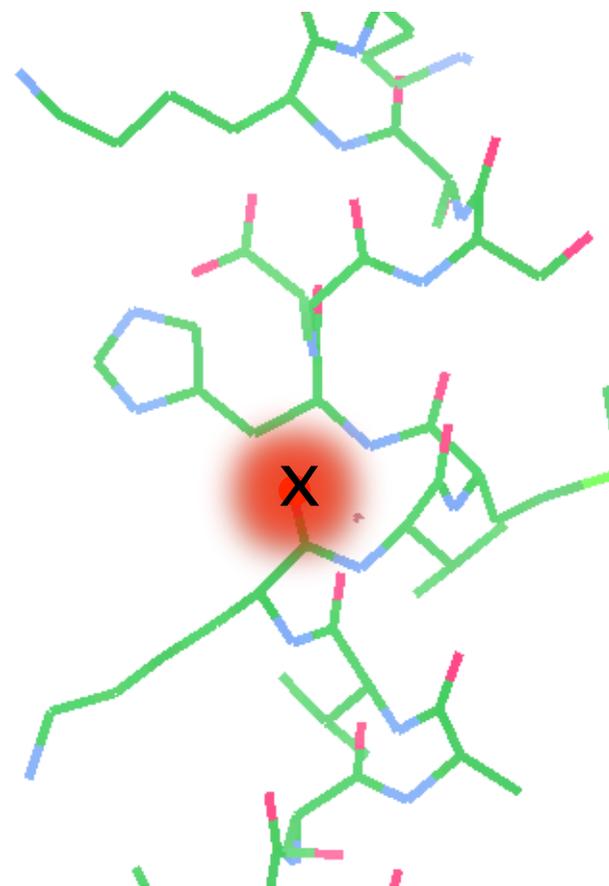
ATOM	5	CB	ASP	A	8	-30.909	9.723	18.264	1.00	33.70	C
ATOM	6	CG	ASP	A	8	-31.252	9.345	16.825	1.00	41.96	C
ATOM	7	OD1	ASP	A	8	-31.072	10.248	15.981	1.00	46.18	O

Model Parameterisation

Standard refinable parameters

Atomic model:

- Position – (x,y,z) coordinates
- Uncertainty – B-factors
- (Occupancies)



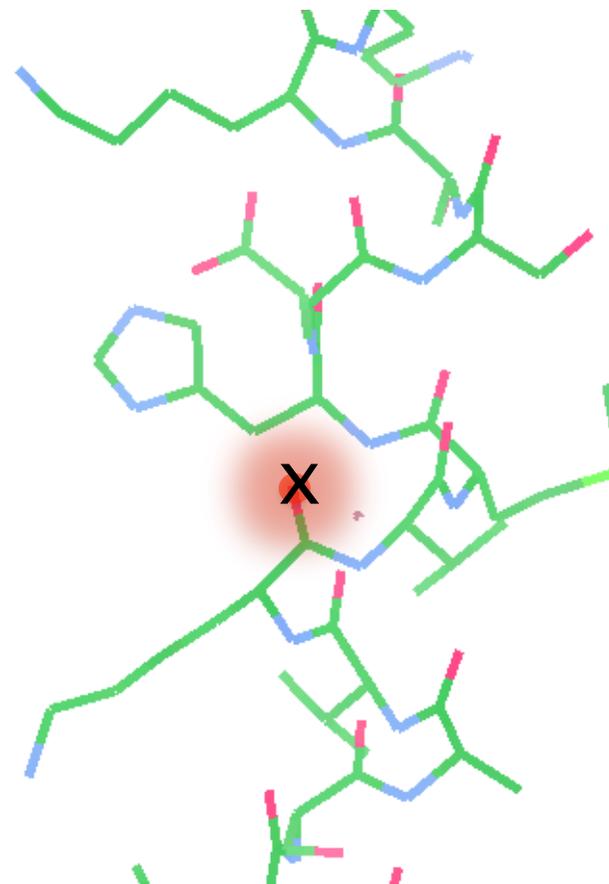
ATOM	5	CB	ASP	A	8	-30.909	9.723	18.264	1.00	33.70	C
ATOM	6	CG	ASP	A	8	-31.252	9.345	16.825	1.00	41.96	C
ATOM	7	OD1	ASP	A	8	-31.072	10.248	15.981	1.00	46.18	O

Model Parameterisation

Standard refinable parameters

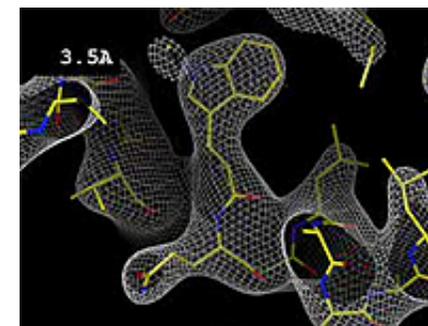
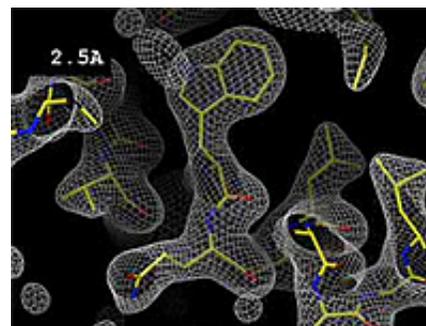
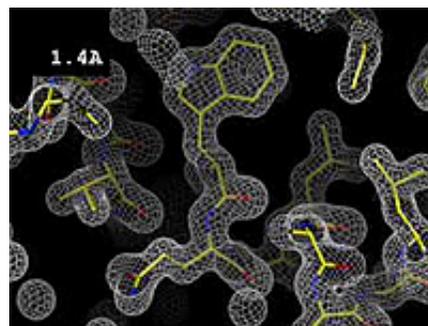
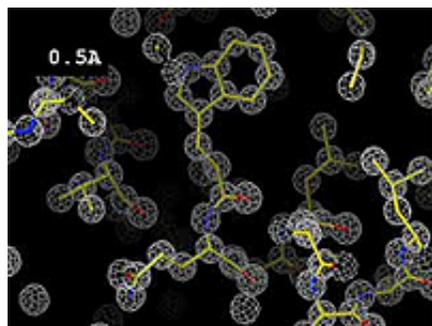
Atomic model:

- Position – (x,y,z) coordinates
- Uncertainty – B-factors
- (Occupancies)



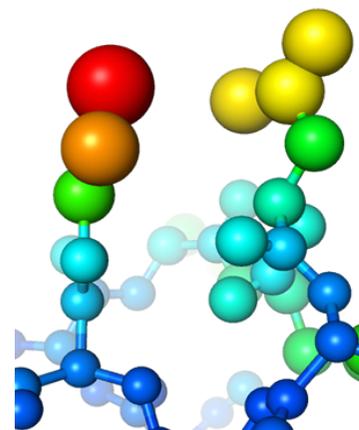
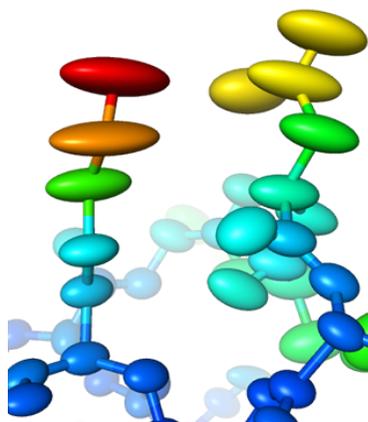
ATOM	5	CB	ASP	A	8	-30.909	9.723	18.264	1.00	33.70	C
ATOM	6	CG	ASP	A	8	-31.252	9.345	16.825	1.00	41.96	C
ATOM	7	OD1	ASP	A	8	-31.072	10.248	15.981	1.00	46.18	O

Model Parameterisation

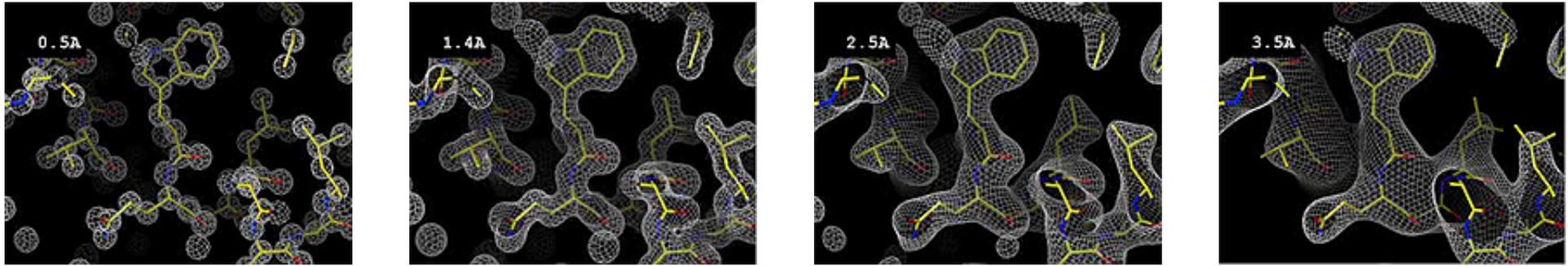


Notes on B-factors

- Atomic B-factors can be modelled in a different way according to the quality of the data: **anisotropic** (6 parameters per atom), **isotropic** (1 parameter per atom), **TLS** (20 parameters per group of atoms)
- B-factors describe relative positional uncertainty



Model Parameterisation



Notes on B-factors

- Atomic B-factors can be modelled in a different way according to the quality of the data: **anisotropic** (6 parameters per atom), **isotropic** (1 parameter per atom), **TLS** (20 parameters per group of atoms)
- B-factors describe relative positional uncertainty
- B-factors are sometimes also referred to as atomic displacement parameters (ADPs) or thermal/temperature factors
- Should not compare atomic B-factors between different models

TLS Groups

Describe rigid body motion – e.g. for chains/domains/subunits

Suitable for medium resolution, when full anisotropy is impossible

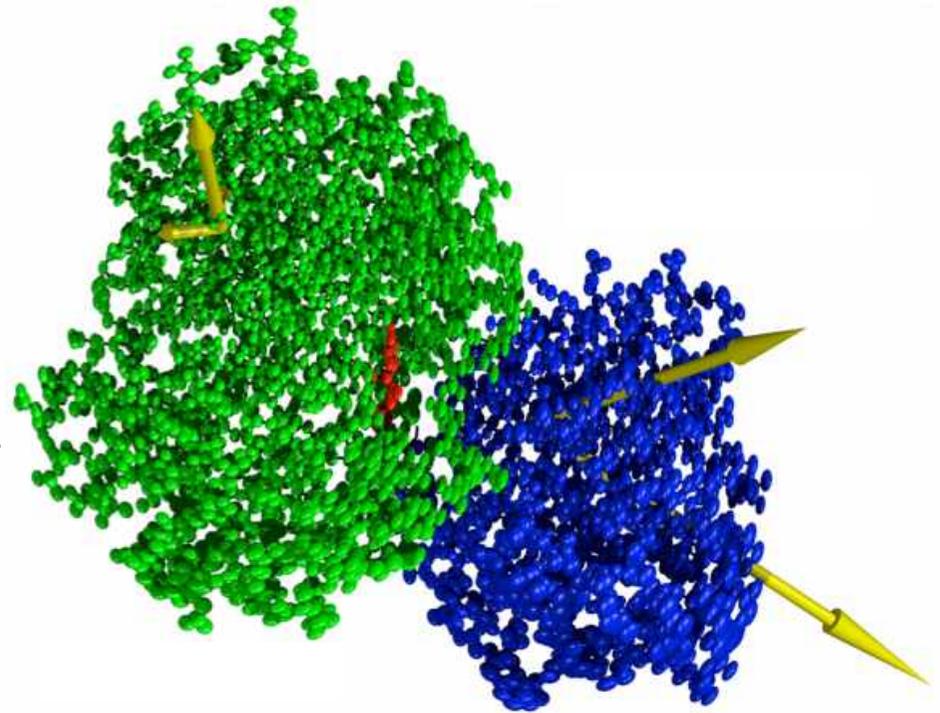
Per group (20 parameters):

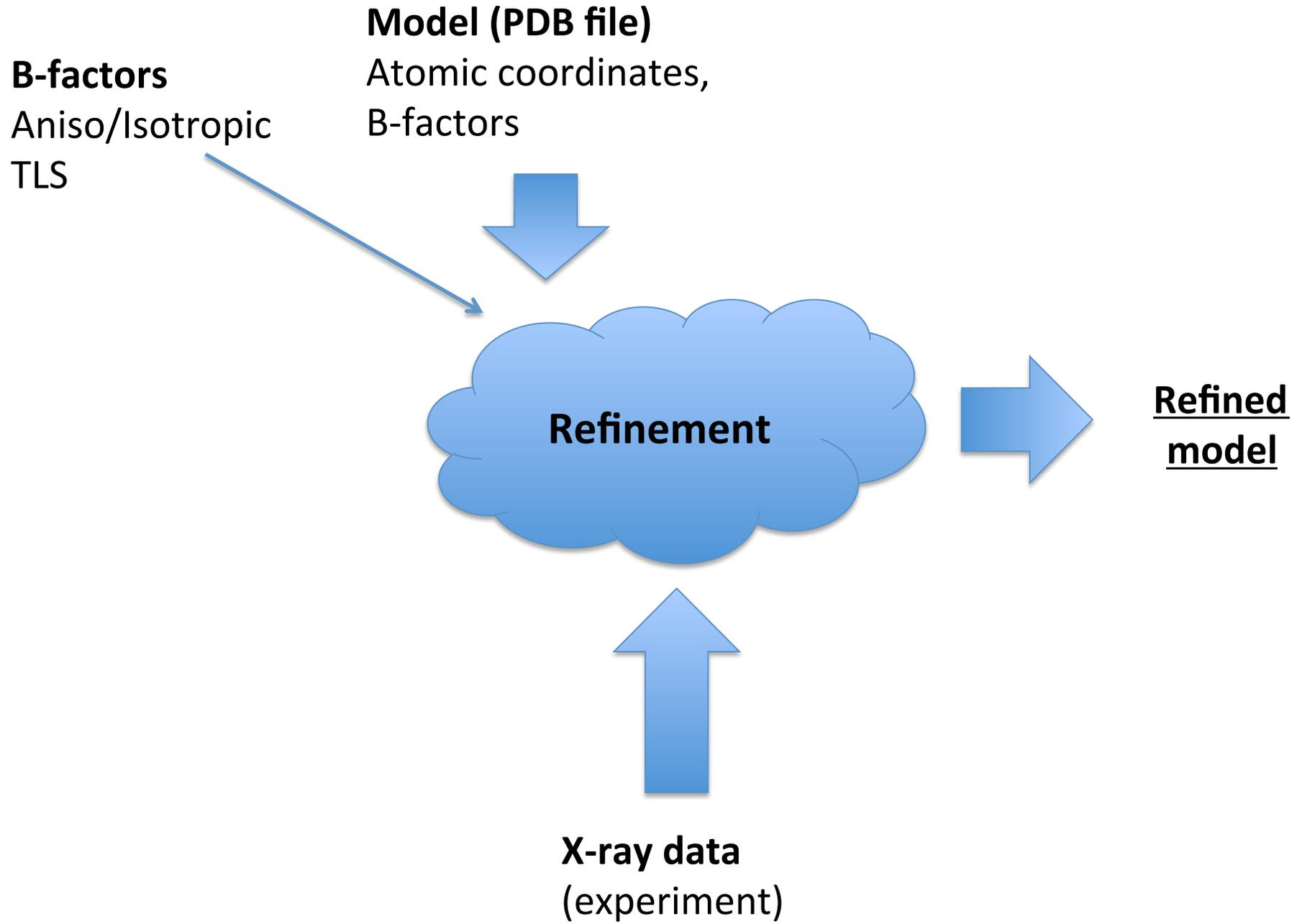
- Translation – 6 parameters
- Libration – 6 parameters
- Screw rotation – 8 parameters

Define groups using CCP4i

or TLSMD webserver:

<http://skuld.bmsc.washington.edu/~tlsmd/>



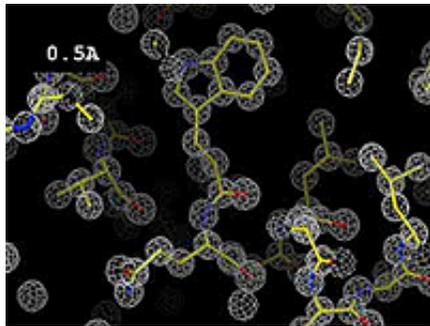


Model Refinement

Refinement strategy will differ for different quality of original data and it can also exploit particular features of your crystal:

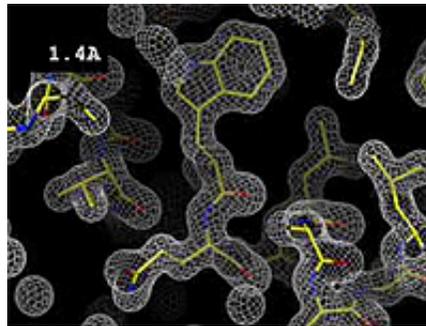


1.0Å



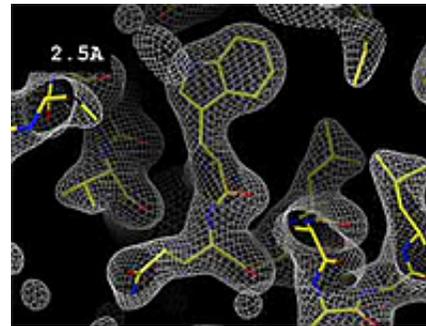
Unrestrained?

2.0Å



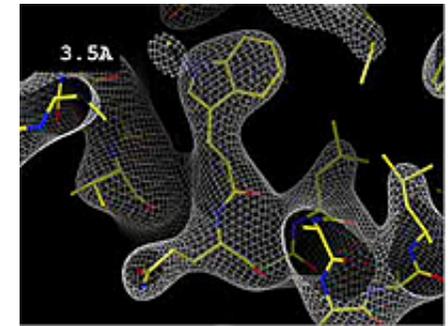
Chemical restraints

3.0Å



Chemical restraints

4.0Å

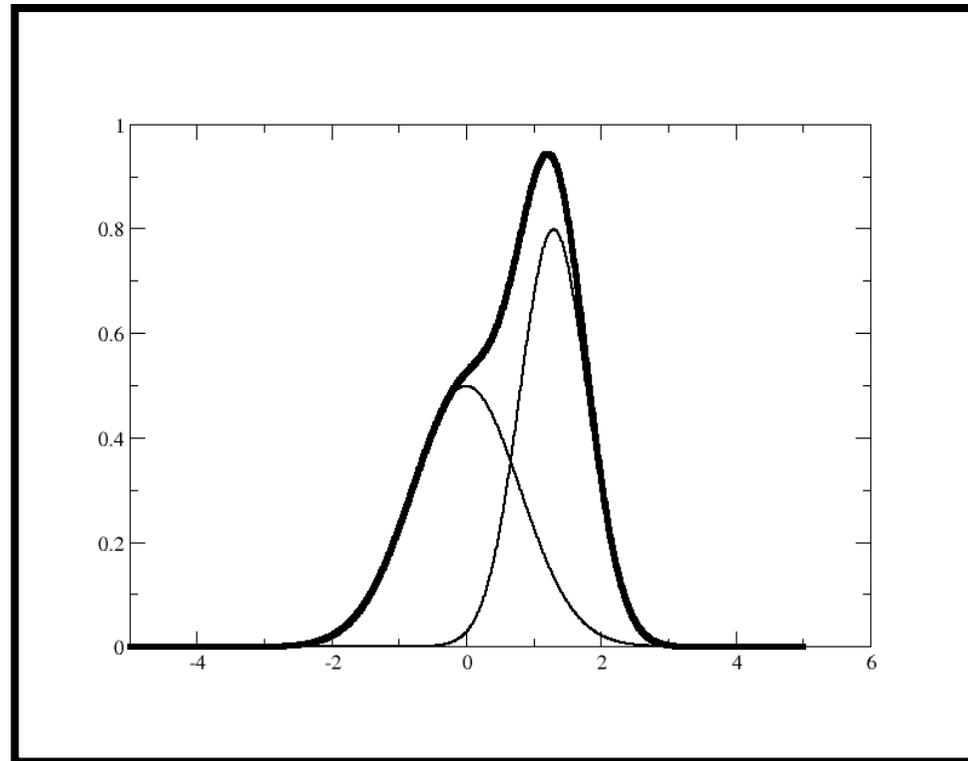


Chemical restraints
Jelly-body restraints
External restraints

Why Restraints?

Example: two-atom ideal case

Distance between atoms 1.3Å. B-factors 20 and 50



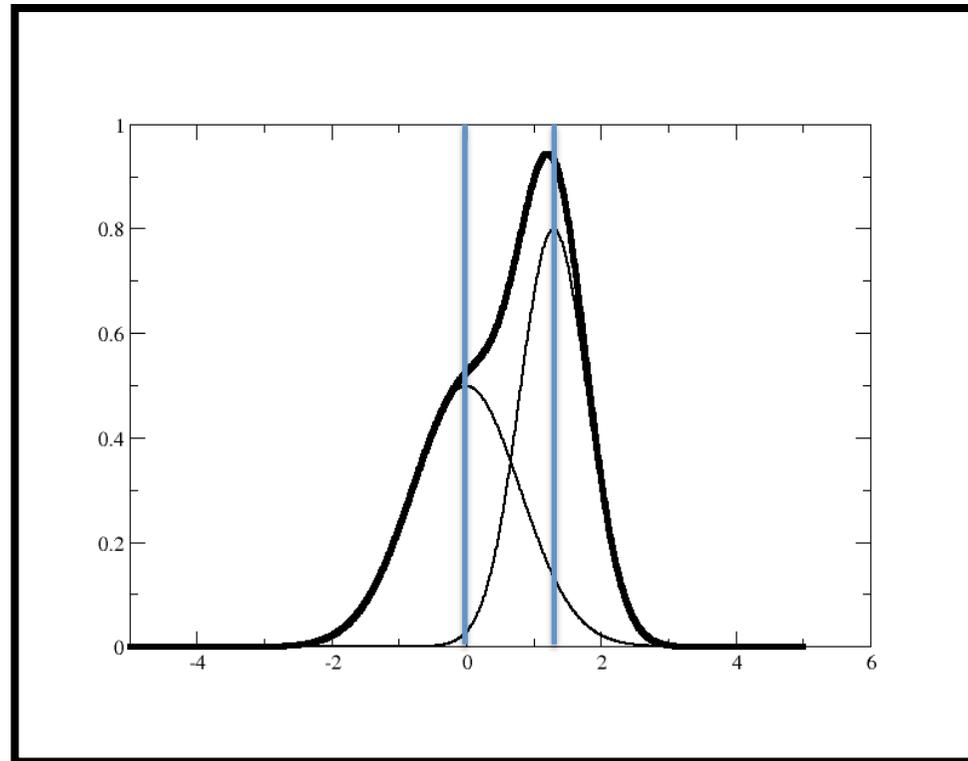
Thin lines – single atoms

Bold line – sum of the two atoms

Why Restraints?

Example: two-atom ideal case

Distance between atoms 1.3Å. B-factors 20 and 50



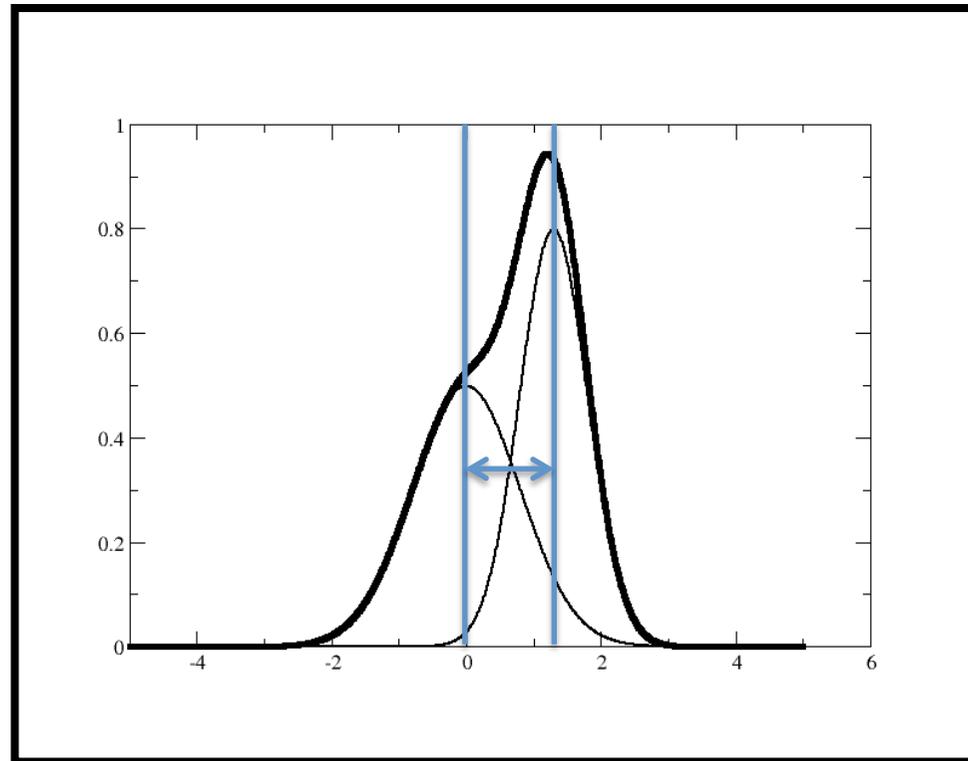
Thin lines – single atoms

Bold line – sum of the two atoms

Why Restraints?

Example: two-atom ideal case

Distance between atoms 1.3Å. B-factors 20 and 50



Thin lines – single atoms

Bold line – sum of the two atoms

Restraints

Standard restraints (used by default) include:

- Bond lengths
- Angles
- Chirals
- Planes
- Some torsion angles
- B-values
- VDW repulsions

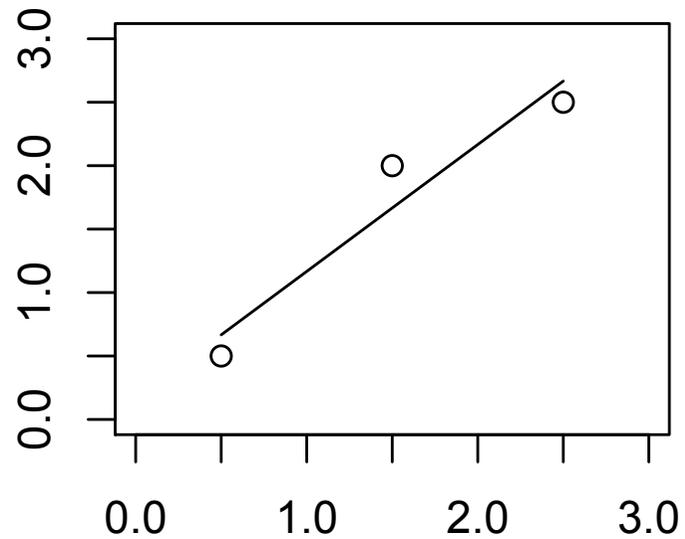
These help to ensure that the model is chemically sensible

Note – we generally deal with restraints, not constraints

Restrains

Why introduce so many restraints?

Answer: to improve the observation:parameter ratio.



Example: Fitting a line $y = a + bx$

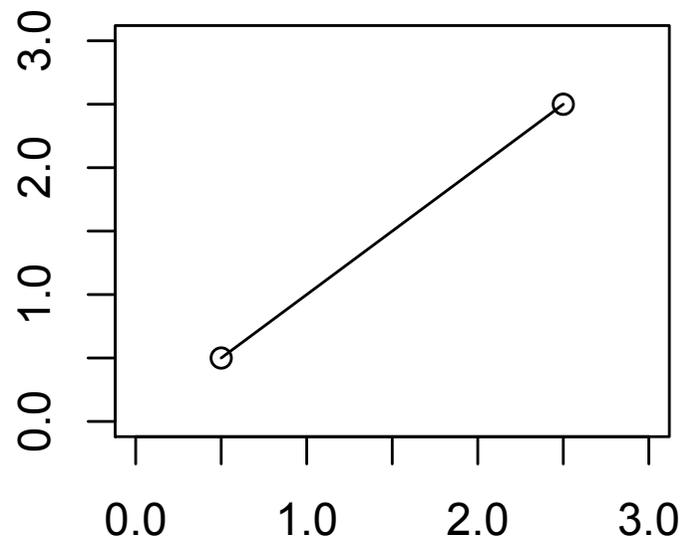
Restrains

Why introduce so many restraints?

Answer: to improve the observation:parameter ratio.

Can fit a line

Line is unreliable



Overfitting
Model Bias

Example: Fitting a line

$$y = a + bx$$

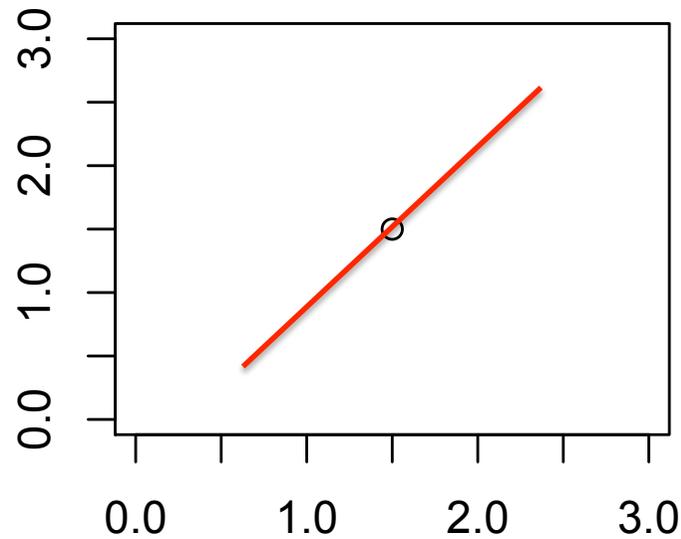
Restraints

Why introduce so many restraints?

Answer: to improve the observation:parameter ratio.

Insufficient observations!

Unstable refinement



Ill-posed problem

Example: Fitting a line

$$y = a + bx$$

Restraints

How to improve the observation:parameter ratio.

1. Reduce number of parameters

N_{obs} is resolution dependent...

High resolution : Anisotropic B-factors – 9 params per atom

Low resolution : Isotropic B-factors – 4 params per atom

- Rigid body refinement – 9 parameters per body

Restraints

How to improve the observation:parameter ratio.

1. Reduce number of parameters
2. Increase number of restraints
 - B-value restraints
 - NCS restraints
 - H-bond and secondary-structure restraints
 - Restraints to homologous known structures
 - Nucleic acid base-pair and base-stacking restraints
 - Jelly-body restraints

NCS

(Non-Crystallographic Symmetry Restraints)

1. NCS constraints

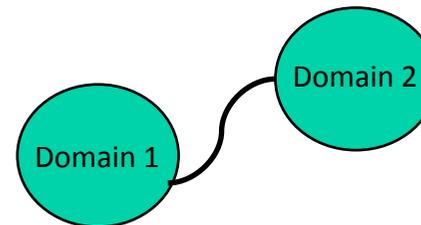
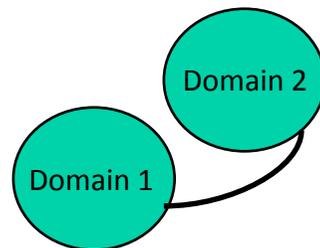
- NCS-related copies are considered to be exactly the same
- Only one set of atomic parameters per molecule

2. Global NCS restraints

- Molecules are superimposed
- Difference between atoms are minimised

3. Local NCS restraints

- Molecules are assumed to be locally similar
- However, they may adopt (slightly) different global conformations
- Restrain differences between local interatomic distances



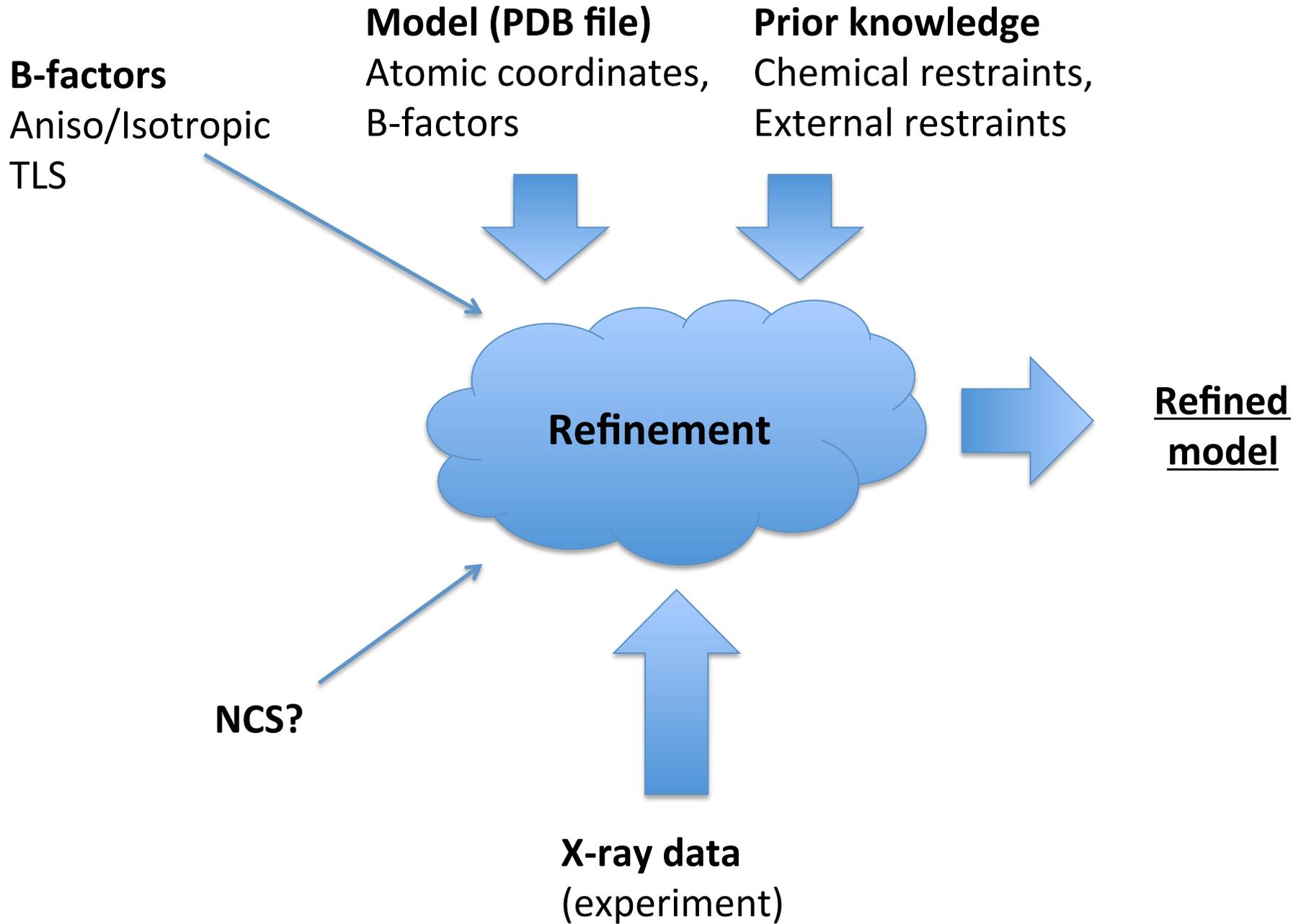
Ligand Refinement

Geometric restraints for protein / nucleic acids are pre-tabulated

Ligands are more complicated

Need a source of prior information

- Common/known structures are dealt with automatically
 - CCP4/REFMAC monomer library has pre-computed descriptions
- New ligands require description (CIF file)
 - New tool - ACEDRG



$$f_{\text{xray}} = -\log[P(\text{obs};\text{model})]$$

likelihood of the data

$$f_{\text{geom}} = -\log[P(\text{model})]$$

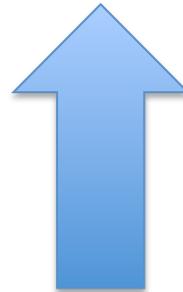
probability of the model

w : relative weighting

Prior knowledge
Chemical restraints,
External restraints



$$f_{\text{tot}} = wf_{\text{xray}} + f_{\text{geom}}$$



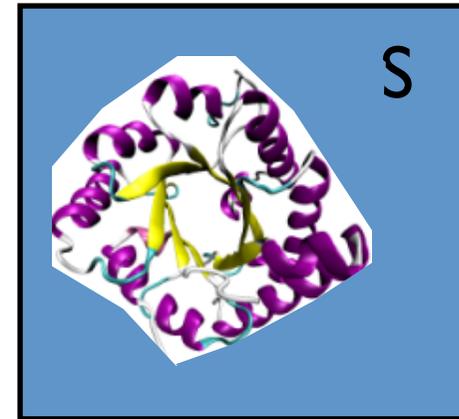
X-ray data
(experiment)

Solvent model

Mask-based bulk solvent correction

Idea:

- Protein region is masked out
- Solvent region is flattened (set to constant)
- Structure factors for solvent are calculated: F_{solvent}



$$F_{\text{total}} = F_{\text{protein}} + \alpha F_{\text{solvent}}$$

α : solvent scale factor

$$k e^{-Bs^2} e^{-s^TUs} (F_{\text{protein}} + \alpha F_{\text{solvent}})$$

k : overall scale factor

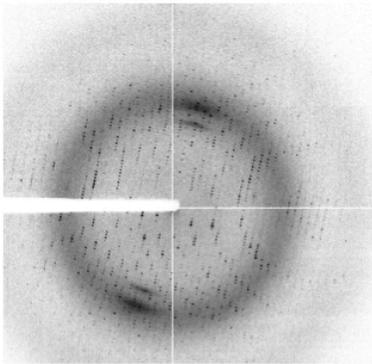
B : overall B-factor

U : overall anisotropic B-factor

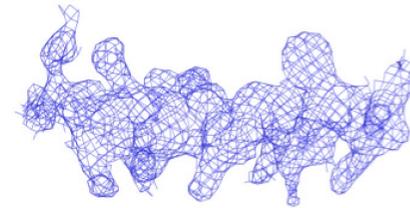
Overall Parameters: Scaling

Problem:

- Observed and calculated amplitudes need to be brought to the same scale so that they can be compared



H, K, L	$ F_{obs} $	ϕ
...		
5, 5, 5	348	-
5, 5, 6	392	-
5, 5, 7	157	-
5, 5, 8	312	-
...		



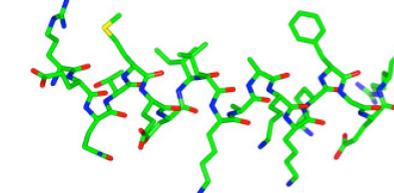
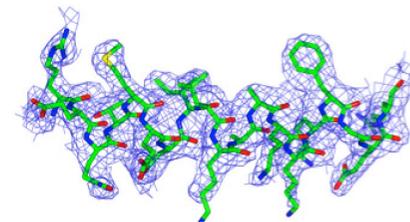
Idea:

Iteratively improve the model, optimising the agreement between

$|F_{obs}|$ and $|F_{calc}|$

Purpose: improve phase estimates: ϕ_{calc}

H, K, L	$ F_{calc} $	ϕ_{calc}
...		
5, 5, 5	355	27°
5, 5, 6	387	8°
5, 5, 7	146	75°
5, 5, 8	340	31°
...		



Overall Parameters: Scaling

Problem:

- Observed and calculated amplitudes need to be brought to the same scale so that they can be compared



Overall Parameters: Scaling

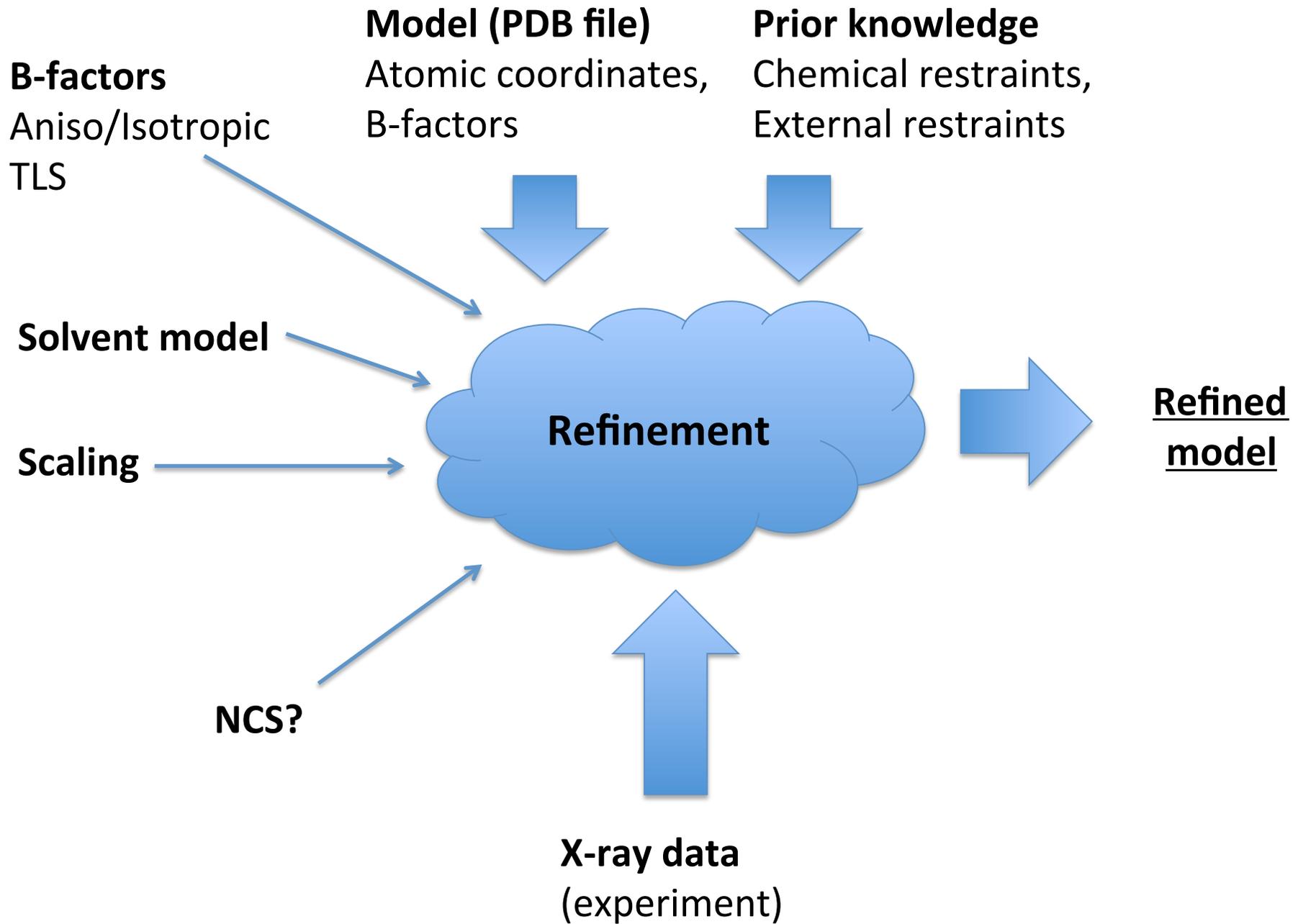
Problem:

- Observed and calculated amplitudes need to be brought to the same scale so that they can be compared

Need to:

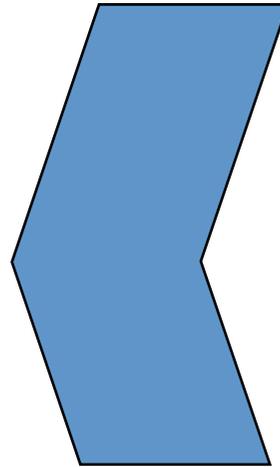
- Modify/scale F_{calc}
- Find a scaling function, with some parameters – “overall parameters”

Scale parameters are optimised in ML refinement, along with all other parameters



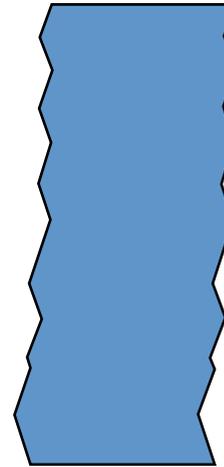
Twin Refinement

twin



yes

polysynthetic
twin



no

A single crystal can
be cut out of the twin:

Need to deal with polysynthetic twin during refinement

Twin Refinement

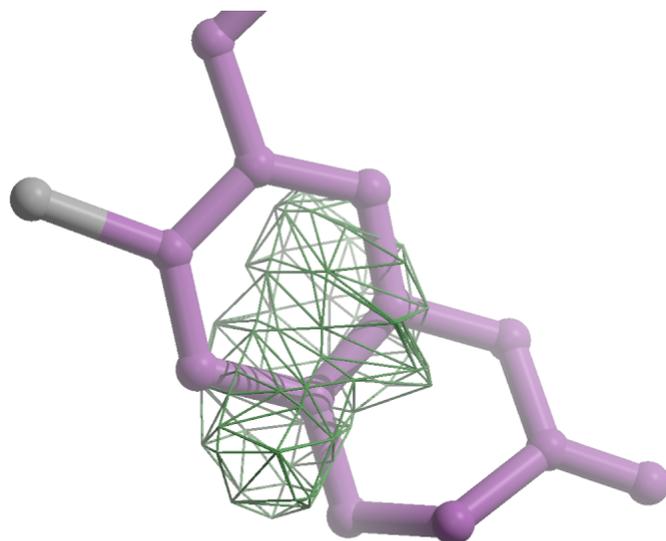
Twin refinement in REFMAC5 is automatic

1. Identify potential twin operators
2. For each operator, calculate R_{merge} (R-factor comparing twin-related intensities)
3. If $R_{\text{merge}} > 0.44$ remove this operator
4. Refine twin fractions
5. Keep only sufficiently large domains (default 7%)

DON'T USE TWIN REFINEMENT IF YOUR R_{free} IS HIGHER THAN 40%!

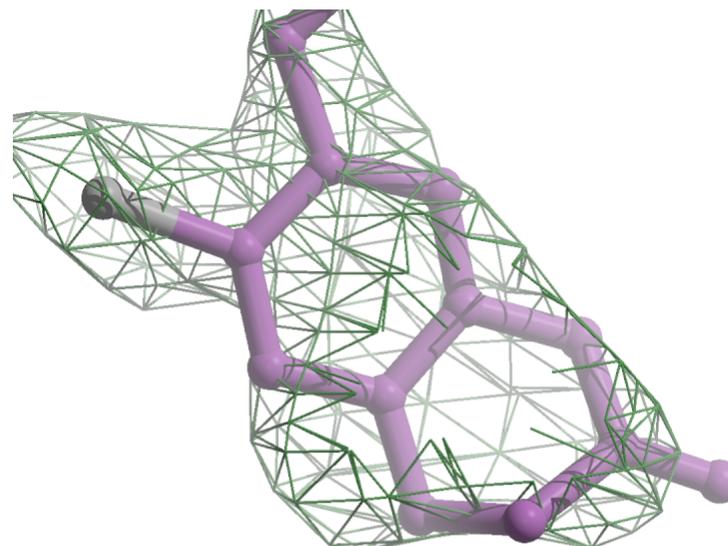
Twin Refinement

Example: Where's the density for my ligand (2.15Å)?



$$R/R_{\text{free}} = 25.5/26.9\%$$

After initial rigid body and restrained refinement



$F_o - F_c$ (3σ)

$$R/R_{\text{free}} = 15.9/16.3\%$$

Re-refine with twin on
(twin fractions: 0.6/0.4)

Borrowed from Ben Bax, ex-GSK

Stages of model building

Early stages (e.g. straight after MR)

- Run many cycles (up to 200) of refinement

Medium stages – during model building

- 10–20 cycles
- Optimise additional parameters like NCS, TLS, etc
- Once your Rfree is lower than 40%, turn on twinning if twinned
- Play with geometry weight parameter

Final stages of refinement

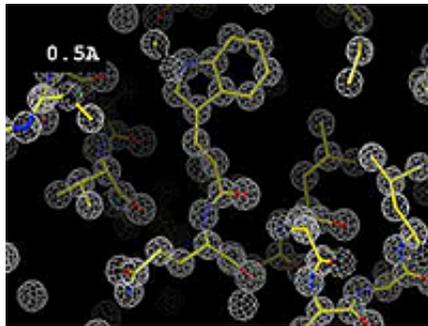
- Even fewer cycles?
- Use optimal additional parameters
- If geometrical quality of the model is not optimal, further play with parameters like geometry weight.

Model Refinement

Refinement strategy will differ for different quality of original data and it can also exploit particular features of your crystal:

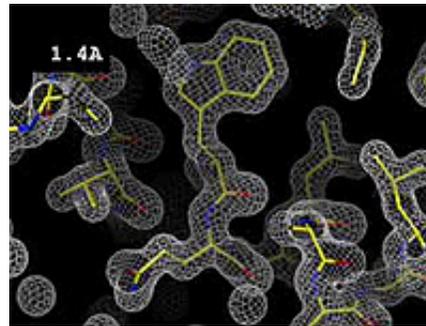


1.0Å



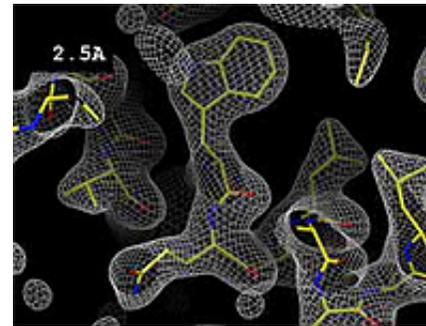
Unrestrained
Anisotropic B
Hydrogens

2.0Å



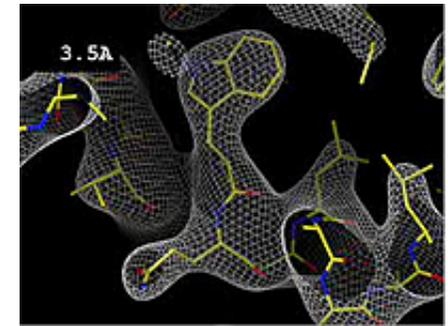
Chemical restraints
Isotropic B
TLS

3.0Å



Chemical restraints
Isotropic B
TLS

4.0Å



Chemical restraints
Jelly-body restraints
External restraints
Isotropic B

- Twinning
- Several similar subunits – NCS (at moderate and low resolution)
- Available homologues for external restraint generation (low resolution)

