

Strukturelle Modellierung  
(Masterstudiengang Bioinformatik)

## Strukturbestimmung mit Röntgenkristallographie

Sommersemester 2013

Peter Güntert

# Introduction

## Myoglobin Struktur



"Vielleicht die bemerkenswerteste Eigenschaft des Moleküls ist seine Komplexität und die Abwesenheit von Symmetrie. Der Anordnung scheinen die Regelmässigkeiten, die man instinktiv erwartet, fast völlig zu fehlen, und sie ist komplizierter als von irgendeiner Theorie der Proteinstruktur vorhergesagt." — John Kendrew, 1958

## Kristallographie: Geschichte

1839, William H. Miller: Miller Indices für Gitterebenen  
1891: 230 Raumgruppen für Kristalle  
1895, Wilhelm Conrad Röntgen: Röntgenstrahlung  
1912, Max von Laue: Röntgenstreuung  
1912, William L. Bragg: Braggsches Gesetz  
1914, Bragg: Kristallstrukturen von NaCl und Diamant  
1937: Dorothy Hodgkin: Kristallstruktur von Cholesterin  
1945: Dorothy Hodgkin: Kristallstruktur von Vitamin B12  
1952: Rosalind Franklin: DNA Röntgenbeugungsdiagramme  
1955: Rosalind Franklin: Tabakmosaikvirus (TMV) Struktur  
1958: John Kendrew: Erste Proteinstruktur (Myoglobin)  
2000: Kristallstruktur des Ribosoms  
2013: > 79'000 Kristallstrukturen in der Protein Data Bank

## Literatur über Kristallstrukturbestimmung

- B. Rupp, *Biomolecular Crystallography*, Garland, 2010.
- W. Massa, *Kristallstrukturbestimmung*, Teubner, 5<sup>2007</sup>.
- C. Branden & J. Tooze, *Introduction to Protein Structure*, Garland, <sup>2</sup>1999.

## Crystallographic structure models versus proteins in solution

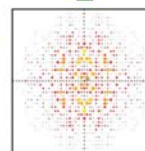
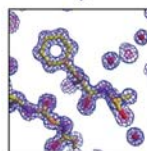
- Protein crystals are formed by a loose periodic network of weak, non-covalent interactions.
- Protein crystals contain large solvent channels. The solvent channels allow relatively free diffusion of small molecules through the crystal and also provide conformational freedom for surface-exposed side chains or loops.
- The core structure of protein molecules in solution as determined by NMR is identical to the crystal structure. Even enzymes generally maintain activity in protein crystals.
- Crystal packing can affect local regions of the structure where surface-exposed side chains or flexible surface loops form intermolecular crystal contacts.
- Large conformational movements destroy crystals and cannot be directly observed through a single crystal structure. Limited information about the dynamic behavior of molecules can be obtained from analysis of the *B*-factors as a measure of local displacement.
- The quality of a protein structure is a local property. Surface-exposed residues or mobile loops may not be traceable in electron density, no matter how well defined the rest of the structure is.

### Challenges of protein crystallography

- Proteins are generally difficult to crystallize and without crystals there is no crystallography. Preparing the material and modifying the protein by protein engineering so that it can actually crystallize is nontrivial.
- Prevention of radiation damage by ionizing X-ray radiation requires cryocooling of crystals and many crystals are difficult to flash-cool.
- The X-ray diffraction patterns do not provide a direct image of the molecular structure. The electron density of the scattering molecular structure must be reconstructed by Fourier transform techniques.
- Both structure factor amplitude and relative phase angle of each reflection are required for the Fourier reconstruction. While the structure factor amplitudes are readily accessible being proportional to the square root of the measured reflection intensities, the relative phase angles must be supplied by additional phasing experiments. The absence of directly accessible phases constitutes the phase problem in crystallography.
- The nonlinear refinement of the structure model is nontrivial and prior stereochemical knowledge must generally be incorporated into the restrained refinement.

### The crystallographic phase problem

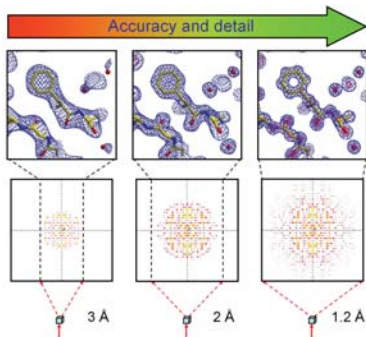
$$\rho(x, y, z) = \frac{1}{V} \sum_{-h}^h \sum_{-k}^k \sum_{-l}^l F_{hkl} \exp[-2\pi i(hx + ky + lz - \alpha_{hkl})]$$



The crystallographic phase problem

In order to reconstruct the electron density of the molecule, two quantities need to be provided for each reflection (data point): the structure factor amplitude,  $F_{hkl}$ , which is directly obtained through the experiment and is proportional to the square root of the measured intensity of the diffraction spot or reflection; and the phase angle of each reflection,  $\alpha_{hkl}$ , which is not directly observable and must be supplied by additional phasing experiments.

### Data quality determines structural detail and accuracy



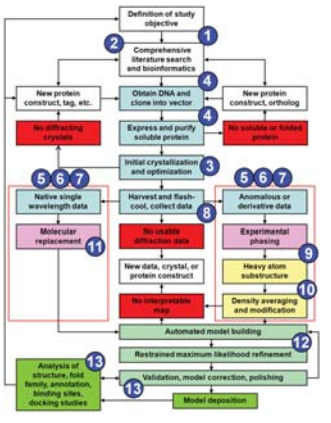
The qualitative relation between the extent of X-ray diffraction, the resulting amount of available diffraction data, and the quality and detail of the electron density reconstruction and protein structure model are evident from this figure: The crystals are labeled with the nominal resolution  $d_{min}$ , given in Å (Angstrom) and determined by the highest diffraction angle (corresponding to the closest sampling distance in the crystal, thus  $d_{min}$ ) at which X-ray reflections are observed. Above each crystal is a sketch of the corresponding diffraction pattern, which contains significantly more data at higher resolution, corresponding to a smaller distance between discernable objects of approximately  $d_{min}$ . As a consequence, both the reconstruction of the electron density (blue grid) and the resulting structure model (stick model) are much more detailed and accurate.

### Kristallstrukturbestimmung

1. Proteinherstellung
2. Kristallisation
3. Messung der Beugungsmuster
4. Datenauswertung
  - a) Bestimmung der Einheitszelle und Raumgruppe
  - b) Phasenbestimmung
  - c) Modellbau
  - d) Verfeinerung der Phasen und der Struktur

### Key stages in X-ray structure determination

The flow diagram provides an overview about the major steps in a structure determination project, labeled with the chapter numbers treating the subject or related general fundamentals. Blue shaded boxes indicate experimental laboratory work, while all steps past data collection are conducted *in silico*.



### Crystallographic computer programs

- Protein crystallography depends heavily on computational methods.
- Crystallographic computing has made substantial progress, largely as a result of abundant and cheap high performance computing.
- It is now possible to determine and analyze complex crystal structures entirely on inexpensive laptop or desktop computers with a few GB of memory. Automation and user interfaces have reached a high level of sophistication (although compatibility and integration issues remain).
- As a result, the actual process of structure solution, although the theoretically most sophisticated part in a structure determination, is commonly not considered a bottleneck in routine structure determination projects.
- Given reliable data of decent resolution (~2.5 Å or better) and no overly large or complex molecules, many structures can in fact be solved de novo and refined (although probably not completely polished) within several hours.
- Automated model building programs—many of them available as web services—have removed much of the tedium of initial model building.

## Key concepts of protein crystallography I

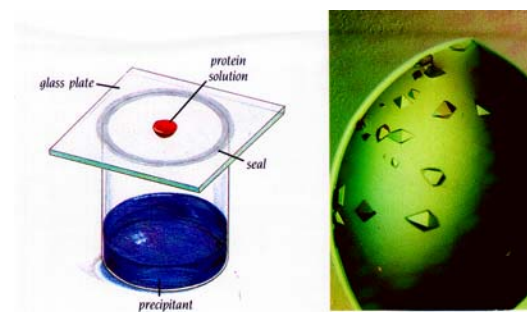
- The power of macromolecular crystallography lies in the fact that highly accurate models of large molecular structures and molecular complexes can be determined at often near atomic level of detail.
- Crystallographic structure models have provided insight into molecular form and function, and provide the basis for structural biology and structure guided drug discovery.
- Non-proprietary protein structure models are made available to the public by deposition in the Protein Data Bank, which holds more than 83 000 entries as of April 2013.
- Proteins are generally difficult to crystallize; without crystals there is no crystallography.
- Preparing the material and modifying the protein by protein engineering so that it can actually crystallize is nontrivial.
- Radiation damage by ionizing X-ray radiation requires cryocooling of crystals, and many crystals are difficult to flash-cool.

## Key concepts of protein crystallography II

- The X-ray diffraction patterns are not a direct image of the molecular structure.
- The electron density of the scattering molecular structure must be reconstructed by Fourier transform techniques.
- Both structure factor amplitude and relative phase angle of each reflection are required for the Fourier reconstruction.
- While the structure factor amplitudes are readily accessible, being proportional to the square root of the measured reflection intensities, the relative phase angles must be supplied by additional phasing experiments.
- The absence of directly accessible phases constitutes the phase problem in crystallography.
- The nonlinear refinement of the structure model is nontrivial and prior stereochemical knowledge must generally be incorporated into the restrained refinement.

# Crystallization

## Proteinkristallisation



## Protein crystallization basics

- Protein crystals are periodic self-assemblies of large and often flexible macromolecules, held together by weak intermolecular interactions. Protein crystals are generally fragile and sensitive to environmental changes.
- In order to form crystals, the protein solution must become supersaturated. In the supersaturated, thermodynamically metastable state, nucleation can occur and crystals may form while the solution equilibrates.
- The most common technique for protein crystal growth is by vapor diffusion, where water vapor equilibrates from a drop containing protein and a precipitant into a larger reservoir with higher precipitant concentration.
- Given the large size and inherent flexibility of most protein molecules combined with the complex nature of their intermolecular interactions, crystal formation is an inherently unlikely process, and many trials may be necessary to obtain well-diffracting crystals.

## The protein is the most crucial factor in determining crystallization success

- Given that a crystal can only form if specific interactions between molecules can occur in an orderly fashion, the inherent properties of the protein itself are the primary factors determining whether crystallization can occur.
- A single-residue mutation can make all the difference between successful crystallization and complete failure.
- Important factors related to the protein that influence crystallization are its purity, the homogeneity of its conformational state, the freshness of the protein, and the additional components that are invariably present, but often unknown or unspecified, in the protein stock solution.

### Hanging drop vapor diffusion

Mix cocktail and protein on glass slide

Turn slide and seal well

Observe for crystal formation

Well with crystallization cocktail (precipitants, additives, detergents, etc. – unlimited combinations possible)

Vapor diffuses into well, concentrations in drop increase

Harvest and mount crystals

© Garland Science 2010

### Solubility phase diagram

Protein concentration ↑

Precipitant concentration →

Pure water ↑

solubility line

decomposition line

clear protein solution single phase stable

metastable solution will eventually separate into protein (maybe in form of crystals) plus saturated solution

precipitate + protein solution two-phase region unstable, spontaneous decomposition

Figure 3-7 A basic solubility phase diagram for a given temperature. The diagram visualizes the general observation that the higher the precipitant concentration in the solution, the lower the maximal achievable protein concentration in the solution and vice versa. Between the solubility line and the decomposition line lies the metastable region representing the supersaturated protein solution, which will eventually—given the necessary kinetic nucleation events—equilibrate and separate into a protein-rich phase (often in the form of precipitate or crystals) and saturated protein solution.

© Garland Science 2010

### Protein solubility versus pH

Protein solubility ↑

pH of solution

Isoelectric point

Figure 3-8 Protein solubility versus pH of protein solution. The protein shown in this example has its solubility minimum at its isoelectric point of ~6.3, where the sum of positive and negative charges (the net charge of the protein) is zero. Even at the isoelectric point, there are still numerous (but net compensating) local charges present on the surface of the protein.

© Garland Science 2010

### Crystal growth

11.2 x 11.2  $\mu\text{m}^2$

27 x 27  $\mu\text{m}^2$

880 x 880 nm<sup>2</sup>

1 x 1  $\mu\text{m}^2$

Figure 3-11 Atomic force microscope images of crystal growth. (Panel A) The atomic force microscope images of the 001 surface of glucose isomerase show the two most common growth patterns observed in crystal growth: step growth starting from 2-dimensional nucleation islands (A, left image) and a double-spiral growth pattern (A, right image). Panel B shows formation of supercritical 2-dimensional nuclei on the 001 surface of cytomegalovirus (CMV), a member of the herpes virus family. As indicated by the arrows, in this case only two virions (B, left image) suffice to generate a critical nucleus from which new step growth commences (B, right image). Images courtesy of Alexander McPherson and Aaron Greenwood, University of California, Irvine.

© Garland Science 2010

### Mosaic crystals

Figure 3-12 Growth of a real mosaic crystal. The schematic drawing shows a crystal growing in a solution of protein molecules (blue spheres). Small impurities (red) and some larger detritus (green squares) are also present in the solution. New molecules attach preferentially to steps and edges (red arrows) and we can recognize a growth defect in the form of a hole; impurities are enclosed at the domain boundaries, and a larger piece of detritus is incorporated at a domain boundary. Individual domains can be substantially misaligned, in this case about 6°; such a highly mosaic crystal would not be useful for diffraction experiments.

© Garland Science 2010

### Crystallization techniques

- The inability to predict *ab initio* any conditions favoring protein crystallization means that, in general, several hundred crystallization trials must be set up in a suitable format and design.
- Crystallization screening experiments are commonly set up manually or robotically in multi-well format crystallization plates.
- The most common procedure for achieving supersaturation is the vapor-diffusion technique, performed in sitting-drop or hanging-drop format. In vapor-diffusion setups, protein is mixed with a precipitant cocktail, and the system is closed over a reservoir into which water vapor diffuses from the protein solution. During vapor diffusion, both precipitant and protein concentration increase in the crystallization drop and supersaturation is achieved.
- As a rule of thumb, low supersaturation favors controlled crystal growth, while high supersaturation is required for spontaneous nucleation of crystallization nuclei. Seeding is a method to induce heterogeneous nucleation at low supersaturation, which is more conducive to controlled crystal growth.

## Robot for automated crystallization

**Figure 3-33 Automated crystallization setup for the small laboratory.** Based on the assumption of modest throughput requirements, and no necessity for full walk-away automation, two low-budget approaches to automation are conceivable: selection of a single system that can prepare crystallization cocktails (perhaps in a limited fashion) and also set up the crystallization plates,<sup>10</sup> or a dual-stage layout using separate cocktail preparation with a generic liquid-handling system followed by a dedicated plate-setup robot.<sup>11</sup> The major reason for separating plate

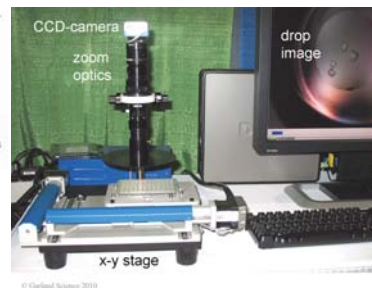


setup from cocktail production is differing requirements for dispensing precision, volume, and speed. First, small volume ( $\mu\text{l}$  to  $\text{nl}$ ), and very accurate (also in geometric terms) dispensing is mandatory for plate crystallization setups, whereas large volume ( $\text{ml}$ ) handling with modest speed and precision requirements suffices for cocktail production. Another advantage of the separation between the cocktail stage and the plate setup is that simple one-to-one dispensing into reservoir wells and drop aliquots followed by protein addition with a single needle dispenser suffices (Figure 3-33) once the cocktails are produced in a 96-well format deep-well block. Deep-well blocks prefilled with crystallization cocktails are also commercially available. In addition, compared with a single-stage setup, failure of one system component does not affect the others. For example, cocktail production can continue while the plate setup robot is inoperative. Figure 3-33 shows a popular robot for 96-well crystallization plate setup.

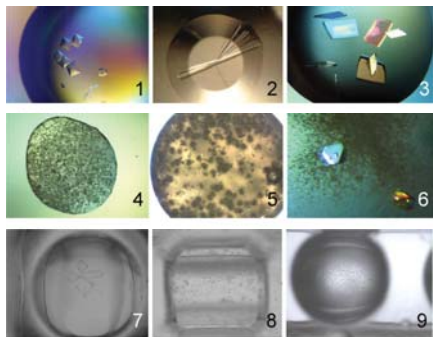
**Figure 3-33 A robot for automated crystallization plate setup.** The Phoenix robot (Art Robbins Instruments) can set up 96 crystallization trials in about one minute. On the left side, a 96-channel syringe dispenser re-arrays 100  $\mu\text{l}$  each 96 prefabricated or purchased crystallization cocktails simultaneously from a standard deep-well block into the reservoirs of an 895-format, 96-well sitting-drop crystallization plate, and places between 1  $\mu\text{l}$  and 100  $\text{nl}$  into the drop shelves or wells. From the right side, a contactless microvalve dispenser nozzle immediately adds the pre-egested protein (stock walk in the well block) rapidly and without contact onto each of the precipitant drops. To minimize evaporation, the plate is then immediately sealed with a sheet of pressure-sensitive adhesive. Taking all losses into account, about 12 to 15  $\mu\text{l}$  of protein stock is required for 96 100  $\times$  100  $\text{nl}$  drops. The robot design has been based on a prototype developed in an academic laboratory setting.<sup>10</sup>

## Crystallization plate imaging

**Figure 3-36 A low-cost automated crystallization plate imaging station.** The crystallization plate is positioned by an x-y translation stage, and a digital zoom camera takes high-resolution images of the crystallization drops. The images taken in about 2 minutes can then be manually inspected on a computer screen, or processed by automated image recognition software. The depicted instrument is the CrysCam microscope manufactured by Art Robbins Instruments.

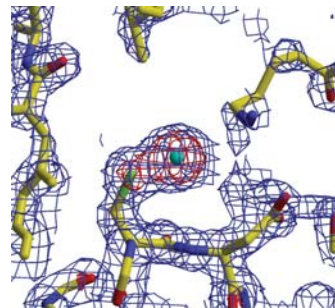


## Crystallization outcomes



**Figure 3-37 Images of crystallization drops with different experimental outcomes.** (1) Perfectly formed octahedral single crystals in hanging drop. (2) Cluster of large needles in microbatch reservoir. (3) Two plates in a hanging drop. (4) Microbatch reservoir in hanging drop. (5) Small crystals in a hanging drop with a microbatch reservoir. (6) Single crystals in a sitting drop. (7) Single crystals in a sitting drop. (8) Single crystals in a sitting drop with a microbatch reservoir. (9) Single crystals in a sitting drop with a microbatch reservoir.

## Heavy atom derivatives



**Figure 3-42 Heavy atom derivatization of a protein.** Shown is the electron density around a gold atom covalently linked to a cysteine residue in the Clostridium tetani neurotoxin.<sup>123</sup> A combination of anomalous and isomorphous signals from gold atoms were used to solve the structure of the ganglioside binding domain of the neurotoxin from bacillus C. tetani, the causative agent of tetanus infections. PDB entry 1a8d

## Heavy atom reagents

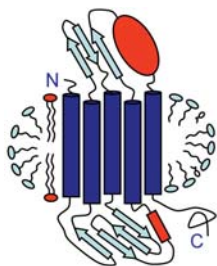
Name	Formula
Platinum potassium chloride, potassium tetrachloroplatinate(II)	$\text{K}_2\text{PtCl}_6$
Aurous potassium cyanide, potassium dicyanoaurate(III)	$\text{KAu}(\text{CN})_2$
Mercuric potassium iodide, potassium tetraiodo mercurate(II)	$\text{K}_2\text{HgI}_4$
Uranyl acetate, uranium(VI) oxyacetate	$\text{UO}_2(\text{C}_2\text{H}_3\text{O}_2)_2$
Mercuric(II) chloride	$\text{HgCl}_2$
Potassium uranyl fluoride, potassium uranium(VII) oxyfluoride	$\text{K}_2\text{UO}_7\text{F}_6$
Para-chloromercuribenzenesulfonate, PCMBs	$\text{Hg}(\text{C}_6\text{H}_4\text{SO}_3)_2$
Trimethyllead acetate	$(\text{CH}_3)_3\text{Pb}(\text{CH}_3\text{COO})_2$
Methylmercuric acetate	$\text{CH}_3\text{Hg}(\text{CH}_3\text{COO})_2$
Ethylmercuric thioisocyanate, thiomersal	$\text{C}_2\text{H}_5\text{HgSC}_2\text{H}_4\text{N}_2$
Hexatantalum tetrabromide	$(\text{Ta}_6\text{Br}_{14})\text{Br}$

**Table 3-1 Selected heavy atom reagents.** The listed reagents are frequently used for derivatization. The top seven entries are historically the most well used, the alkylated compounds below and the powerful Ta-clusters are more recent and very successful derivatization reagents. Many more are listed in the heavy atom data bank<sup>124</sup> and in the review by M.A. Rould<sup>125</sup>. All these substances are quite toxic when ingested because they bind to proteins and taking corresponding precautions is prudent. The uranium salts are generally prepared from natural uranium (0.7%  $^{235}\text{U}$  or depleted uranium ( $^{235}\text{U}$ ), which both are only a weak  $\alpha$ -particle source.

## Less than 1% of all deposited protein structures are membrane protein structures

- About a third of all expressed human proteins are presumed to be membrane proteins, and over 60% of all current drug targets are membrane receptors. Their primary functions include transport of material and signals across cell membranes as well as motor functions.
- Despite membrane proteins being a significant class of proteins, it was nearly 30 years, and 195 deposited protein structures, after Kendrew's first myoglobin structure in 1958 that the first integral membrane protein structure, the photosynthetic reaction center isolated from the bacterium *Rhodospirillum rubrum*, was published in 1985. That research led to a Nobel Prize for crystallographic work being awarded to Johann Deisenhofer, Hartmut Michel, and Robert Huber in 1988.
- In early 2007, there were 242 coordinate entries of 122 different membrane proteins out of 35100 total entries in the PDB, still a factor of 1/145 disfavoring the membrane proteins. Clearly, membrane protein crystallization remains a major challenge for crystallography.

### Resolubilized membrane protein



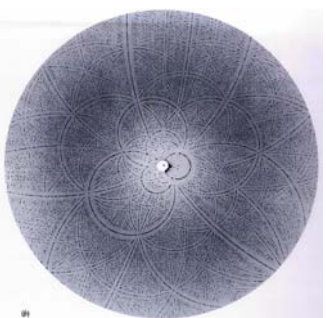
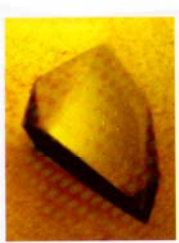
- Membrane phospholipid
- Detergent
- Amphiphile

**Figure 3-43** Resolubilized multi-pass, polytopic transmembrane protein with its associated detergent collar. In addition to the detergent collar, membrane fragments are often associated and co-solubilized with the transmembrane stem, as sketched on the left side of the membrane collar. Small amphiphile molecules are often added to fine-tune the size of the membrane collar for subsequent crystallization, as shown at the right side of the membrane collar.

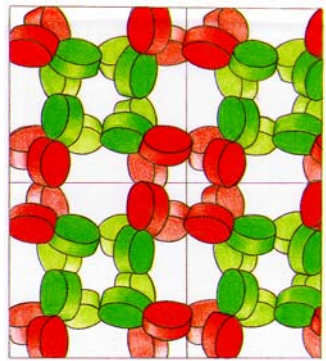
© Garland Science 2008

# Crystals

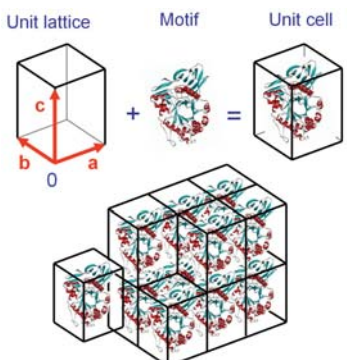
### Kristall und Beugungsmuster



### Proteinkristall



### Unit lattice + Motif = Unit cell



**Figure 5-24** Assembly of a primitive triclinic 3-dimensional crystal from unit cells. In analogy to the 2-dimensional case, the unit lattice is filled with a motif, and the crystal is built from translationally stacked unit cells. The basis vectors form a right-handed system  $[0, \mathbf{a}, \mathbf{b}, \mathbf{c}]$ .

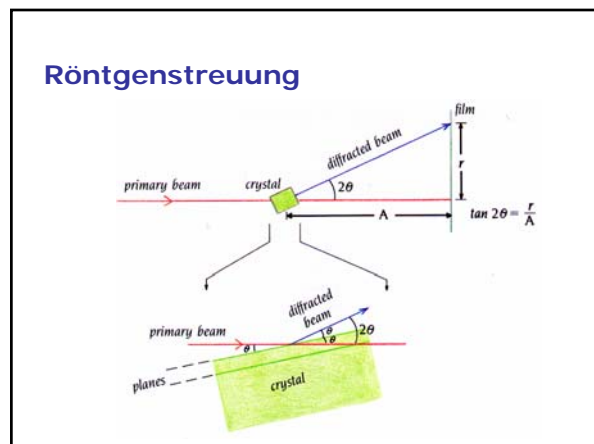
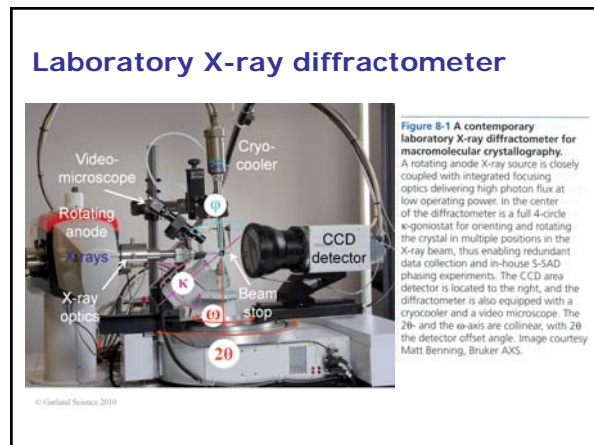
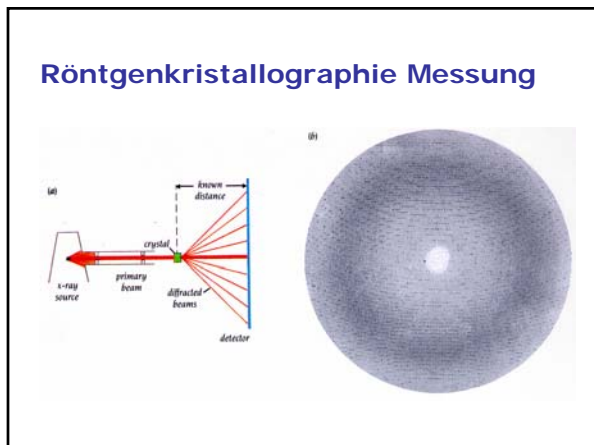
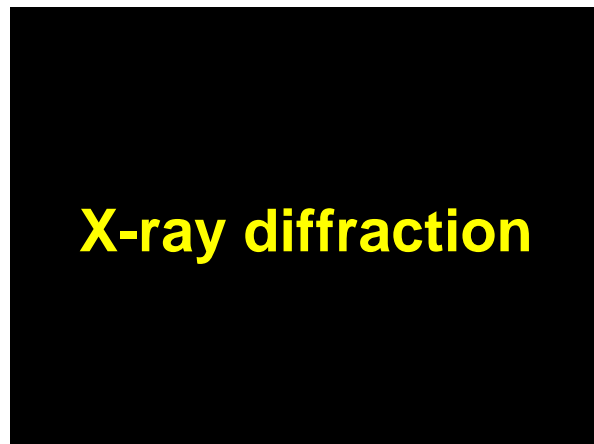
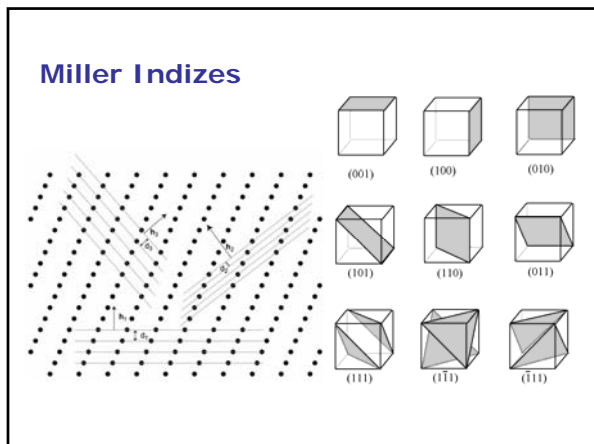
### Unit cell parameters

The three basis vectors of a unit lattice  $[0, \mathbf{a}, \mathbf{b}, \mathbf{c}]$  extend from a common origin in a right-handed system; that is, if going counterclockwise from basis vector  $\mathbf{a}$  to basis vector  $\mathbf{b}$ , the third basis vector  $\mathbf{c}$  points upwards (Figure 5-25). The vector product  $\mathbf{a} \times \mathbf{b}$  generates a third vector  $\mathbf{c}$  perpendicular to  $\mathbf{a}$  and  $\mathbf{b}$ , and the vector product  $\mathbf{a} \times \mathbf{b}$  is *positive defined* in a right-handed system. The magnitude of this vector,  $|\mathbf{a} \times \mathbf{b}|$ , is equal to the area spanned by the vectors  $\mathbf{a}$  and  $\mathbf{b}$ . The unit cell volume  $V_{uc}$  is given by the triple vector product,  $V_{uc} = \mathbf{a} \cdot (\mathbf{b} \times \mathbf{c})$ .

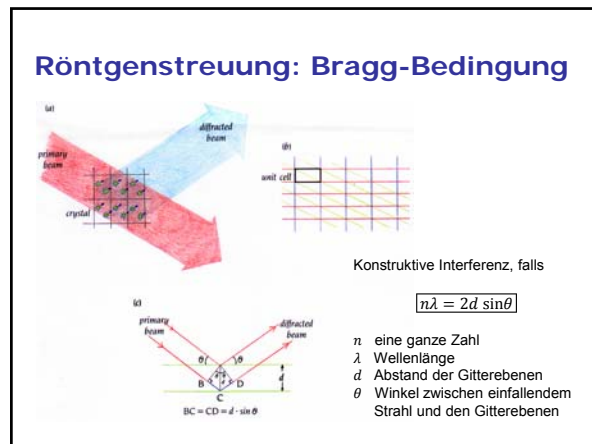
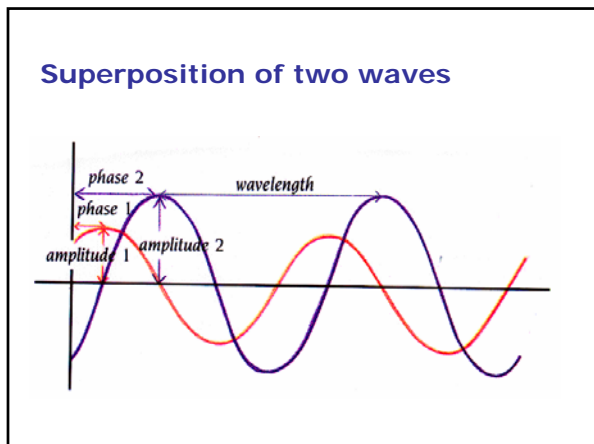
The angle between  $\mathbf{a}$  and  $\mathbf{b}$  is  $\gamma$ ; the angle between  $\mathbf{b}$  and  $\mathbf{c}$  is  $\alpha$ , and the angle between  $\mathbf{a}$  and  $\mathbf{c}$  is  $\beta$ . Similarly, the plane spanned by  $\mathbf{a}$  and  $\mathbf{b}$  is denoted as  $C$ , the plane between  $\mathbf{b}$  and  $\mathbf{c}$  is  $A$ , and the plane between  $\mathbf{a}$  and  $\mathbf{c}$  is labeled  $B$ .

The length of a unit cell vector is given by its norm:  $|\mathbf{a}| = a$ ,  $|\mathbf{b}| = b$ , and  $|\mathbf{c}| = c$ . The cell dimensions and angles are the six cell parameters (or cell constants)  $a, b, c, \alpha, \beta$ , and  $\gamma$ .









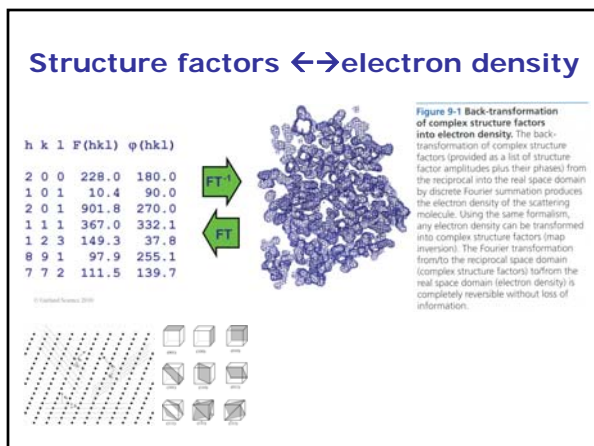
# Fourier transform

### Fourier transform relates structure factors and electron density

$$F(\mathbf{k}) = \int_R \rho(\mathbf{r}) e^{2\pi i \mathbf{r} \cdot \mathbf{k}} d\mathbf{r}$$

$$\rho(\mathbf{r}) = \int_{R^*} F(\mathbf{k}) e^{-2\pi i \mathbf{r} \cdot \mathbf{k}} d\mathbf{k}$$

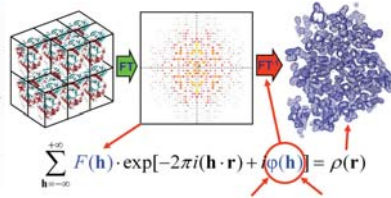
$\rho(\mathbf{r})$  electron density at position  $\mathbf{r}$  in real space  $R$   
 $\rho(\mathbf{r}) \in \mathbb{R}$  is real  
 $F(\mathbf{k})$  structure factor at position  $\mathbf{k}$  in reciprocal space  $R^*$   
 $F(\mathbf{k}) \in \mathbb{C}$  is complex with (measurable) amplitude  $|F(\mathbf{k})|$  and (not measurable) phase  $\alpha(\mathbf{k})$ , i.e.  
 $F(\mathbf{k}) = |F(\mathbf{k})|e^{i\alpha(\mathbf{k})}$



# Phases

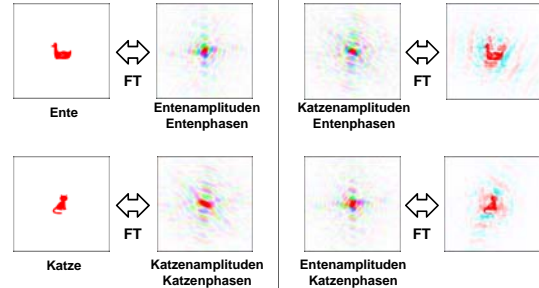
## The crystallographic phase problem

**Figure 9-15 The crystallographic phase problem.** The measurable component of the Fourier transform of the crystal is only the scalar structure factor amplitude  $|F(\mathbf{h})|$  proportional to the square root of  $I(\mathbf{h})$ . The missing phases  $\phi(\mathbf{h})$  must be supplied by additional phasing experiments or in the form of model phases via molecular replacement. The two necessary Fourier coefficients in the back-transform formula are emphasized in blue.



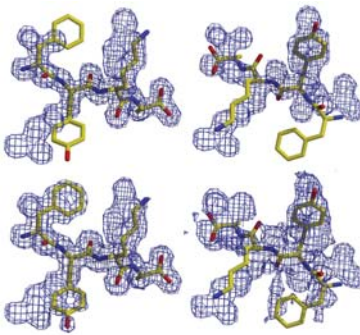
© Garland Science 2010

## Fourier Transformation: Phasen und Amplituden



<http://www.ysbl.york.ac.uk/~cowtan/fourier/>

## Phase bias in electron density maps



**Figure 9-18 Phase bias in electron density maps.** The upper panels show a mutant peptide (Phe-1 to Lys-10) and the same peptide rotated (leading to reverse chain direction) simply superimposed on the electron density of the original Val-Arg-1 to Ala peptide. The lower panels show the electron density reconstructed using the diffraction data from the new models above, but using the old starting phases from the original peptide. The result is quite disturbing: the shape of the electron density is still dominated by the starting model, and only weak outlines of the correct molecule density are visible. In the lower left panel, not even the direction of the peptide could be assigned for the reversed peptide with any certainty.

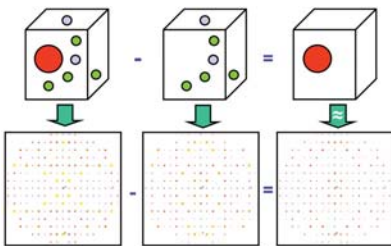
© Garland Science 2010

## Determination of phases

- **Ab initio phasing (direct methods):** Exploit theoretical phase relationships. Requires high resolution ( $< 1.4 \text{ \AA}$ ) data.
- **Heavy atom derivatives (multiple isomorphous replacement; MIR):** Crystallize the protein in the presence of several heavy metals without significantly changing the structure of the protein nor the crystal lattice.
- **Anomalous X-ray scattering at multiple wavelengths (multi-wavelength anomalous dispersion; MAD):** Incorporation of Seleno-methionine.
- **Molecular replacement:** Use structure of a similar molecule as the initial model.

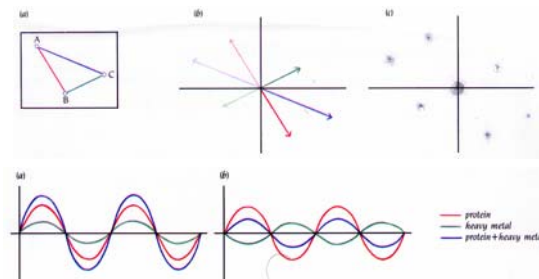
## Isomorphous difference data

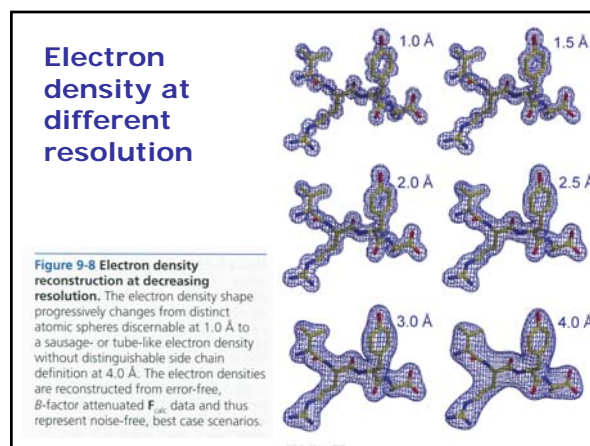
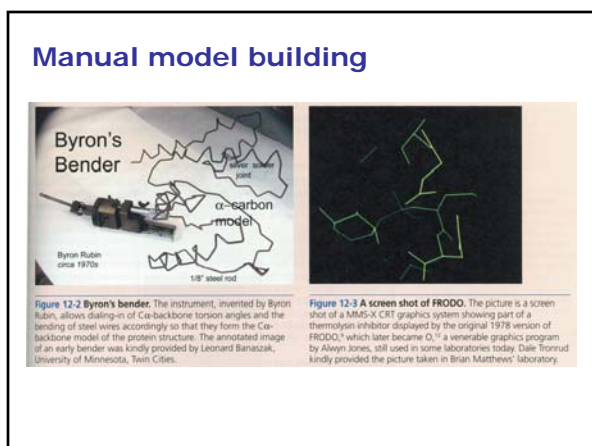
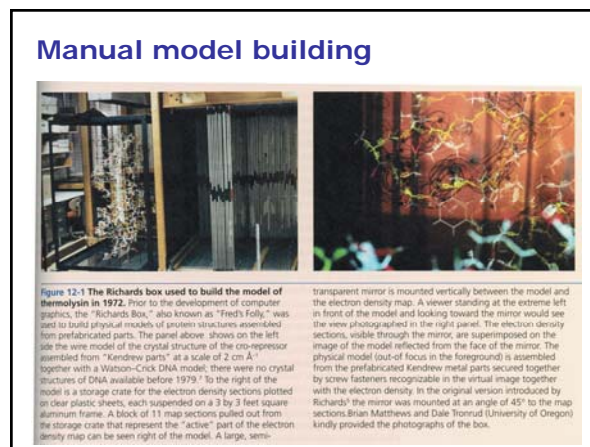
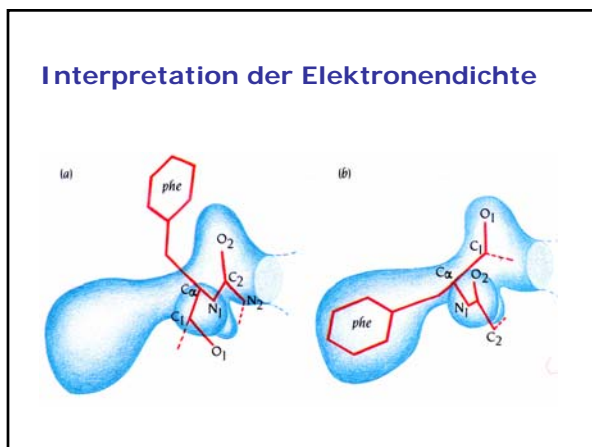
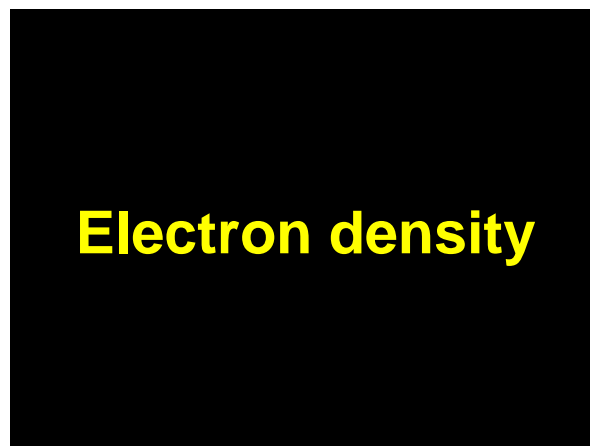
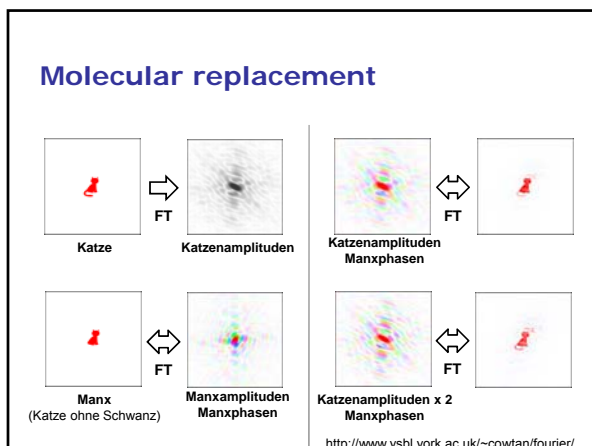
**Figure 10-1 The concept of isomorphous difference data.** The top line shows the gedankenexperiment in real space of subtracting a native protein crystal from an exactly isomorphous derivative crystal. The light atoms "cancel" out, and only the heavy marker atom remains in the difference crystal. While we cannot produce a real difference crystal, we can very well obtain a "difference diffraction pattern" from the differences between experimental data of the derivative and the native protein. The difference diffraction pattern has the same reciprocal dimensions and thus the same number of reflections, but represents the much simpler scenario of the "difference crystal."



© Garland Science 2010

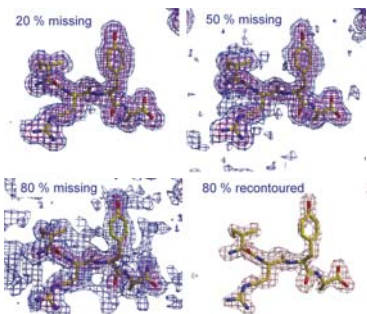
## Multiple isomorphous replacement (MIR)





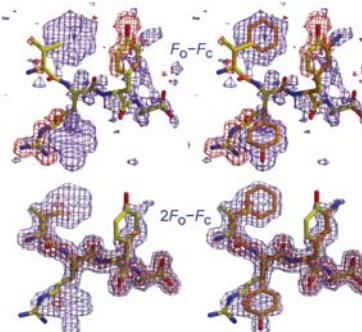
## Effect of omitted data

**Figure 9-18 Effect of randomly omitted data.** The panels show the effect of an increasing amount of randomly missing data. In the top left panel with 20% of the data randomly deleted, there is barely a difference noticeable compared with the maps generated from complete data shown in Figure 9-8. Even when the reconstruction misses 30% and 80% of data, the molecule is still traceable despite the increase in noise. About 800 out of 4000 reflections are all that is left in the reconstruction of the bottom electron density. The density in the bottom right panel is recontoured 80% missing at a higher  $\sigma$ -level, and the molecule is still traceable. Comparing the bottom left and bottom right panels emphasizes the importance of selecting a suitable density level for model density visualization and model building.

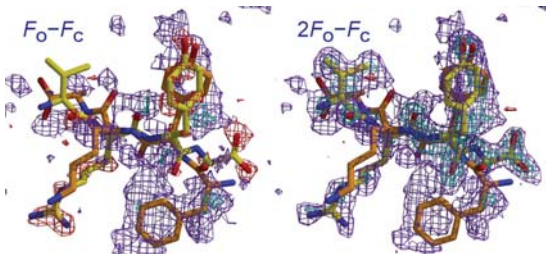


## Difference maps

**Figure 9-19  $F_o - F_c$  difference maps and  $2F_o - F_c$  maps.** The  $F_o - F_c$  difference maps in the top panel show negative difference density (where there should not be any density) in red and positive difference density (where there should be density) in purple. Both regions correctly reveal the difference between the starting model (yellow sticks) and the correct model (orange sticks, shown in the right panels). The sensitive difference maps are thus particularly valuable for detailed model correction. The  $2F_o - F_c$  maps in the bottom panel can be interpreted as a combination of the difference map ( $F_o - F_c$  map) and a  $2F$  (model) map. The  $2F_o - F_c$  map is contoured to amplify positive density and is well suited to early model building stages. Another common color scheme is red for negative difference density and green for positive.



## Poor start phases → poor electron density maps



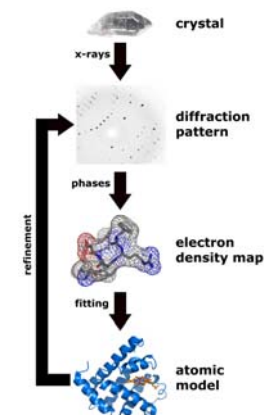
**Figure 9-20 Poor starting phases give poor electron density maps.** In the case of the reverse traced mutant peptide (orange sticks) as the true structure providing the intensities, no basic map type is able—despite good 1.5 Å data—to give sufficient clues as to how to correct the starting model (yellow sticks) that provided the phases. The difference map informs us in some parts about what is wrong, but none of the maps has sufficient reconstructive power to produce an outline of the correct orange molecule.

## Key concepts of model building

- The key to successful protein structure modeling is the cycling between local real space model building and model correction and global reciprocal space refinement.
- The molecular model is built in real space into electron density using computer graphics.
- Local geometry errors remaining after real space model building are corrected during restrained reciprocal space refinement by optimizing the fit between observed and calculated structure factor amplitudes.
- Successive rounds of rebuilding, error correction, and refinement are needed to obtain a good final protein model.
- While experimental electron density maps constructed from poor phases will be hard to interpret, an initial experimental map will not be biased toward any structure model.
- In contrast, when molecular replacement models are the sole source of phases, the electron density maps will be severely biased, and the map will reflect the model features.

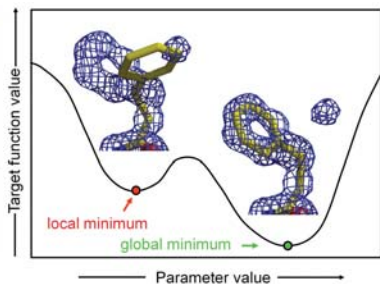
# Refinement

## Strukturermittlung



### Local minima during refinement

**Figure 12-7 Local minima and radius of convergence.** The figure visualizes the concept of trapping in a local minimum for a real space scenario. The  $C_{\alpha}$  atom of the magenta Phe ring is trapped in the electron density of a water molecule, in which it happens incidentally to fit quite well. In such cases, a refinement program may not be able to proceed upwards over the "activation" barrier—or may allow only limited positional parameter shifts—that prevent the large movement of the entire ring out of the partial density until it snaps into the correct electron density. Increased ability to overcome local minima by allowing "upwards movement" during parameter search implies higher radius of convergence and higher probability to approach the global minimum, generally at the cost of more computation and lower accuracy.



© Garland Science 2010

### X-ray crystallography: R-factor

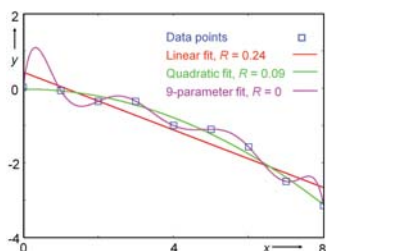
- Measures agreement between measured data (reflections) and 3D structure
- Definition: Relative difference between structure factors,  $F(hkl)$ , that were observed ( $F_{obs}$ ) and back-calculated from the 3D structure ( $F_{calc}$ ):

$$R = \frac{\sum |F_{obs}| - |F_{calc}|}{\sum |F_{obs}|} \quad \text{with} \quad I_{hkl} \propto |F(hkl)|^2$$

$I_{hkl}$  = intensity of reflection ( $hkl$ )

- Perfect agreement:  $R = 0$
- Good protein X-ray structure:  $R < 0.2$
- Random structure:  $R \approx 0.6$

### Over-fitting



**Figure 12-8 Fitting and over-fitting of a function.** The data points are measurements of the drop of a diffractometer pushed over a cliff, taken at constant time intervals  $x$ . The linear 2-parameter model (red line) describes the data poorly, but a quadratic 3-parameter fit (green graph) clearly describes the data very well, as it represents the physically correct model (a parabolic function describing the trajectory of a dropping object). We can further improve the fit (but not the model) by adding more parameters, and a 9-parameter polynomial (magenta) perfectly fits the data. Despite the perfect fit, the model is definitely nonsense, because the trajectory of the falling object takes upward turns, which is physically impossible. Following Bayesian reasoning the model, despite describing the data well, can be rejected based on a vanishingly small prior probability. In multi-parametric models such as crystal structures, over-fitting is unfortunately much less obvious, and cross-validation is a necessary practice.

### X-ray: Free R-factor

- Use, say, 90% of the data (reflections) for the structure determination
- Use the remaining 10% to compute the  $R$  value → "free"  $R$  value, obtained from independent data
- Detects errors better than conventional  $R$ -factor
- Each reflection influences whole electron density
- Many reflections → No problem to omit 10% of the reflections from the structure determination

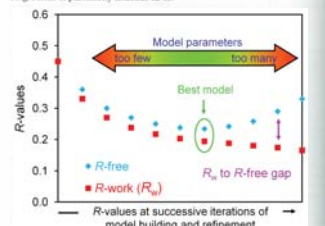
Brünger, A. T. (1992). Free  $R$  value: a novel statistical quantity for assessing the accuracy of crystal structures. *Nature* 355, 472-475.

### Cross-validation: R-free value

**Figure 12-9 Cross-validation R-value (R-free).** Before the first refinement step, the experimental data are split into a small test set (~5% of reflections) that is never used in refinement and the working data set. After each successive round of model rebuilding and completion, the current model is refined to convergence and both  $R$ -work and  $R$ -free are plotted for the corresponding refinement run. Both  $R$ -values improve progressively as the model becomes more complete and more parameters are introduced. At a certain stage in refinement, the model will be optimal, and introduction of further parameters into the model (often by unguided over-modeling of the discrete solvent or rigid side chain conformations) will not improve the model. At this point,  $R$ -free will stop improving, and with further overfitting  $R$ -free will even start to increase again while  $R$ -work keeps dropping (overfitting in the drawing, see Figure 12-4). For an actual example, see Figure 12-41. Given proper weighting, the best model will be the model with the lowest  $R$ -free (or to be precise, the highest log-likelihood). The gap between  $R$ -work and  $R$ -free is only a secondary mark of overfitting; it depends on a variety of parameters (Section 12.2), and observed values of the  $R$ -free/ $R$ -work ratio show a large variance (Figure 12-41). Note that each individual reciprocal space refinement run itself must be allowed to reach convergence—stopping an individual refinement run when its  $R$ -free reaches a transient minimum is bad practice (Sidebar 12-6).

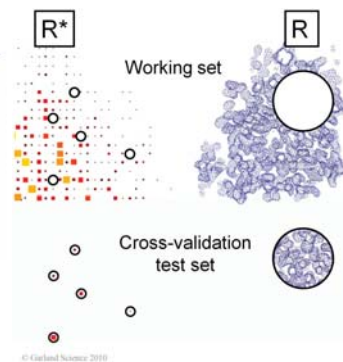
$$R_{work} = \frac{\sum |F_{obs} - F_{calc}|}{\sum |F_{obs}|} \quad \text{and} \quad R_{free} = \frac{\sum |F_{obs} - F_{calc}|}{\sum |F_{obs}|} \quad (12.30)$$

Axel Brünger introduced the  $R$ -free value<sup>10</sup> and has shown that  $R$ -free is related to the mean phase error<sup>11</sup> and is therefore a measure for the phase accuracy and thus for model quality, in contrast to the working  $R$ -value. A change to the model that improves its descriptions of physical reality will therefore also improve the fit to the excluded data, while purely cosmetic overparameterization will only lower  $R$ -work and not the cross-validation  $R$ -free (Figure 12-9). This can be loosely interpreted in terms of hypothesis testing: If the model refines as well without elaborate parameters as with them—determined by a lack of improvement in  $R$ -free—then the elaborate model is not any better and must be rejected on grounds of parsimony (Sidebar 12-3).



### Cross-validation in reciprocal and real space

**Figure 12-10 Cross-validation in reciprocal space and in real space.** A subset of unique reflections is set aside (the test or cross-validation set) before the model is refined and is excluded from any further refinement. The model is then refined against the working data set and the progress of the refinement is tracked against the test data set. In a similar fashion, the electron density or model in a questionable region can be removed, the model is again refined (often combined with a bias removal step), and the omitted region is inspected for new electron density. The figure layout follows an idea by A. T. Brünger.<sup>10</sup>



© Garland Science 2010

## Data-to-parameter ratio for X-ray protein structure determination

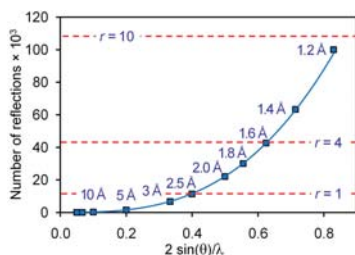


Figure 12.11 Data-to-parameter ratio for protein structures. The graph shows the number of reflections as a function of resolution (in units of  $\lambda$ ). The red dashed lines are drawn at numbers of reflections that correspond to a given data-to-parameter ratio  $r$ . The number of reflections approaches the number of refined parameters for positional and individual  $B$ -factor refinement around 2.5 Å. Below that resolution, unrestrained refinement is underdetermined, and only at atomic resolution does the redundancy of measured data become high enough ( $r > 10$ ) that unrestrained refinement becomes even remotely conceivable. The redundancy levels are generally valid for proteins with a solvent content around 50%. Tighter packing means a smaller unit cell and thus fewer reflections compared with loose packing. For torsion angle only refinement, the  $r/p$  ratio is slightly better (Section 12.2); therefore it is often the only available refinement protocol for low resolution (below  $\sim 3.5$  Å) structures.

## Key concepts of refinement I

- During refinement the parameters describing a continuously parameterized model are adjusted so that the fit of discrete experimental observations to their computed values calculated by a target function is optimized.
- Observations can be experimental data specific to the given problem, such as structure factor amplitudes, or general observations that are valid for all models.
- Stereochemical descriptors valid for all models such as bond lengths, bond angles, torsion angles, chirality, and non-bonded interactions are incorporated as restraints to improve the observation-to-parameter ratio of the refinement.
- The most accurate target functions are maximum likelihood target functions that account for errors and incompleteness in the model.
- Various optimization algorithms can be used to achieve the best fit between parameterized model and all observations, which include measured data and restraints.

## Key concepts of refinement II

- The radius of convergence for an optimization algorithm describes its ability to escape local minima and approach the global minimum, generally with increased cost in time and lower accuracy.
- Indiscriminate introduction of an increasing number of parameters into the model can lead to overparameterization, where the refinement residual measured as linear  $R$ -value still decreases, but the description of reality, i.e., the correct structure, does not improve.
- The evaluation of the residual against a data set excluded from refinement provides the cross-validation  $R$ -value or  $R$ -free. If parameters are introduced that do not improve the phase error of the model,  $R$ -free will not decrease any further or may even increase.
- Refined models carry some memory of omitted parts, which can be removed by slightly perturbing the coordinates and re-refining the model without the questionable part of the model.
- The known geometry target values for bond lengths, bond angles, and torsion angles as well as planarity of certain groups can be regarded as additional observations contributing to a higher data-to-parameter ratio.

## Key concepts of refinement III

- In addition, geometry targets constitute prior knowledge that keeps the molecular geometry in check with reality during restrained refinement.
- The geometry targets, chirality values, and non-bonded interactions are implemented as stereochemical restraints and incorporated into the target function generally in the form of squared sum of residuals in addition to the structure factor amplitude residual.
- The structure factor amplitude residual is commonly called the X-ray term (or X-ray energy) and the restraint residuals the chemical (energy) term.
- In terms of maximum posterior estimation, geometry target values and their variance define the prior probability of our model without consideration or knowledge of the experimental (diffraction) data.
- Geometric relations and redundancies between identical molecules in the asymmetric unit can be exploited through NCS restraints.
- Particularly at low resolution, strong NCS restraints are an effective means of stabilizing and improving the refinement.

## Key concepts of refinement IV

- In the early stages of model building, experimental phase restraints are also an effective means to stabilize and improve the refinement.
- The data-to-parameter ratio in protein structures is greatly increased through the introduction of stereochemical restraints.
- A protein of 2000 non-hydrogen atoms has about 8000 adjustable parameters and about the same number of restraints.
- At 2 Å about 15 000 to 25 000 unique reflections are observed for a 2000 nonhydrogen atom protein, which yields a total data to parameter ratio of about 2-3 at 2 Å.
- Anisotropic  $B$ -factor refinement consumes 5 additional parameters per atom, and is generally not advisable at resolutions  $< 1.4$  Å.
- The most difficult point in the parameterization of macromolecular structure models is accounting for correlated dynamic or static displacement.
- Isotropic  $B$ -factors are inadequate to describe any correlated dynamic molecular movement, and anisotropic  $B$ -factors, except at very high resolution, lead to overparameterization of the model.

## Key concepts of refinement V

- Molecular and lattice packing anisotropy can also affect diffraction, and adequate correction by anisotropic scaling, or in severe cases additional anisotropic resolution truncation, is necessary.
- Maximum likelihood target functions that account for incompleteness and errors in the model are superior to basic least squares target functions, particularly in the early, error-prone stages of refinement.
- Maximum likelihood target functions are implemented in REFMAC, Buster/ TNT, and CNS as well as the PHENIX/ cctbx programs, together with all commonly used restraint functions including phase restraints, which is of advantage at low resolution or in the early stages of refinement.
- Optimization algorithms are procedures that search for an optimum of a nonlinear, multi-parametric function.
- Optimization algorithms can be roughly divided into analytic or deterministic procedures and stochastic procedures.
- Deterministic optimizations such as gradient-based maximum likelihood methods are fast and work well when reasonably close to a correct model, at the price of becoming trapped in local minima.

### Key concepts of refinement VI

- Stochastic procedures employ a random search that also allows movements away from local minima. They are slow but compensate for it with a large radius of convergence.
- Evolutionary programming as used in molecular replacement or simulated annealing in refinement is a stochastic optimization procedure. This is generally of advantage if we do not know (MR) or are far from (initial model refinement) the correct solution.
- Deterministic optimizations can be classified depending on how they evaluate the second derivative matrix. They generally descend in several steps or cycles from a starting parameter set (model) downhill toward a hopefully but not necessarily global minimum.
- Energy refinement of a molecular dynamics force field and torsion angle refinement are two parameterizations that are used together with the stochastic optimization method of simulated annealing.
- In molecular dynamics the target function is parameterized in the form of potential energy terms and the development of the system is described by equations of motion. In torsion angle parameterization, the structure model is described by its torsion angles, which requires fewer parameters than coordinate parameterization.

### Key concepts of refinement VII

- Both molecular dynamics and torsion angle parameterization are often combined with simulated annealing optimization, where the molecular system is perturbed and returns to equilibrium according to an optimized slow cooling protocol.
- Dummy atom placement and refinement is used for discrete solvent building, model completion, and phase improvement in general.
- Dummy atoms are placed in real space in difference electron density peaks, the new model is refined unrestrained in reciprocal space, and in the new map poorly positioned atoms are removed and new ones placed again.
- Dummy atom refinement can be combined with multi-model map averaging where it forms the basis of bias minimization protocols and the automated model building program ARP/wARP.

### Model building and refinement practice I

- Building of a model into an empty map begins with the tracing of the backbone.
- Tracing is aided by density skeletonization, followed by placement of C $\alpha$  atoms into positions where side chains extend from the backbone.
- The sequence is docked from known atom positions from the heavy atom substructure or sequences of residues of characteristic shapes.
- The initial model is refined in reciprocal space with geometric restraints and phase restraints, and the next map is constructed from maximum likelihood coefficients.
- The model is then further completed and refined in subsequent rounds with increasing X-ray weights while tracking *R*-free and stereochemistry. Nuisance errors are removed after analysis in a polishing step.
- Automated model building programs greatly simplify model building, and auto-built models often only need to be completed and polished. Autobuilding programs follow similar steps as manual model building and employ pattern recognition algorithms to identify residues.

### Model building and refinement practice II

- Rebuilding poor initial molecular replacement models can be aided by a first step of torsion angle-simulated annealing (TA-SA) refinement.
- The large radius of convergence of TA-SA facilitates the necessary large corrections and escape from local minima. Also, before automated model rebuilding and correction, TA-SA can improve the amount and quality of the model that is automatically rebuilt.
- In low resolution structures the backbone can be traced correctly, but the sequence may be shifted. Such register errors can be hard to detect from electron density shape alone and are usually detected by poor side chain interactions or unusual environment.
- A common mistake leading to overparameterization of the model is overbuilding of the solvent. Discrete water molecules should have hydrogen bonded contact(s) to other solvent molecules or to protein.
- Poorly placed waters tend to drift away during refinement because of lack of density and restraints and often end up far away from other molecules and with high *B*-factors.

### Model building and refinement practice III

- Binding sites have a tendency to attract various detritus from the crystallization cocktail, and will therefore often contain some weak, unidentifiable density that can be (wishfully) mistaken for desired ligand density.
- Plausible binding chemistry, ligand conformation, and independent evidence are necessary to avoid misinterpretation.
- The three major criteria for abandoning refinement and rebuilding are:
  - (i) No more significant and interpretable difference density in  $mF_{\text{obs}} - DF_{\text{calc}}$  maps remains.
  - (ii) No more unexplained significant deviations from stereochemical target values and from plausible stereochemistry remain.
  - (iii) The model makes chemical and biological sense.
- Global measures such as absolute values of *R* and *R*-free (or the level of boredom) do not determine when refinement is finished.

### Literatur über Kristallstrukturbestimmung

- **B. Rupp, *Biomolecular Crystallography*, Garland, 2010.**
- **W. Massa, *Kristallstrukturbestimmung*, Teubner, 2007.**
- **C. Branden & J. Tooze, *Introduction to Protein Structure*, Garland, 1999.**

### Skript

[www.bpc.uni-frankfurt.de/guentert/wiki/index.php/Teaching](http://www.bpc.uni-frankfurt.de/guentert/wiki/index.php/Teaching)