# Chemical Crystallography and Structural Chemistry
(VO 270287)
Lecture 8
28th May 2020

Dr. Tim Grüne

Centre for X-ray Structure Analysis

Faculty of Chemistry

University of Vienna

tim.gruene@univie.ac.at

# Previous Lecture

1. Model building

2. Refinement

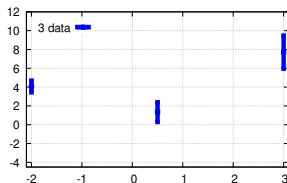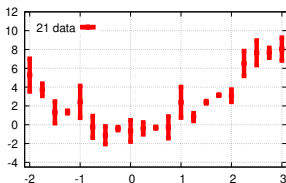3. Constraints & Restraints

# Today's Lecture

1. Example of constraints and data:parameter ratio

2. Validation

# Example: Stabilisation through restraints

Two hypothetic measurements:

Experiment 1: high resolution, 21 pairs of measurements $(x_1, y_1), \ldots, (x_{21}, y_{21})$ and errors $\sigma_1, \ldots, \sigma_{21}$

Experiment 2: low resolution, 3 pairs of measurements $(x_1, y_1), \ldots, (x_3, y_3)$ and errors $\sigma_1, \ldots, \sigma_3$

# Example: Stabilisation through restraints

Testing two models:

**Model 1:** $g(x) = g_2 x^2 + g_1 x + g_0$

**Model 2:** $h(x) = h_3 x^3 + h_1 x + h_0$

Either model has three parameters, $g_0, g_1, g_2$ and $h_0, h_1, h_3$ respectively. These parameters correspond to *e.g.* the model coordinates $(x_i, y_i, z_i)$, or the ADPs $U_i$.

We will fit both models to the data to find out which model better represents the data.

# Example: Stabilisation through restraints

Least-squares-minimisation:

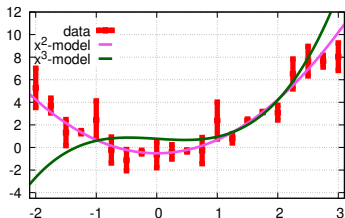$$\text{minimise:} \sum_{i=1}^{N} \frac{1}{\sigma_i^2} (y_i - g(x_i))^2 \qquad \text{model 1}$$

$$\text{minimise:} \sum_{i=1}^{N} \frac{1}{\sigma_i^2} (y_i - h(x_i))^2 \qquad \text{model 2}$$

- Experiment 1: $N = 21$ data points

- Experiment 2: $N = 3$ data points

We will start with the high resolution experiment 1

# Example: Stabilisation through restraints

experiment 1: high resolution; high data to parameter ratio = 21:3=7



Model 1: $1.2x^2 + 0.0x - 0.5$
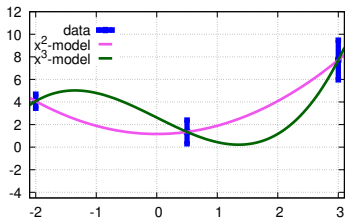rmsd: 1.07
Model 2: $0.5x^3 - 0.3x - 0.8$
rmsd: 4.74

The *root mean square deviation* rmsd between model and data corresponds to the crystallographic $R1$ value.

The lower rmsd 1.07 clearly favours model 1. The pink curve also visually fits the data better than the green curve.

# Example: Stabilisation through restraints

experiment 2: low resolution, low data to parameter ration = 3:3 = 1



model 1: $0.7x^2 + 0.0x + 1.2$
rmsd: 0
model 2: $0.5x^3 - 2.7x - 2.6$
rmsd: 0

When there are as many parameters as data points, any model can be fitted
perfectly to the data. We cannot distinguish between the two models

# Example: Stabilisation through restraints

experiment 2: low resolution with constraint

For some reason we know that the data must pass through the point $(0,0)$. For the two models this means

$$0 = g(0)$$
$$= g_2 * 0^2 + g_1 * 0 + g_0$$
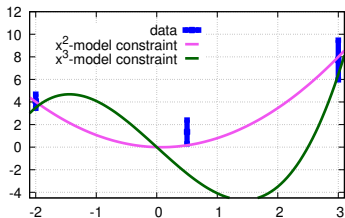$$\Rightarrow g_0 = 0$$

and analogously

$$h_0 = 0$$

The constraint reduced the number of parameters, only two parameters per model

# Example: Stabilisation through restraints

experiment 2: low resolution with constraint



model 1: $0.9x^2 - 0.1x$
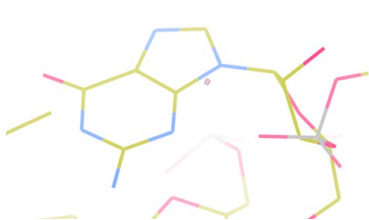rmsd: 1.13
model 2: $0.8x^3 - 4.9x$
rmsd: 3.7

Due to the constraint, data to parameter ratio = 3:2 = 1.5. Now there is an *rmsd*, and it favours (again) the first model.
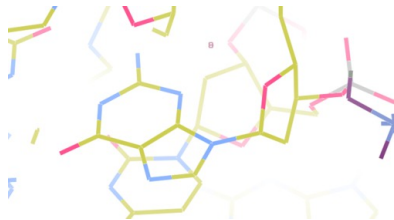
# Summary & model building

- (*cf.* phase problem)

- phases are calculate from the model

- model phases and observed data yield the electron density map, and electron difference map

- model building improves the model in large steps

- refinement optimises the model against the data

- medium resolution data or poor quality data require restraints and constraints in order to create a chemically sensible model

**Model quality and data quality: structure validation**

# Atom coordinates $\neq$ model accuracy



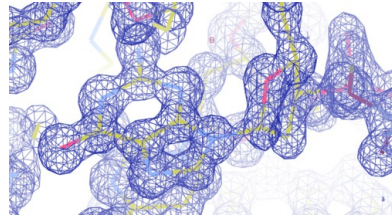Guanine model in ribosome, data resolution 3.1Å

Guanine model in Z-DNA, at resolution 1.0 Å

The coordinates of the model do no reveal the data quality, nor the model quality.

# model coordinates = interpretation of data



Guanine model **with map** in ribosome, data resolution 3.1Å



Guanine model **with map** in Z-DNA, at resolution 1.0 Å

Only in combination with the data can we judge the model quality

# Once more: data to parameter ratio

Example Ciprofloxacin ($a = 9.5$Å, $b = 9.9$Å, $c = 11.0$Å, $\alpha = 94.2°$, $\beta = 100.2°$, $\gamma = 91.3°$)

- $FC_{17}N_3O_9H_{30}$: $60 \times 9 = 540$ Parameter

**data resolution 0.43 Å:** 26'308 reflections $\hat{=}$ 48.7 data points per parameter: very high, reliable refinement

**data resolution 0.8 Å:** 2'926 reflections $\hat{=}$ 5.4 data points per parameter: medium, refinement needs checking

# Once more: data to parameter ratio

Example Ribosome ($a = 401.4$Å, $b = 401.4$Å, $c = 175.9$Å, $\alpha = \beta = \gamma = 90°$, $P4_12_12$)
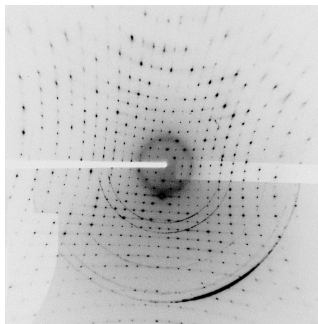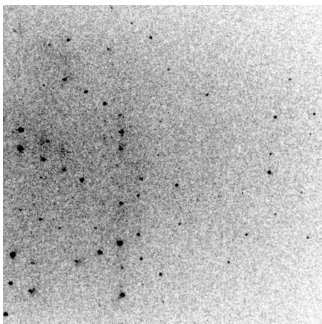
- PDB ID 1J5E: 51'atoms atoms = 207'768 parameters

- data resolution 3.05 Å 238'205 reflections

$$\frac{238'205}{207'768} = 1.15$$

Even at such low data to parameter ratio can a reasonable model be built and refined. It is important to be aware of differences in the interpretation of the data

# Quality indicators

# Example data quality

# Important quality indicators

**$R_{meas}$** relative difference between:

>   1. measured data
>   2. calculated data

**data completeness** : fraction of measured data w.r.t. theoretically possible data

**multiplicity** (*alias*: *redundancy*): how often every unique reflection was measured (on average)

**signal strength** $I(hkl)/\sigma_{I(hkl)} < 1$: noise

**$CC_{1/2}$**   1. split data set into two random halves
    2. calculated correlation coefficient between symmetry equivalent reflections

## R-values for data

The classic data quality indicator is $R_{int}$, alias $R_{merge}$ or $R_{sym}$:

$$R_{int} = \sum_h \sum_j \frac{|I_{hj} - \langle I_h \rangle|}{\langle I_h \rangle}$$

$R_{int}$ mathematically increases with multiplicity, although data quality improves with multiplicity

$R_{int}$ is typically shown in publications. It is, however, obsolete and should not be published. $R_{meas}$ *alias* $R_{r.i.m.}$ should be published instead:

$$R_{meas} = \sum_h \frac{n_h}{n_h - 1} \sum_j \frac{|I_{hj} - \langle I_h \rangle|}{\langle I_h \rangle}$$

# Example data statistics (`XPREP`)

```
  Resolution    #Data #Theory %Complete Redundancy Mean I Mean I/s Rmerge
  Inf - 2.46      196     197     99.5      39.27    215.01  110.27  0.0300
 2.46 - 1.13     1762    1825     96.5      14.86     75.32   42.01  0.0453
 1.13 - 0.89     1972    2123     92.9       8.71     25.52   19.00  0.0895
 0.89 - 0.77     2007    2258     88.9       6.81     10.84   10.39  0.1425
 0.77 - 0.69     1864    2499     74.6       3.37      5.66    5.76  0.1885
 0.69 - 0.62     2108    3360     62.7       2.24      2.88    3.29  0.2890
 0.62 - 0.57     1929    3542     54.5       1.44      1.51    1.79  0.4191
 0.57 - 0.54     1123    2367     47.4       1.10      0.90    1.14  0.5593
 ----------------------------------------------------------------------
 0.64 - 0.54     3720    7014     53.0       1.43      1.47    1.76  0.4170
  Inf - 0.54    12961   18171     71.3       5.08     20.64   13.61  0.0514

Merged [A], lowest resolution = 11.49 Angstroms
```
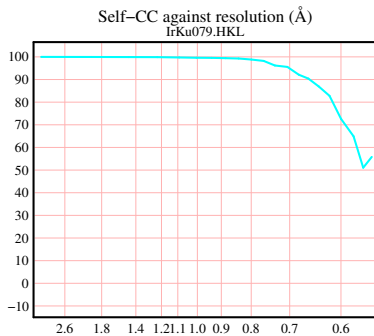
# CC1/2, and resolution cut-off

- CC1/2 should be close to 100% throughout resolution range

- where CC1/2 drops below 70%, noise becomes significant, and data at higher resolution can be excluded from refinement

- $I/\sigma(I)$ should be about 2, where CC1/2 about 70%

- $I/\sigma(I)$ should be about 1, where CC1/2 about 40% (in cases very resolution cut-off is critical)

# Example CC$_{1/2}$, and resolution cut-off



CC1/2 *vs.* data resolution

# R-values for the model

$$R = R1 = \sum_h \frac{||F_h(data)| - |F_h(model)||}{|F_h(data)|}$$

weighted R-value:

$$wR = \sum_h \frac{|w_h|F_h(data)| - |F_h(model)||}{w_h|F_h(data)|}$$

weighted intensity based R-value:

$$wR2 = R_B = \sqrt{\sum_h \frac{|w_h(I_h(data) - I_h(model))^2|}{w|I_h(data)|^2}}$$

**Small molecules**: $R1$ of the refined model 2-5 %.

universität
wien

## `GooF`

*Goodness of Fit*

$$GooF = \sqrt{\frac{\sum_h w_h \left(F_h{}^2(data) - F_h{}^2(model)\right)^2}{n - p}}$$

- Takes number of parameters ($p$) and number of data ($n$) into account

- Ideally $\approx 1$, increases with worse model

# model: residual density

SHELXL calculates the "highest peak" and "deepest hole" in the electron
density map. Units are electrons, e.g. at the beginning of model building:

```
Electron density synthesis with coefficients Fo-Fc


Highest peak    4.95  at  0.5434  0.9981  0.3231  [  0.04 A from RU01
Deepest hole  -3.34  at  0.0057  0.4976  0.3299  [  0.99 A from RU02
~~~~~~

Mean =    0.00,   Rms deviation from mean =    0.34 e/A^3
~~~~~~~~~~~~~
```

## model: residual density

SHELXL calculates the "highest peak" and "deepest hole" in the electron density map. Units are electrons, e.g. for the refined model:

```
Electron density synthesis with coefficients Fo-Fc

Highest peak    0.50  at  0.6610  0.1969  0.4278  [  0.69 A from C006
Deepest hole   -1.22  at  0.2635  0.6156  0.2132  [  0.04 A from P003
^^^^^^

Mean =    0.00,   Rms deviation from mean =     0.06 e/A^3
^^^^^^^^^^^^
```

# checkCIF (PLATON web-based)



Every published structure *should* have a checkCIF report. There are different alert levels of decreasing severity. Reviewers typically require that a structure should **not** contain A- or B-alerts.

## Summary

- A model without data does not reflect data quality

- Data quality: data resolution, multiplicity, R-values, $I/\sigma_I$, $CC_{1/2}$

- Model quality: R1-values, GooF, residual density

- available for everyone: checkCIF `http://checkcif.iucr.org` (with or without data)

- *ALERT levels* A, B, …

- (Analogously for macromolecular structures: `http://molprobity.biochem.duke.edu/`)

# End of lecture