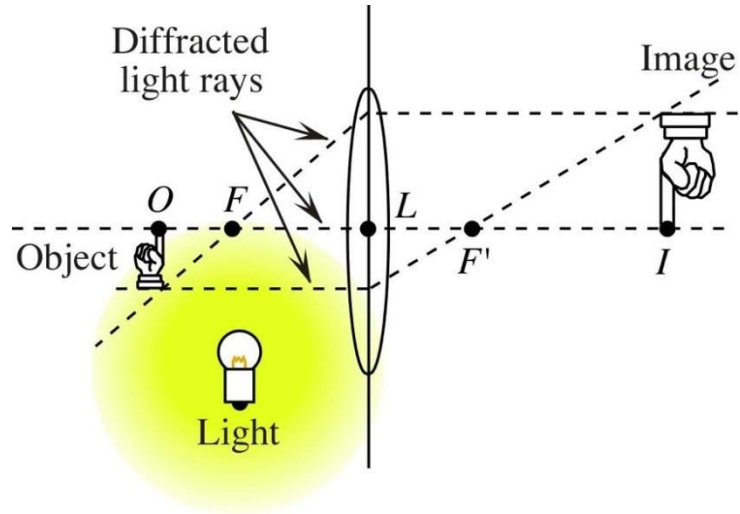# LECTURE 9: X-ray Crystallography

# Why X-Ray Crystallography?

- What information can we get from crystallographic structures?
  - 3D structure of the molecules of interest
    - Relative locations of atoms of proteins (**atomic coordinates**)
    - Location of binding/active site(s)
  - Temperature factors: the magnitude of the thermal motions of individual atoms or groups of atoms (for low resolution structures).
    - Local quality of the model
  - "Snap-shots" of the protein in different functional states
    - With a cleverly-chosen ligands or substrate/transition state/product analogs

- X-ray crystallography can be applied to:
  - Mechanistic enzymology
  - Structure-based drug design/optimization
  - Understanding biological phenomenon
  - Comparing X-ray and NMR structure
  - etc.

## (Light) microscope:



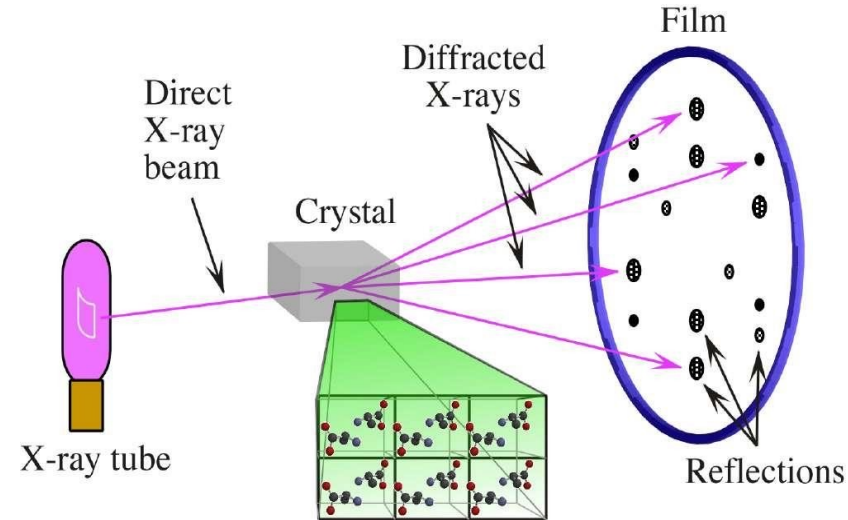## X-ray diffraction of a crystal



**Limitations:**

Object needs to be larger than the wavelength of the light (visible light 400-700 nm, atoms = 0.15 nm apart)

X-rays(0.08-0.6 nm) cannot be focussed by lenses
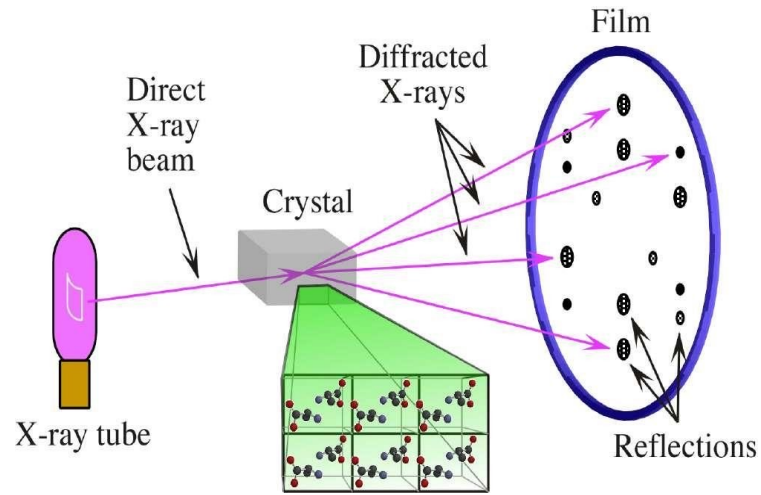
Molecules are very weak scatterers

A crystal contains many molecules in identical orientation

Diffracted X-rays of individual molecules 'add up' (positive interference) to produce strong reflections

Computers can simulate a lens and reconstruct the image (Fourier transform)

# X-ray diffraction of a single protein crystal

- The wavelength of X-ray is in the same size regime as covalently bonded atoms (~1.5 Å).

- There are two fundamental reasons why we can't simply take a picture of a protein:
  1. X-rays(0.08-0.6 nm) cannot be focused by lenses
  2. A single protein molecule is a very weak scatterer

- The solutions to these problems are:

1. Arrange the molecules into a 3D repeating array – a crystal.
   - Protein molecules are arranged in precisely the same orientation, which amplifies x-rays scattering. In addition, the radiation scattered from a repeating array can cause constructive interference – **diffraction**.

2. Use a "mathematical lens" – the Fourier Transform – to reform an image of the protein molecule.
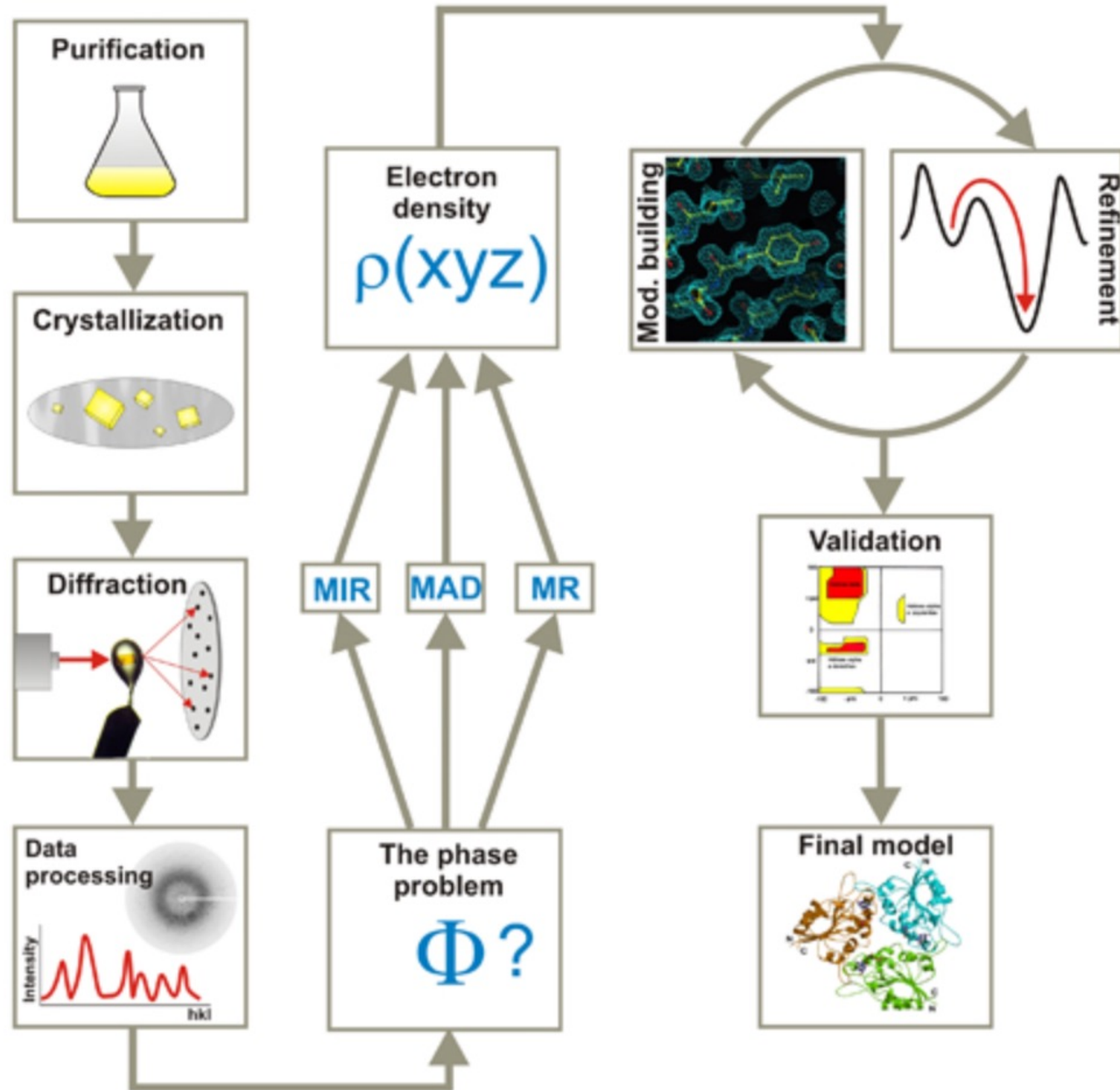
# Outline of A Complete Structure Determination by X-ray Crystallography

**Electron density**

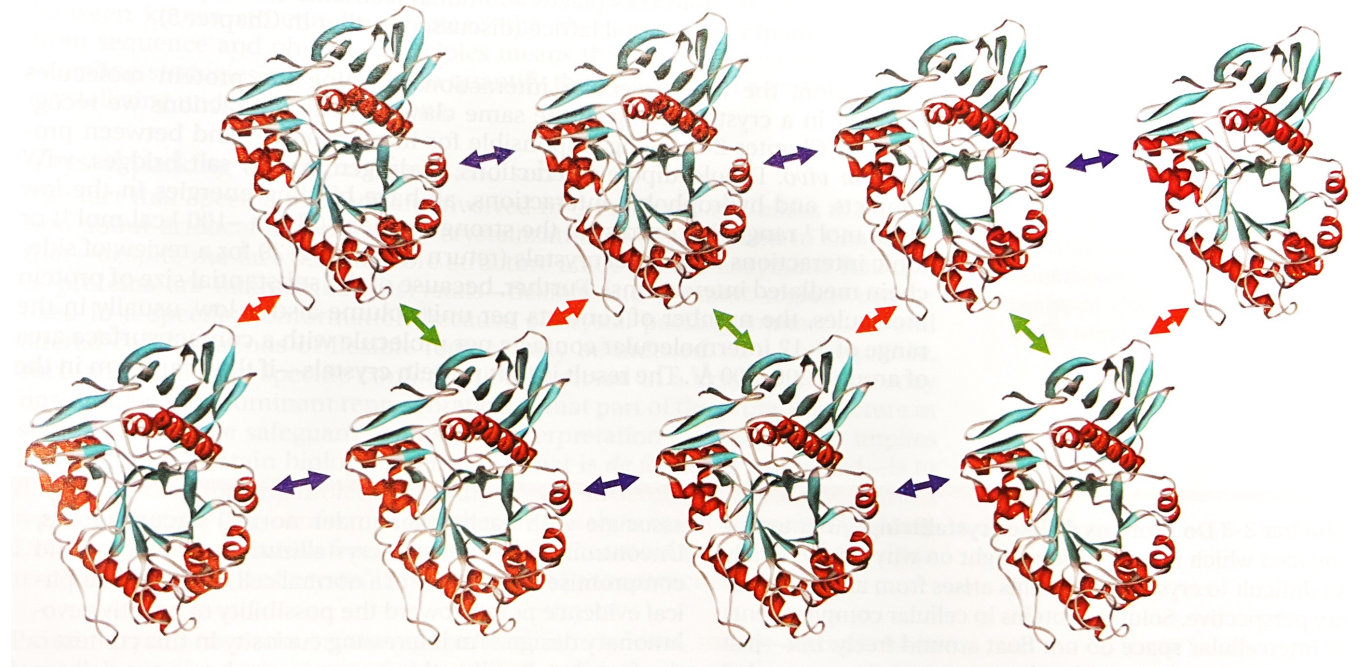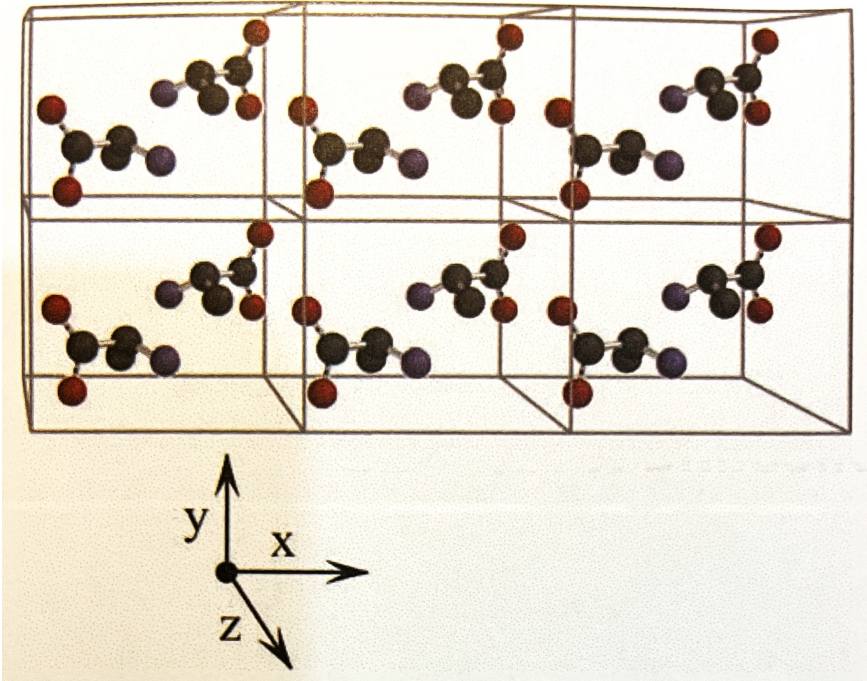$$\rho(xyz) = \frac{1}{V} \sum_{\substack{hkl \\ -\infty}}^{+\infty} |F_o(hkl)| \cdot e^{-2\pi i[hx+ky+lz-\phi_c(hkl)]}$$

Magnitude

Phase

Purification

Crystallization

Diffraction

Data processing

Electron density ρ(xyz)

MIR  MAD  MR

The phase problem Φ?
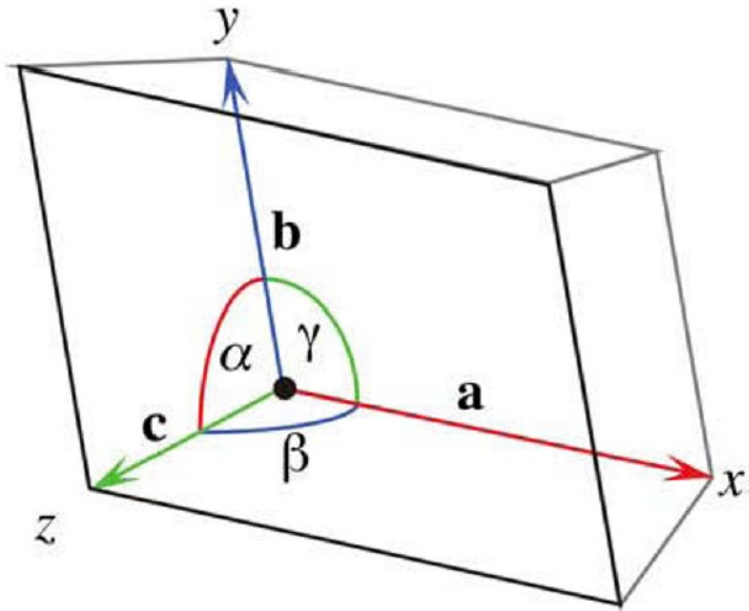
Mod. building

Refinement

Validation

Final model

# What are Protein Crystals?

- A crystal of any molecule (ions, sugar, or proteins) is simply an ordered array of molecules making specific intermolecular contacts to form a repeating, 3-dimensional lattice.
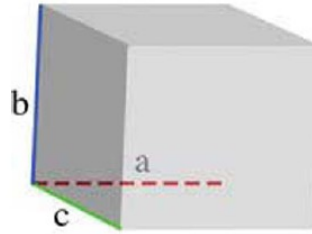  - Crystals are stacks of bricks called **unit cell**.

# Unit Cell

How many different ways to arrange points in space where each point would have an identical "atmosphere"?
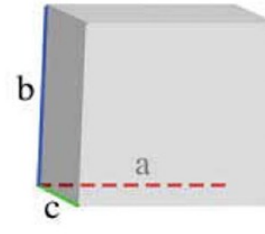


General (triclinic) unit cell, with edges **a, b, c** and angles α, β, γ
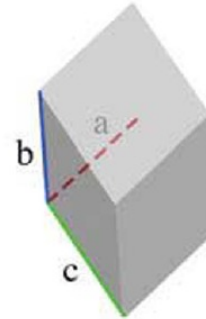
**Cubic**
$a=b=c,$
$\alpha=\beta=\gamma=90°$
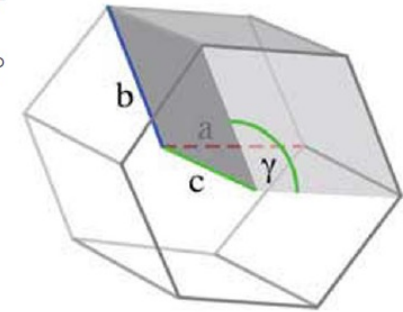
**Tetragonal**
$a=b\neq c,$
$\alpha=\beta=\gamma=90°$

**Orthorhombic**
$a\neq b\neq c,$
$\alpha=\beta=\gamma=90°$

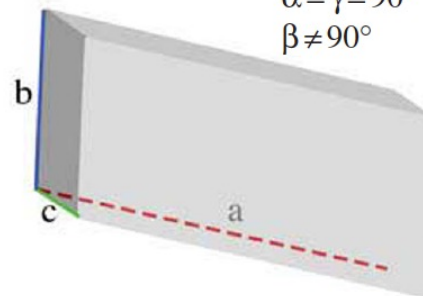**Rhombohedral**
$a=b=c,$
$\alpha=\beta=\gamma\neq 90°$

**Hexagonal**
$a=b=c,$
$\alpha=\beta=90°$
$\gamma=120°$

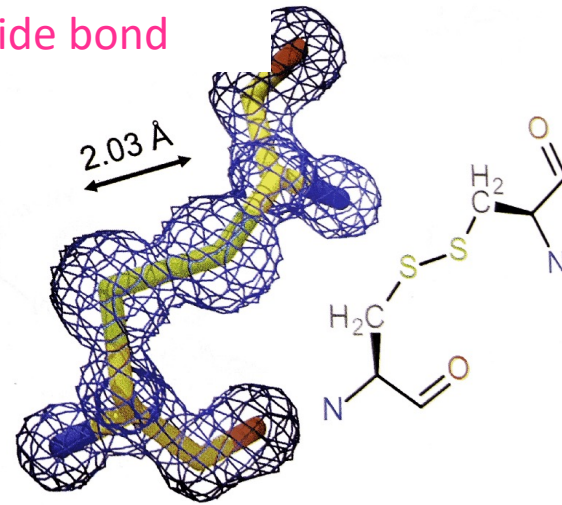**Monoclinic**
$a\neq b\neq c,$
$\alpha=\gamma=90°$
$\beta\neq 90°$

**Triclinic**
$a\neq b\neq c,$
$\alpha\neq\beta\neq\gamma\neq 90°$

# Interactions in Crystals

- Ionic interaction, H bond, hydrophobic, covalent bond.


Disulfide bond


Charge-$\pi$ stacking


H bond


Ionic interaction


Solvent network

# What are Protein Crystals?

- Protein crystals do not act like more familiar crystals (table salt, etc.)

  - Protein crystals contain ~40-60% solvent (water, ions, small molecules, PEG, detergent, etc.) – we need large crystals for diffraction.

  - Protein molecules are so big but the number of interactions between them is relatively small, which makes protein crystals extremely fragile – everything damages them: mechanical stress, osmotic shock, temperature changes, etc.

# X-ray Diffraction

- **X-ray Scattering** is a change in the direction of motion of X-ray beam because of a collision with electron cloud of the atoms in protein.

- **X-ray Diffraction** is scattering where the intensity is modulated by interference from other scatterers held in a regular array (crystal).



© Encyclopædia Britannica, Inc.

- Conditions that Produce Diffraction (constructive interference): Bragg's Law

# Crystal Planes



- In crystals, diffraction is treated as if it were reflection from sets of equivalent, parallel planes of atoms in a crystal.
  - Each spot in the diffraction pattern is called a *reflection*.

- Planes in a crystal can be specified using a notation called Miller indices [*hkl*], where *h*, *k*, and *l* are reciprocals of the plane with the *x*, *y*, and *z* axes.

# Bragg's Law in Protein Crystals



$$2d_{hkl} \sin \theta = n\lambda$$

- **$d_{hkl}$**: interplanar spacing between parallel planes

- θ: reflected angle

- n: an integer

- λ: X-ray wavelength

Wave 1 travels the same distance as Wave 2 plus an added distance BC+BD (2BC)

# What do we need for a crystallographic experiment?

1) Well-Ordered crystals
2) Size: bigger size for enough diffracting material
   – More proteins packed into the crystal

How big and how well-ordered? Depends on your goal.

- A totally unknown protein: low resolution (3-3.5 Å) model is significant to the field.

- Structure-guided drug design: better ordered and often bigger crystals for high resolution data.

# Growing Crystals: Protein Sample

- Not all proteins can crystallize

- Tags or no tags?

- Purity/contaminants
  - Freshness, conformational states

- Quality control
  - Gel filtration, static/dynamic light scattering (to look for aggregate protein), Thermal shift assay, etc.

- Concentration
  - generally, 10-15 mg/mL is the starting point. But the range is huge (1 mg/mL to 100 mg/mL. The "right" concentration is determined empirically.

- The composition of the protein storage buffer
  - Identity of buffer, pH, presence/amounts of salts/additives/stabilizers, etc.

# Growing Crystals – Protein Phase Diagram[1]

**Protein Phase Diagram**



- Crystals grow because it becomes thermodynamically favorable for the protein to come out of solution. But the kinetics of the process must be sufficiently slow.

- Set up conditions where protein is soluble to begin with.

- Move to the area in the phase diagram that favors protein coming out of solution but not aggregating (metastable).

# Growing Crystals: Assay



A  Layer of immersion oil  
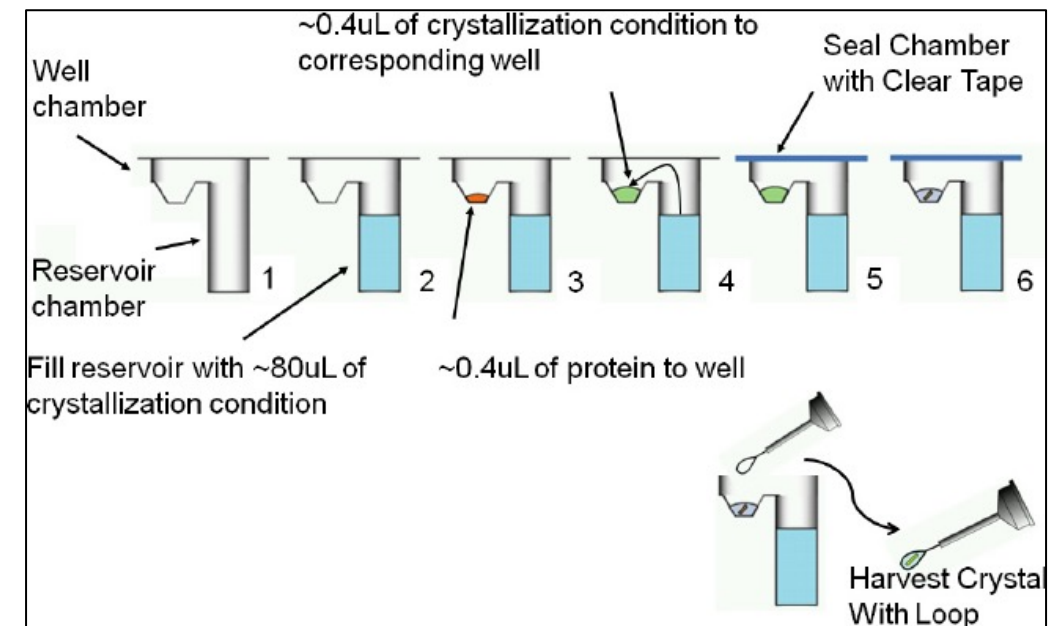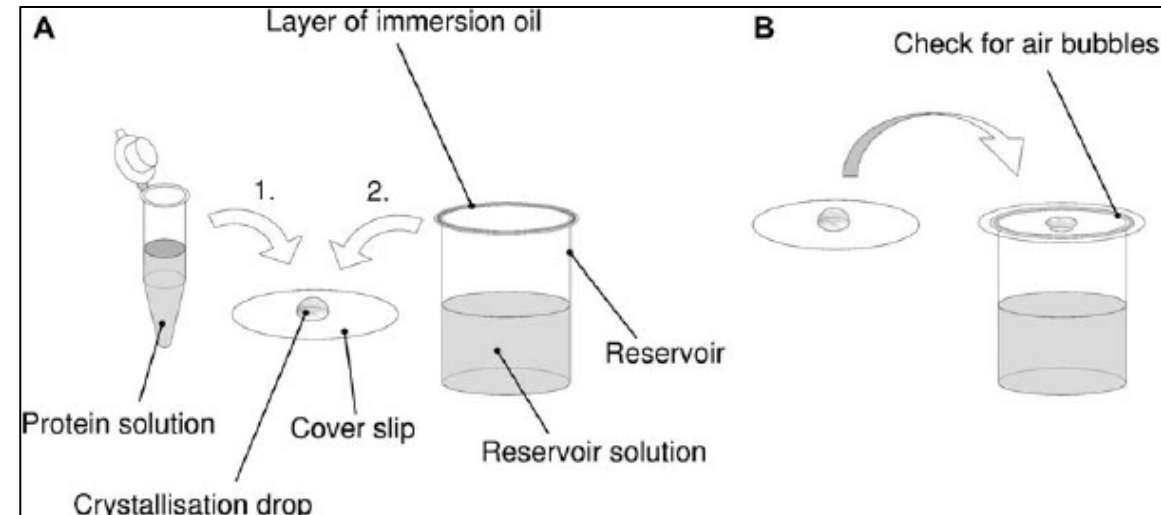B  Check for air bubbles  
Protein solution  Cover slip  Reservoir  
Crystallisation drop  Reservoir solution

- **Vapor diffusion** (hanging or sitting drops): most common
  - Protein and crystallization solutions are mixed together either on a coverslip or in the well of a sitting drop support.
  - The crystallization drop is then sealed inside a chamber with a much greater volume of the crystallization solution.
  - Over time, water leaves the drop until the concentrations of precipitant in the drop and well equalize.



Well chamber  
~0.4uL of crystallization condition to corresponding well  
Seal Chamber with Clear Tape  
Reservoir chamber  
Fill reservoir with ~80uL of crystallization condition  
~0.4uL of protein to well  
Harvest Crystal With Loop

Armour, B. L., Barnes, S. R., Moen, S. O., Smith, E., Raymond, A. C., Fairman, J. W., Stewart, L. J., Staker, B. L., Begley, D. W., Edwards, T. E., Lorimer, D. D. **J. Vis. Exp**. (76), e4225, doi:10.3791/4225 (2013).
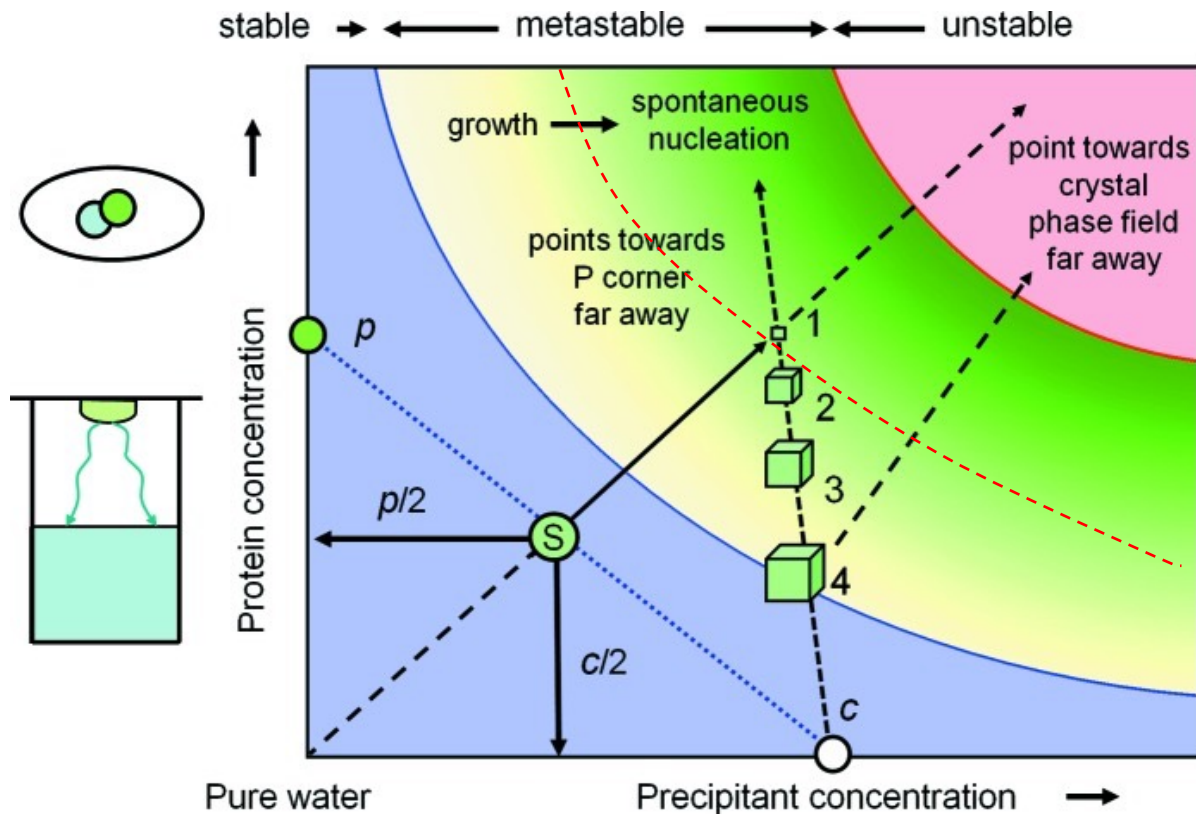
- All crystallization(reservoir) solutions making up the screen all have
  1. a **precipitating agent** – salt (NaCl or $(NH_4)_2SO_4$) or polymer (like polyethylene glycol (PEG).
  2. A **buffer**
     - Note: pH is important in crystallization because it influence the distribution of charges across the surface. Charged residues are particularly important for forming crystal contacts.
  3. An **Accessory/additive compound** (often salts, but can be small organic molecules)

$\Rightarrow$ Some volume of protein solution is mixed with a volume of the crystallization solution and the experiment is sealed and watch for days/weeks/months.

# Growing Crystals – Protein Phase Diagram₂

**Protein Phase Diagram**



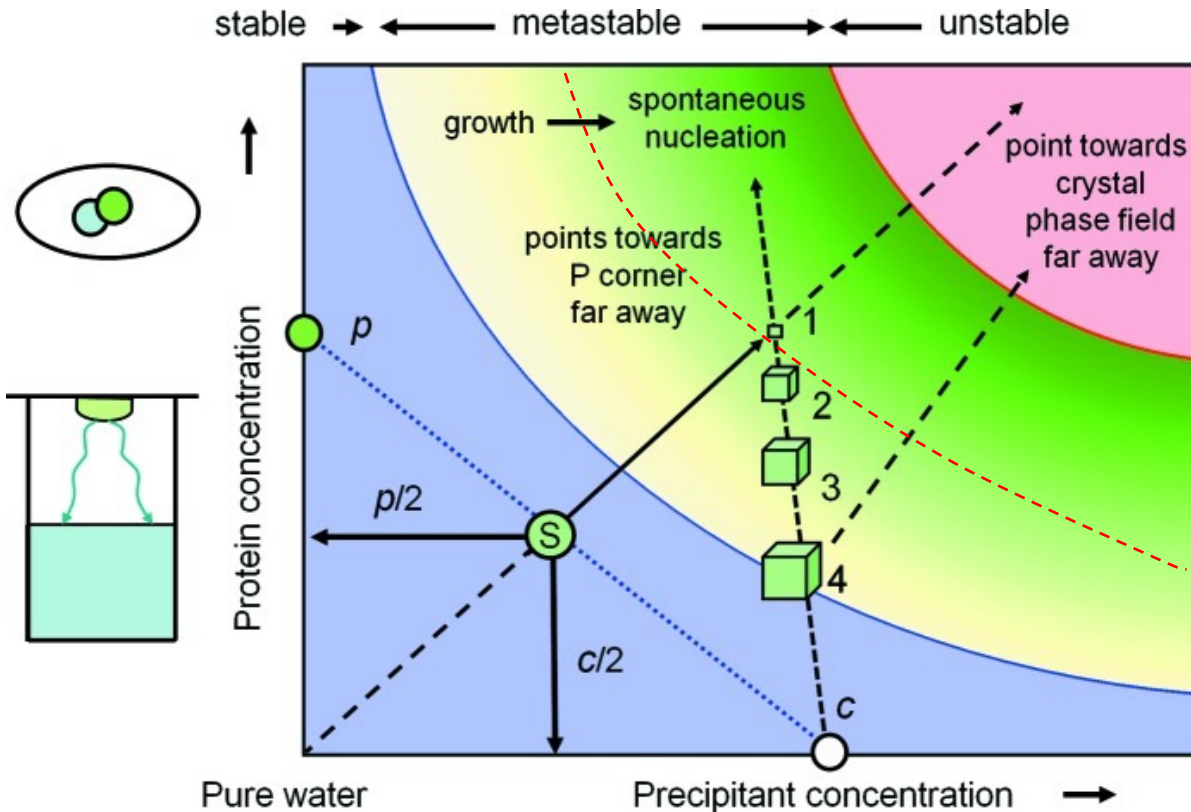- As vapor diffusion proceeds, both the protein and precipitant concentrate until both reach their stock value, which puts the protein just over the imaginary line (beyond which crystals can spontaneously nucleate).

- As the small number of nuclei grow into single crystals, the protein concentration falls below the nucleation line and into the growth region of the metastable zone.
  - This is what we want.

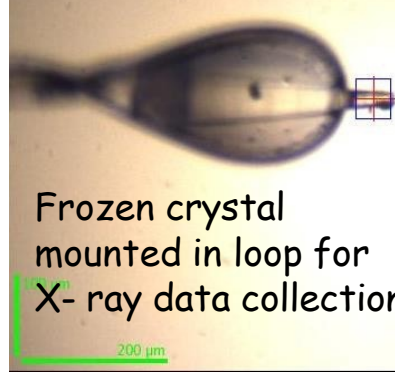# Growing Crystals – Protein Phase Diagram₃

**Protein Phase Diagram**



- **Think/Discuss:**
1. What will happen if the equilibration moves far into the nucleation zone?

2. If the protein or precipitant concentration is low (just over the boundary into metastable zone), what will we see?

**Note:** The phase diagram is helpful in terms of understanding the various experiments and their outcomes, but in practice it is very difficult to reconcile these "what should happen" scenarios with what actually happens in real life! There are too many variables in crystallization for such simple explanations to hold up.

# Set Up

**X-ray data are measured on frozen crystals (~100K)**

Frozen crystal mounted in loop for X- ray data collection

200 µm

50keV Electrons

Focussing Mirrors (or Monochromator)

Chi-Circle

ω

Area Detector(s)

Crystal

χ

φ

Primary X-ray Beam

Focussed Beam

ω

Rotating Anode (Cu)
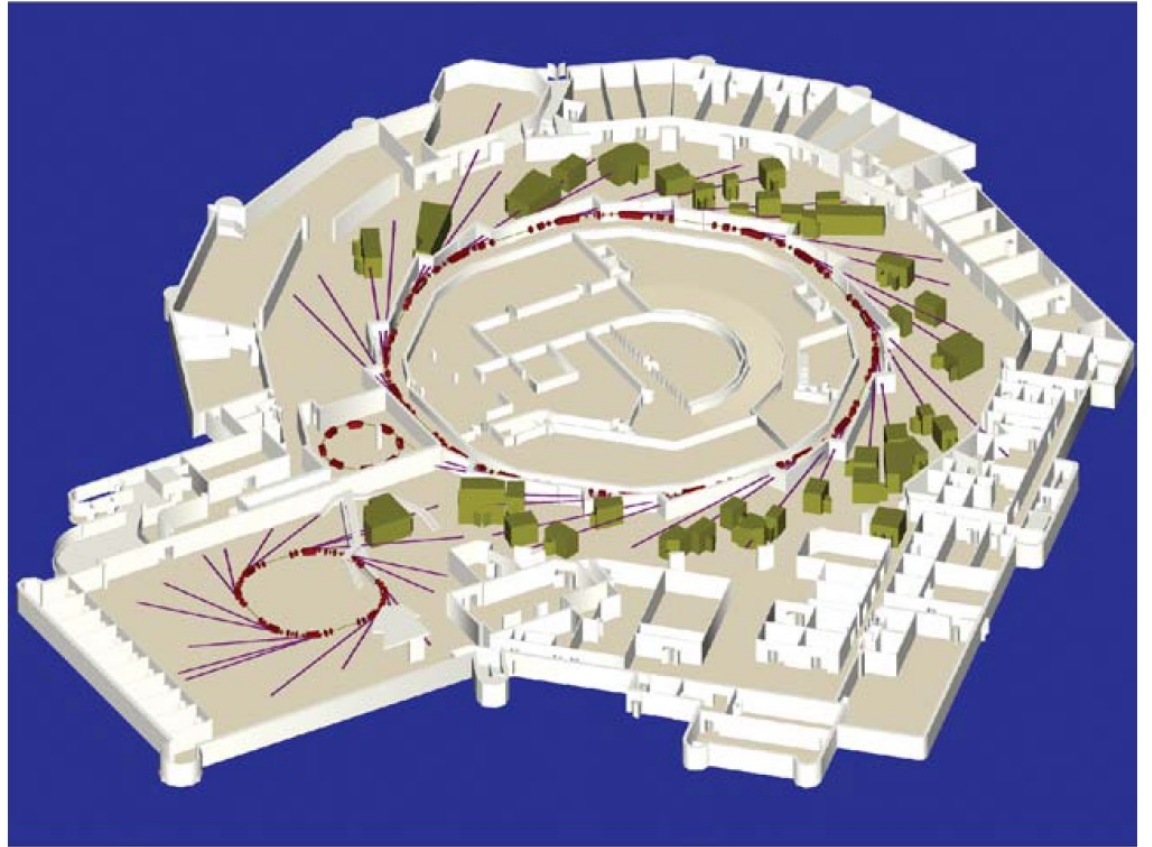
θ

4-Circle Goniometer (Eulerian or Kappa Geometry)

- **X-ray Source**: X-ray tubes, rotating anode tubes, or **particle storage rings**
- **Goniometer** (holder)
- Detectors: X-ray films, CCD cameras, or Multiwire detectors.

# For high quality X-ray data collection extremely intense synchrotron beam lines are used

In the giant particle storage ring, electrons or positrons circulate at the nearly speed of light. The charged electrons emits energy (synchrotron radiation) when forced in into curved motion, and in accelerators, the energy is emitted as X-rays.
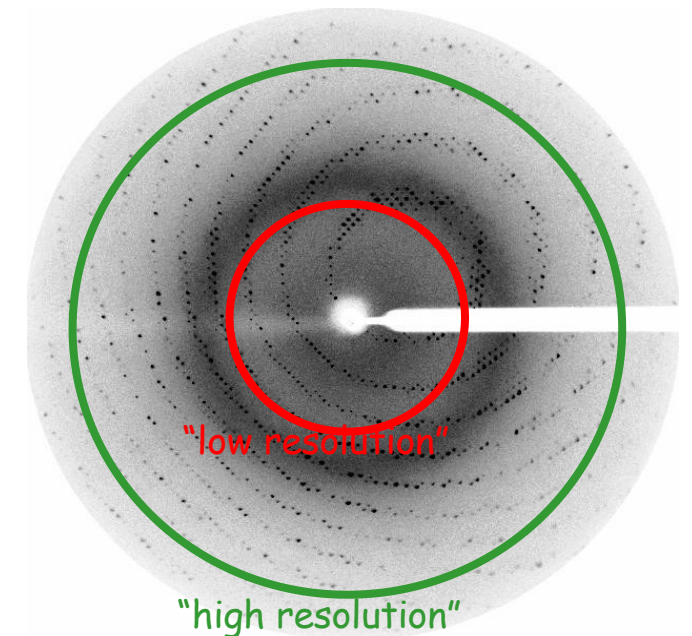
Accessory devices tangential to the storage ring provide powerful monochromatic X-rays at selectable wavelengths.



National Synchrotron Light Source, Brookhaven National Laboratory, Brookhaven NY. (Left) Aerial view of exterior. (Right) Interior floor plan

# Data Collection

- This is the last physical part of the experiment – everything after this is computational.

- The result of this step is a set of 90-hundreds of diffraction images that (hopefully) cover the entire 3D diffraction pattern.

- Two Pieces of Data
  - The position of a reflection point on the reciprocal lattice, given by coordinates *h, k, l*. Determined by the direction reflected.
  - The intensity of the reflection.

- The degree of order in the crystal determines the quality of the diffraction data and ultimately the quality of the final atomic model.



"low resolution"

"high resolution"
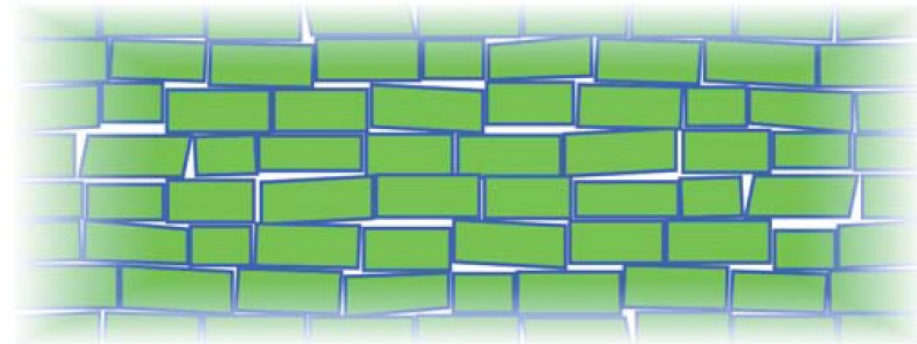
# Mosaic Spread



**Figure 3.2** ▶ Crystals are not perfectly ordered. They consist of many small arrays in rough alignment with each other. As a result, reflections are not points, but are spherical or ovoid, and must be measured over a small angular range.
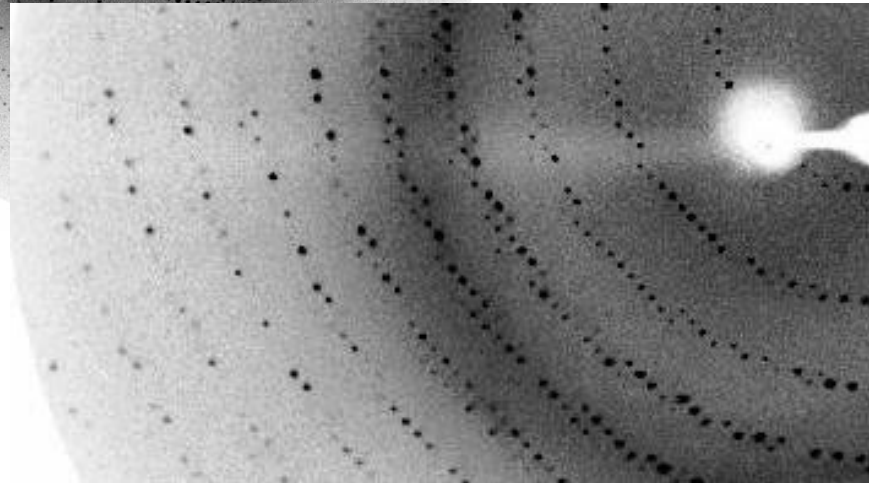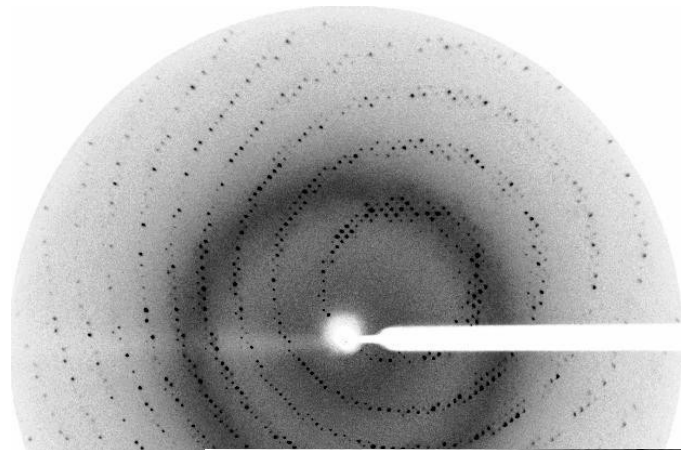
- Mosaic spread refers to the misalignment of separate domains within a crystalline lattice. We can think of blocks of perfectly ordered crystalline lattice that are fused together in a slightly haphazard arrangement to form a single crystal. This misalignment causes the diffracted x-rays to be slightly divergent (not perfectly collinear). As a result, the recorded reflection intensities (signal) are blurred over a large area of the detector that has intrinsic noise (diffuse x-ray scatter in air, noise in the x-ray detection photochemistry and electronics, etc.).

- The signal-to-noise ratio is lower for a mosaic crystal in comparison to a well-ordered crystal.
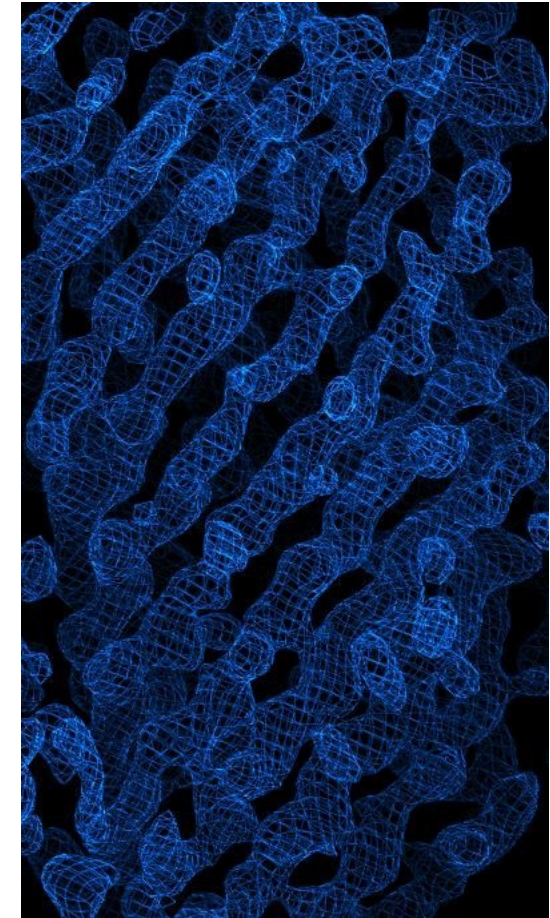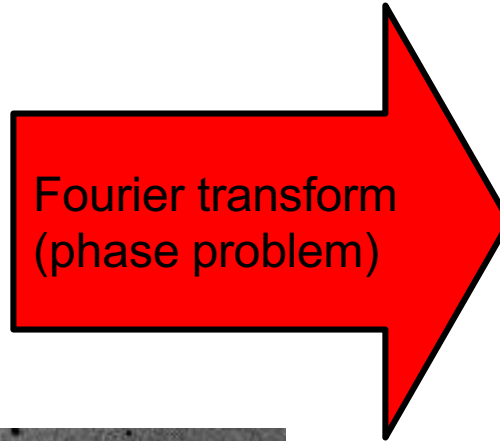
# Data Processing

Three Phases:

1.  **Indexing** – determining the unit cell dimensions and symmetry, as well as the orientation of the crystal with respect to the beam

2.  **Integration** – Integrate the intensities for each reflection.

    ❖ The product of this step is a file with Miller indices (h. k, l) and intensity for each measured reflection.

3.  **Scaling** – Scaling attempts to accounts for error in the intensity measurements.

    ❖ The product of this step is a file with **unique** Miller indices (h. k, l), intensity, and error for each measured reflection.

# Phasing



Raw data: Thousands of intensities of reflections

Fourier transform (phase problem)

Electron density

# Phase Problem

Electron density

Magnitude

Phase

$$\rho(xyz) = \frac{1}{V} \sum_{\substack{hkl \\ -\infty}}^{+\infty} |F_o(hkl)| \cdot e^{-2\pi i[hx+ky+lz-\phi_c(hkl)]}$$

- From diffraction to electron density map requires Fourier Transform.
- Magnitude is related to the intensity.
- But it's impossible (hard) to extract phase angle from diffraction pattern directly, which is known as **phase problem**.
- To address this:
  - Isomorphous replacement
  - Anomalous scattering
  - Molecular replacement

# Isomorphous Replacement

- Heavy atoms contribute to some reflections strongly.
- Insert a heavy metal atom into crystal protein and locate in diffraction pattern and in the cell.  Use the location of metal ion to find the phase angle for the other protein atoms.

- Requirements:
  - Add atom with the same unit cell size.
  - Cannot disturb protein structure.
  - Often use Hg, Pt, Au.
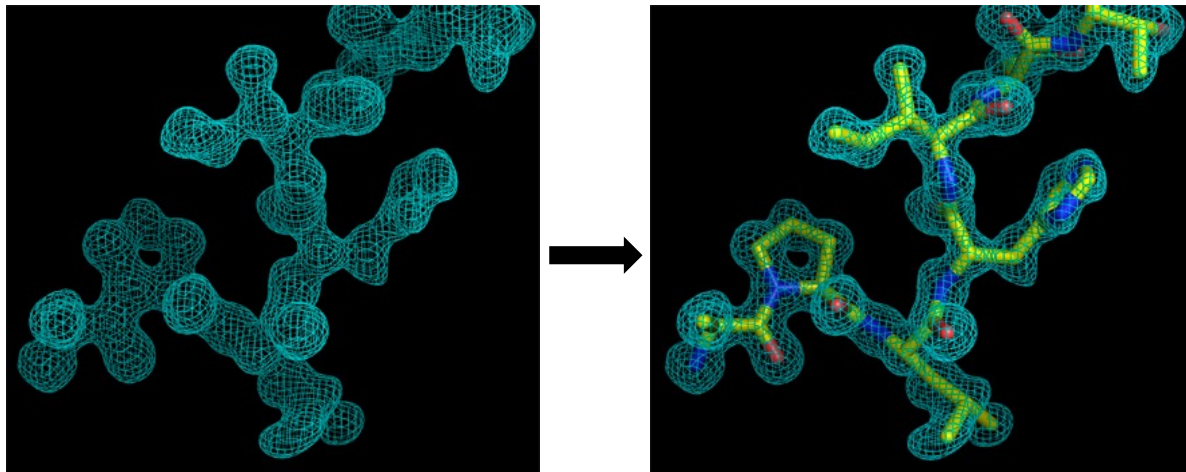
# Anomalous Scattering

- Heavy atoms can absorb X-rays of specific wavelength.
- When the X-ray wavelength is near the heavy-atom absorption edge, a fraction of the radiation is absorbed by the heavy atom and reemitted with altered phase.

- During protein production, Selenomethionine, instead of Met, is added to the special media.

- Single-wavelength anomalous diffraction phasing (SAD)
- Multi-wavelength anomalous diffraction phasing (MAD)

# Molecular Replacement

- Use the phases from structure factors of a known protein as initial estimates of phases for a new protein.

- Place a model of known protein in the unit cell of the new protein is called **molecular replacement**.


- sequence homology of known protein is similar to the new protein

- Apo-enzyme when study conformational changes with ligand bound.

- Domains of different know proteins

# Model Building

- Determination of initial phases allows us to calculate an **initial electron density map** that is used to build the first rough model of the protein
  - **$2|F_o|-|F_c|$**: "real map" that is used to add atoms to the model.
  - **$|F_o|-|F_c|$**: difference map that highlights sections where the model differs greatly from the experimental data.
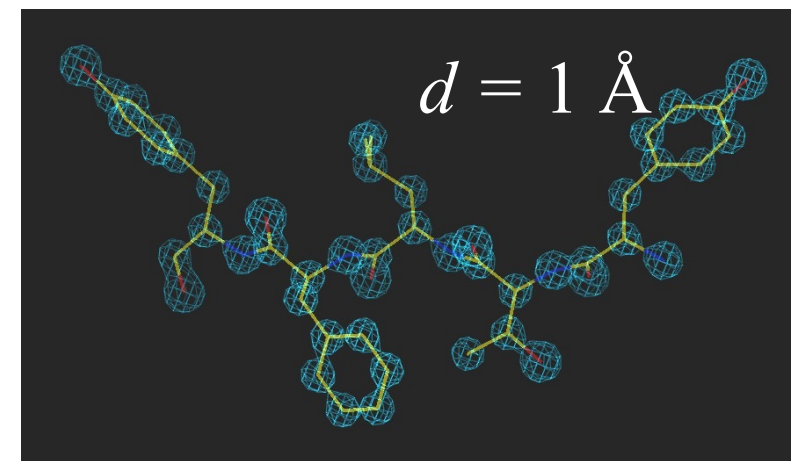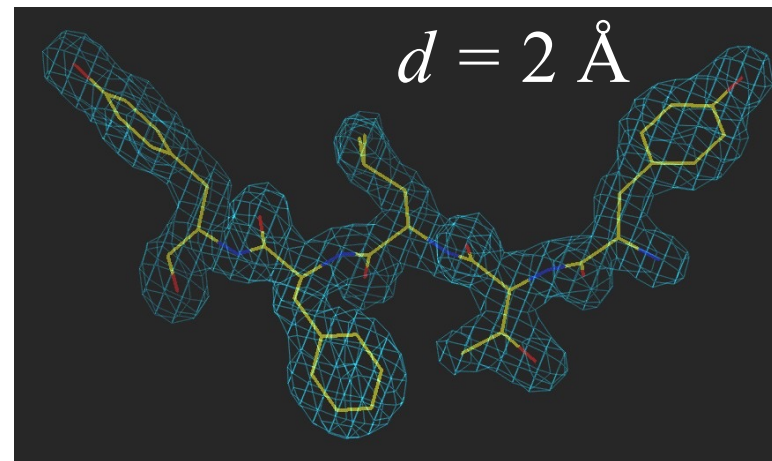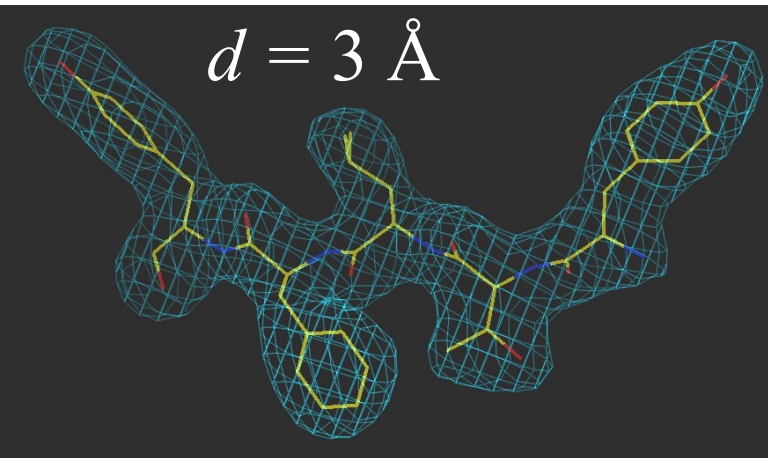


- Either by hand or **automatically**, depending on the resolution.
- Refine the model to make the final model agree closely with the experimental data and known stereochemical data.
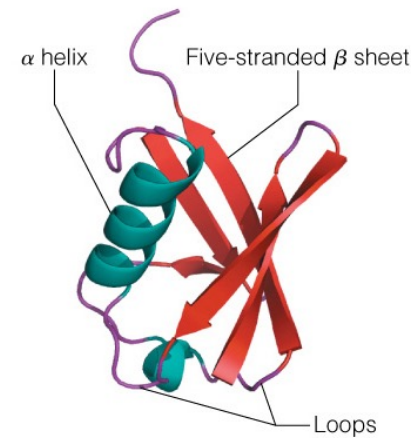
# Atomic Resolution

- Crystallographers express the resolution of a structure in terms of distance.
  - If the resolution limit is smaller than the atoms distance, atoms will appear as separate maxima and their positions can be obtained with high precision
  - If the resolution limit is more than the atoms distance $\Longrightarrow$ fused electron density
    - i.e., C-C bond length is 1.5 Å



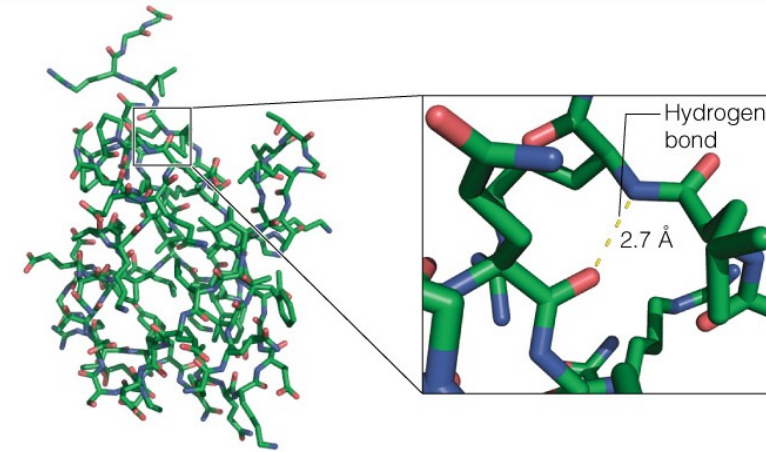$d = 3$ Å $\qquad$ $d = 2$ Å $\qquad$ $d = 1$ Å

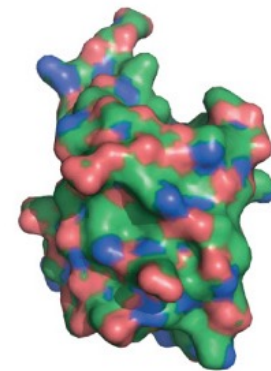- Proteins can be represented by different models:
  - Cartoon "ribbon" structure
    - Highlights secondary structure
    - Typically ignores side chains
  - Stick model
    - Highlights side chains
    - Hard to see secondary structure.
  - Surface model
    - Highlights overall "shape" or tertiary structure.
    - Can be color coded by charge, polarity
  - Space filling model
    - Highlights the electron shell overlaps
    - no longer possible to see the angles between bonds



(a) A cartoon model of the protein backbone. An $\alpha$ helix (cyan), is packed against a five-stranded $\beta$ sheet (red) composed of parallel and antiparallel strands. Loops are shown in magenta.

(b) A stick model showing the locations of all atoms (excluding H atoms). C atoms are green, N atoms are blue, and O atoms are red. The inset shows a hydrogen bond of 2.7 Å between main-chain atoms.

Atom coloring is the same as in panel (b).

# Judging The Quality of Crystallographic Structures[1]

**Data statistics:**

1. Completeness – How much of the diffraction pattern was actually recorded. Both the overall completeness and the completeness for high resolution data are important. (overall: ~100%, high resolution shell: ~80%)

2. Multiplicity (redundancy): average number of independent measurements of each reflection in a crystallographic data set. (~ 3.8 multiplicity)

3. $R_{merge}$ or $R_{symm}$ – These "R factors" are a measure of the error in the data – literally the difference between any measurement of a particular reflection and the average of multiple measurements of that reflection (the same reflections are measured multiple times). (Lower $R_{merge}$ is better)

4. Signal to noise ratio ($I/\sigma$) – ratio of the average intensity of the diffraction spots and the background intensity
   - At least 2.0 for the highest resolution data.

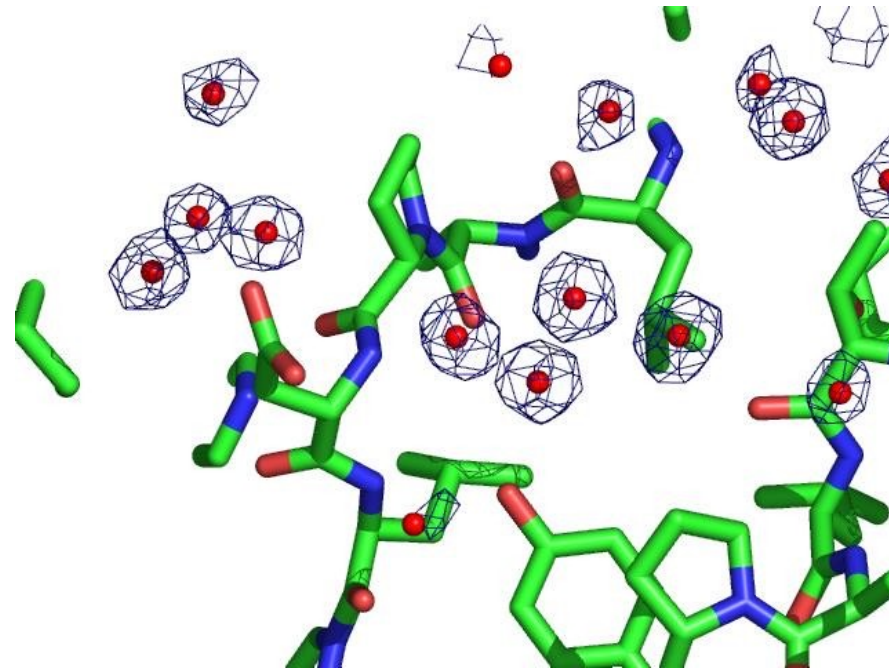# Judging The Quality Of Crystallographic Structures$_2$

**Model quality statistics:**

1. Crystallographic R factors ($R_{cryst}$ or $R_{free}$) – measure the difference between the model and the experimental data. The difference between should be less than 5% at any resolution.

2. Deviations from ideality – The differences between the bond parameters and standard values should be < 0.01 Å and 1.5$^o$ for angles.

3. B factors (temperature factors) – measure the degree of thermal motion of each atom. Show mobile and/or disordered parts of the model.

Things you - as a potential user of crystallographic data - should know about crystals and crystal structures

# Two types of solvent: Ordered and Disordered

- Ordered water molecules show up as discrete blobs of electron density in contact with the protein or with other ordered water molecules

- Disordered water regions show up as featureless (flat) electron density

# PDB files:

- Basically, just simple text files
- At the top: information about the crystal:
  - Which proteins/ligands etc
  - Crystallization conditions
  - How was the structure solved
  - The resolution
  - Some useful statistics to judge the quality of the crystal
  - How to get from the structure to the biological unit
  - Remarks about missing bits etc.
- Crystal parameters: cell dimensions/space group
- A list of all atoms in the structure

# A crystal structure according to the protein data bank (PDB)

occupancy

x,y,z coordinates (Å)

```
ATOM      25  N    ASP A 928      19.062    9.157   35.067  1.00   4.73              N
ATOM      26  CA   ASP A 928      19.770   10.123   34.232  1.00   4.58              C
ATOM      27  C    ASP A 928      19.075    9.938   32.899  1.00   4.56              C
ATOM      28  O    ASP A 928      19.074    8.824   32.351  1.00   5.39              O
ATOM      29  CB   ASP A 928      21.259    9.776   34.071  1.00   3.13              C
ATOM      30  CG   ASP A 928      22.112   10.245   35.233  1.00   5.52              C
ATOM      31  OD1  ASP A 928      21.693   11.114   36.025  1.00   5.42              O
ATOM      32  OD2  ASP A 928      23.239    9.742   35.349  1.00   7.93              O
ATOM      33  N    VAL A 929      18.417   10.985   32.405  1.00   3.68              N
ATOM      34  CA   VAL A 929      17.726   10.864   31.125  1.00   4.63              C
```
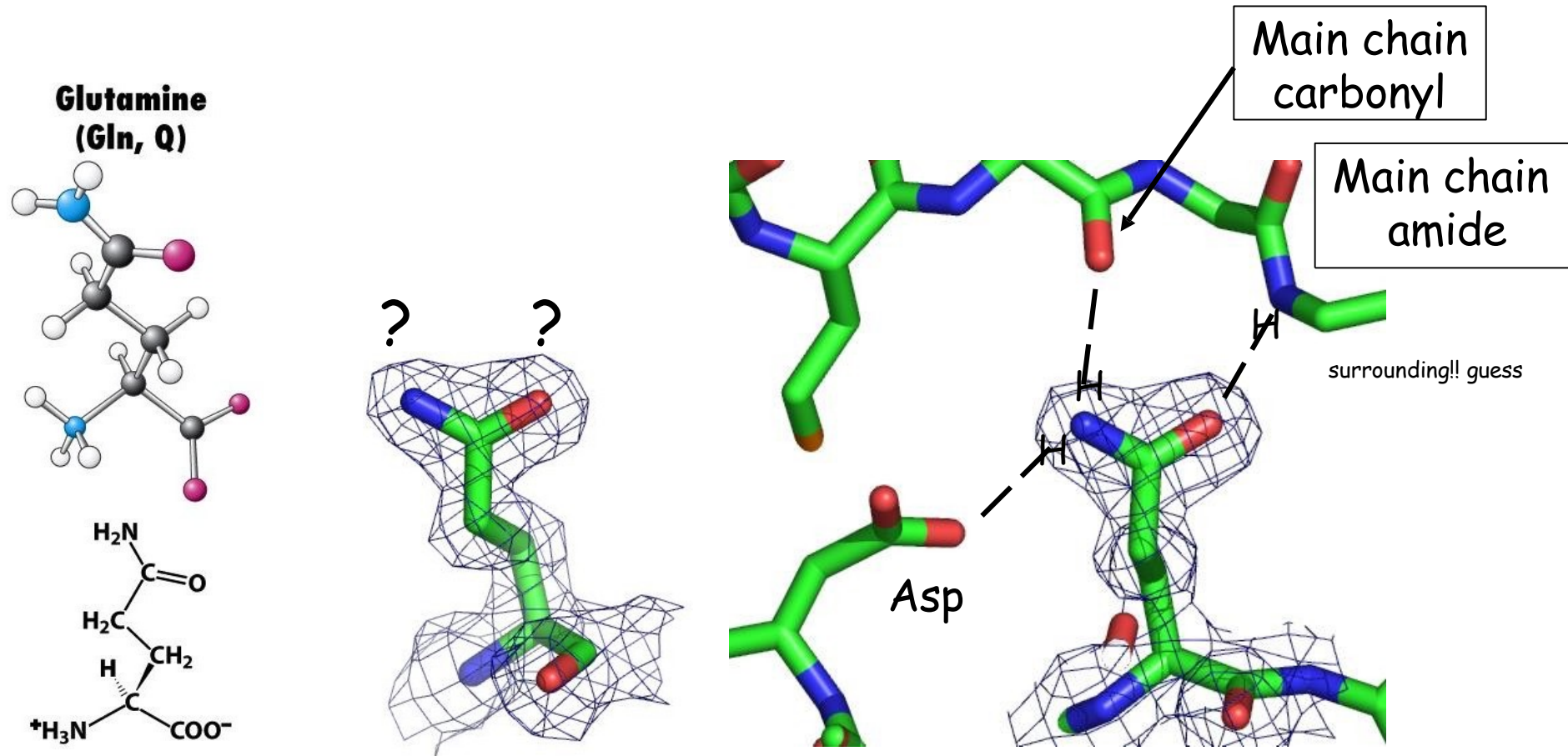
Isotropic B-factor or temperature factor is a measure of the mobility of an atom

$B$ (Å²) = $8\pi^2 \langle u^2 \rangle$, where $\langle u^2 \rangle$ is the mean square atomic displacement

# "Occupancy"

- The occupancy "$n_j$" of atom "j": is a measurement of the fraction of molecules in the crystal in which atom j occupies the position specified in the model.

- If all molecules in the crystal are precisely identical, then occupancies for all atoms are 1.00.

- Occupancy is necessary because occasionally two or more distinct conformations are observed for a small region like a surface side chain. For example, if the **two conformations occur with equal frequency, then atoms involved receive occupancies of 0.5 in each of their two possible positions.**

- So, "occupancies" is an estimates of the frequency of alternative conformations, giving some additional information about the dynamics of the protein molecule.

# Position of N and O atoms in Gln (and Asn) side chain must be inferred from hydrogen bonding network



Glutamine
(Gln, Q)

?   ?

Main chain carbonyl

Main chain amide

surrounding!! guess

Asp

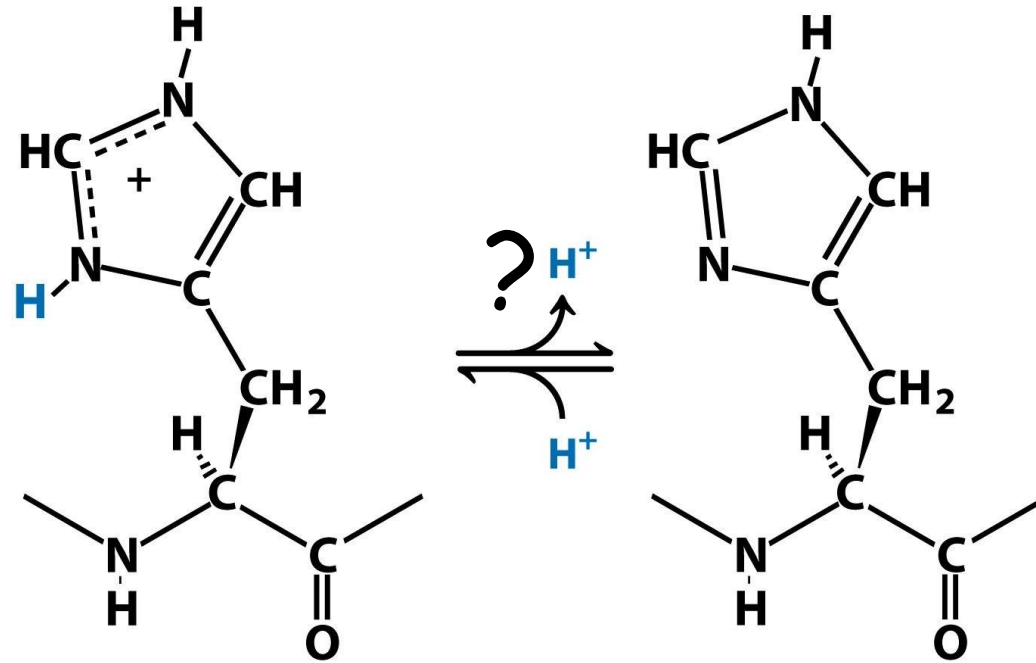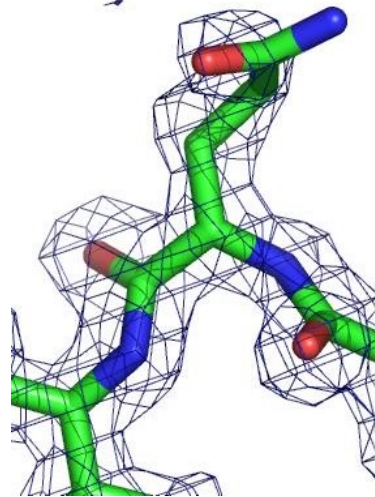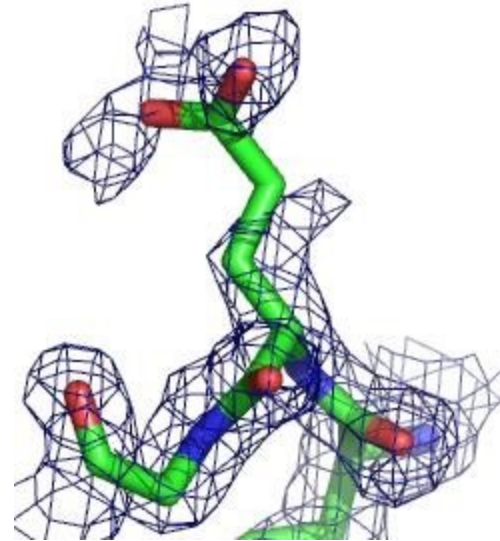# The same holds for the orientation and protonation of the imidazole ring of histidines



Figure 2-15
Biochemistry, Sixth Edition
© 2007 W. H. Freeman and Company

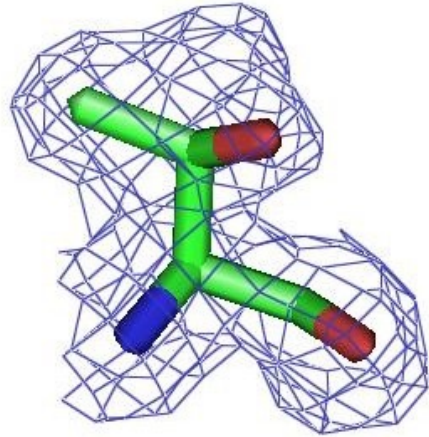# A pdb file may contain residues for which no, or only limited electron density is visible
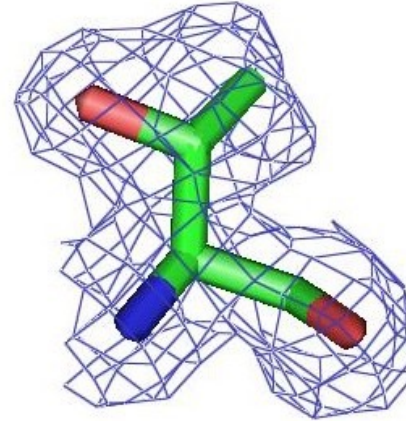


No density for amide N of glutamine



Break in side chain density of glutamate

# Sometimes the electron density suggests two side chain conformations but may only one is modeled in the pdb file
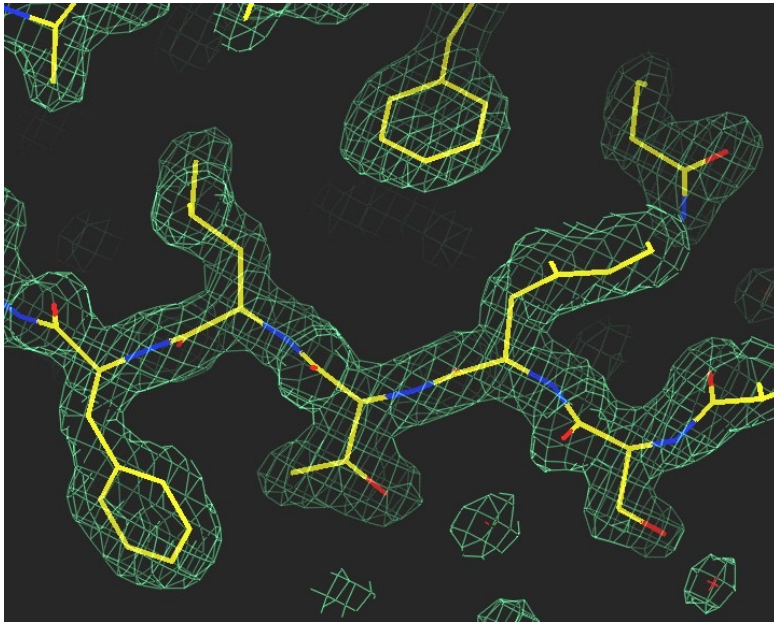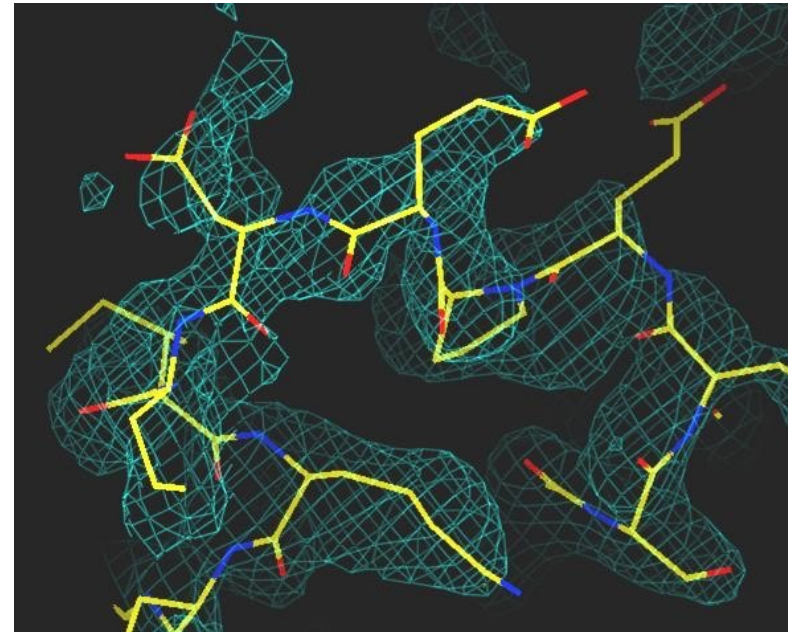
Threonine side chain conformation present in pdb file

Alternative conformation that is also compatible with electron density

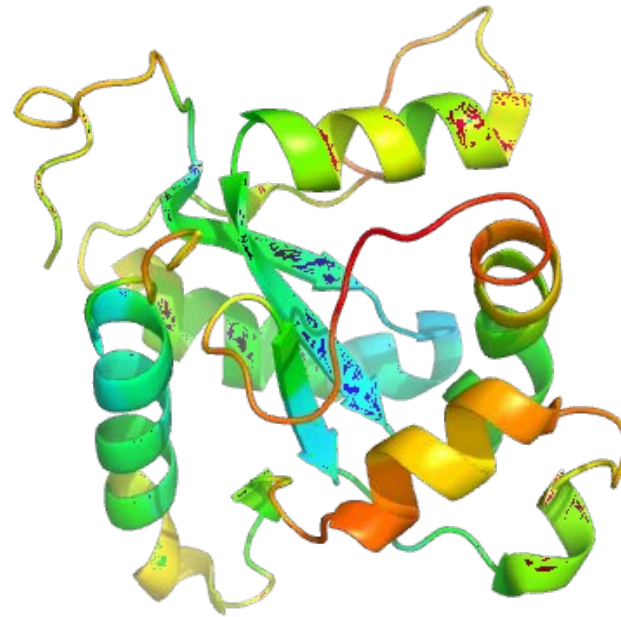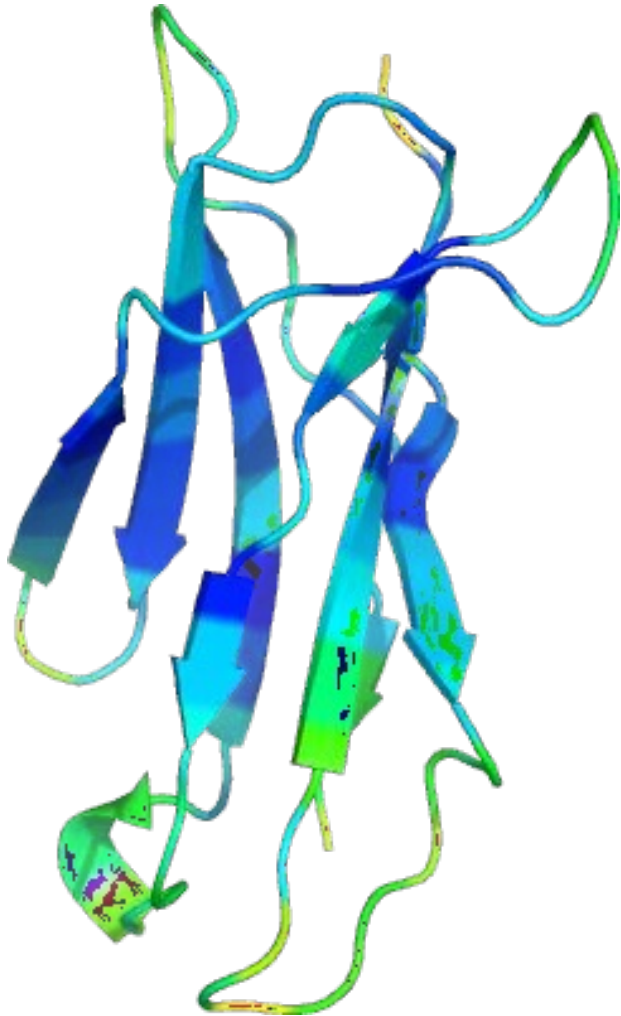# The interpretation of dynamic loops in the pdb file may be tentative



Well defined β-strand in the core of a protein: atomic positions are reliable

Flexible loop at the surface of a protein: atomic positions are not well defined

# Look at B-factor distribution!

**Protein coloured by B-factor**:

Well defined regions have low B-factors (blue/green)

Poorly defined/more mobile regions have high B-factors ( yellow/orange/red)

# A protein molecule is dynamic

- The electron density is a spatial average over all molecules in the crystal and a time average over the duration of the X-ray data measurement

- Multiple discrete conformations of a residue in different molecules are superimposed.

- Damage caused by X-rays may change the protein (mainly breaking of disulfide bonds)

- A crude description of dynamics is provided in the pdb file as the isotropic "B-factor"

- Some dynamical aspects evident in the electron density are lost in the pdb file

# Reading a crystallography paper:

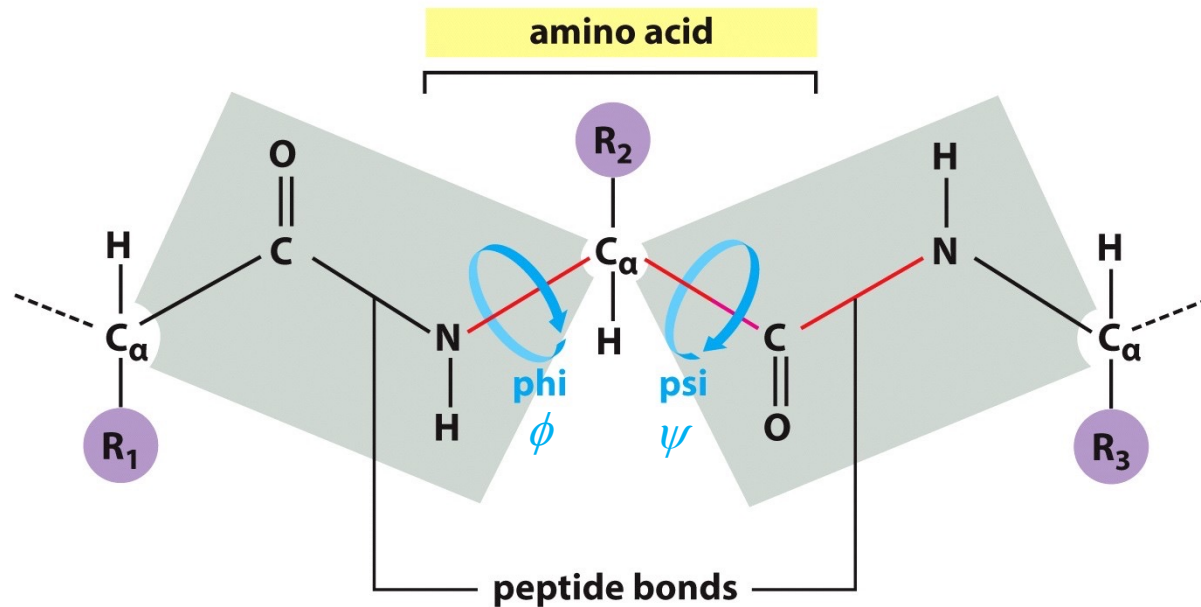**Table 1. Crystallographic Data Collection and Refinement Statistics**

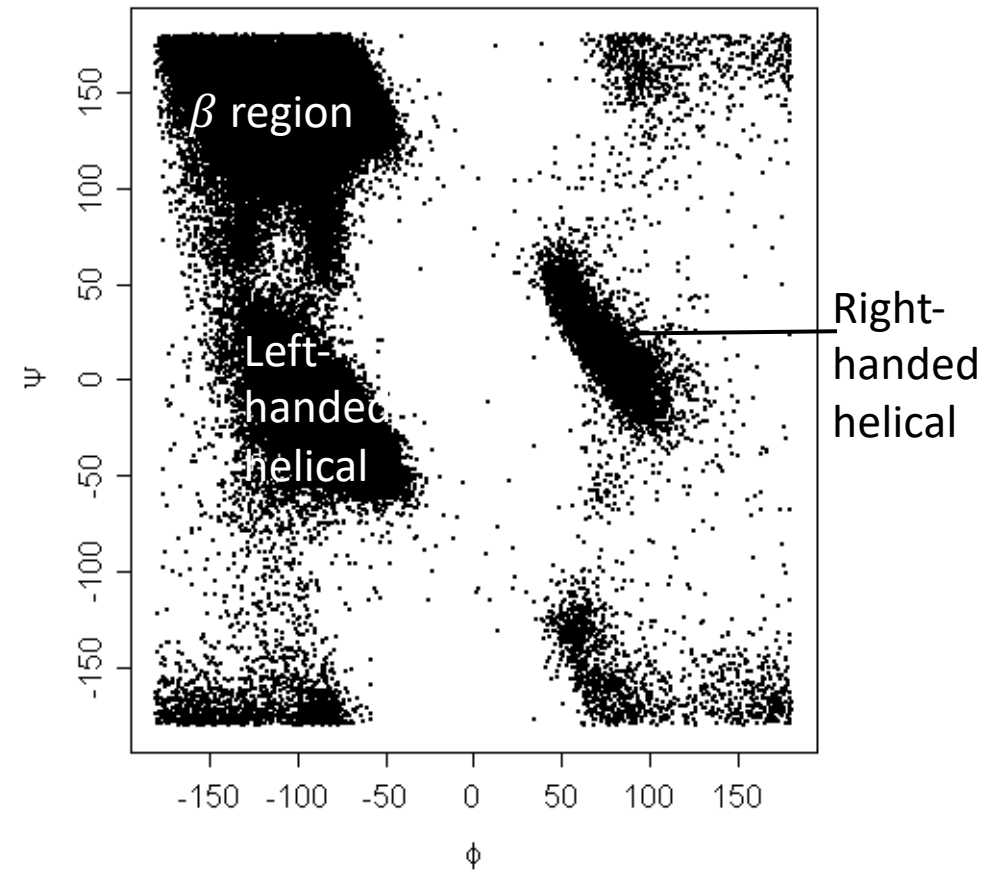| | SwMppP | SwMppP·D-Arg |
|---|---|---|
| resolution (Å) (last shell)[a] | 41.45−2.10 (2.14−2.10) | 44.53−2.25 (2.29−2.25) |
| wavelength (Å) | 0.97852 | 0.97856 |
| no. of reflections | | |
| observed | 1532933 (60186) | 256371 (12494) |
| unique | 106188 (5236) | 79816 (3984) |
| completeness (%)[a] | 100.0 (100.0) | 90.2 (91.2) |
| $R_{merge}$ (%)[a,b] | 0.106 (0.734) | 0.106 (0.726) |
| multiplicity | 14.4 (11.5) | 3.2 (3.1) |
| $\langle I/\sigma(I)\rangle^a$ | 30.7 (5.2) | 10.0 (1.9) |
| Model Refinement | | |
| no. of reflections in the working set | 100876 | 76747 |
| no. of reflections in the test set | 5256 | 3015 |
| $R_{cryst}$ ($R_{free}$) | 0.148 (0.177) | 0.162 (0.197) |
| no. of residues | 1417 | 1404 |
| no. of solvent atoms | 903 | 693 |
| no. of TLS groups | 29 | 33 |
| average B factor (Å²)[c] | | |
| protein atoms | 32.0 | 34.1 |
| ligands | 30.9[d] | 43.9[d] |
| solvent | 35.6 | 36.6 |
| root-mean-square deviation | | |
| bond lengths (Å) | 0.013 | 0.015 |
| bond angles (deg) | 1.395 | 1.565 |
| coordinate error (Å) | 0.17 | 0.22 |
| Ramachandran statistics (favored/allowed/outliers) (%) | 98.3/1.7/0 | 98.4/1.6/0 |

[a]Values in parentheses apply to the high-resolution shell indicated in the resolution row. [b]$R = \sum(||F_{obs}| - \text{scale} \times |F_{calc}||)/\sum|F_{obs}|$. [c]Isotropic equivalent B factors, including the contribution from TLS refinement. [d]In the unliganded SwMppP structure, "ligands" refers to the bound Cl ions, while in the D-Arg complex structure, it refers to the D-Arg-PLP unit.

**Judge the quality of the data:**

- $R_{merge}$: 0.05-0.10 good, 0.1-0.15 acceptable
- I/σ = signal/noise >2.0
- Completeness
- Redundancy
- $R_{work}$/$R_{free}$:
    - difference < 0.05,
    - $R_{work}$≈ resolution/10
- Deviations of known geometry

# Ramachandran Plot

- shows frequency of ($\phi$, $\psi$) observed for residues in folded proteins



The torsional angles of each residue define the geometry of its attachment to its two adjacent residues, so the torsional angles determine the conformation of the residues and the peptide.
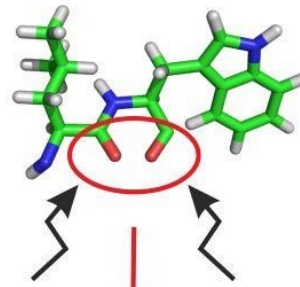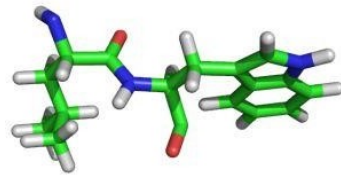


Each dot represents an amino acid

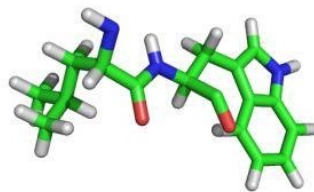The **green/yellow regions** correspond to conformations where there are no steric clashes

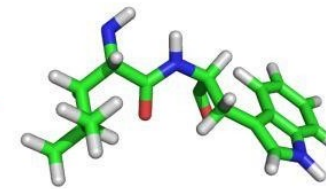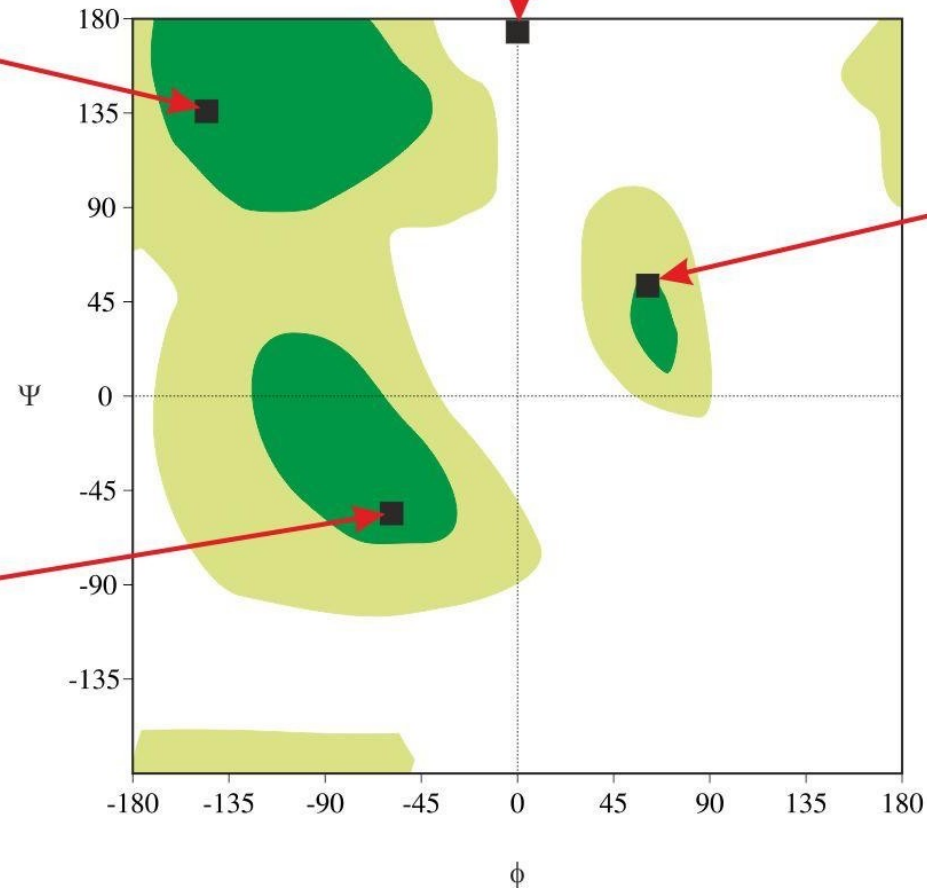**White regions**: sterically disallowed for all amino acids except glycine

steric distortion

antiparallel β-sheet

left-handed α-helix

right-handed α-helix

Ψ

φ