

Determination of Biological Structures

Workshop on Computer-Aided Drug Discovery
Graz, 5.9.2022

Karl Gruber
Institute of Molecular Biosciences, University of Graz

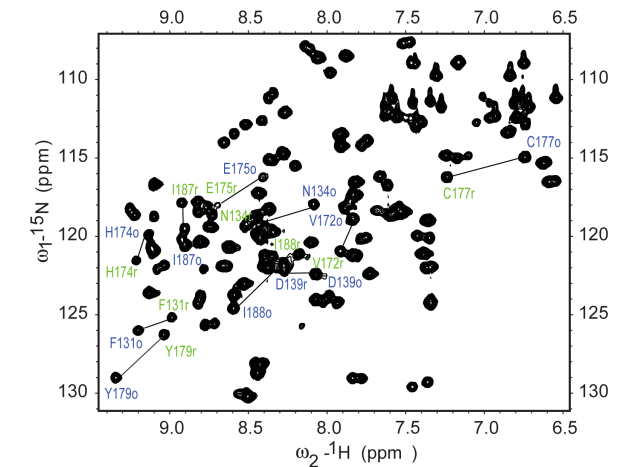
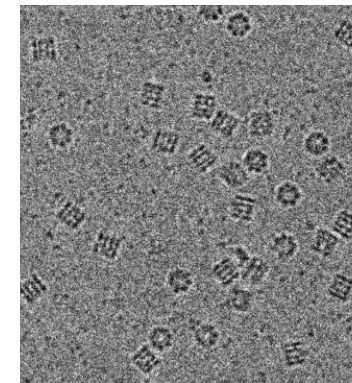
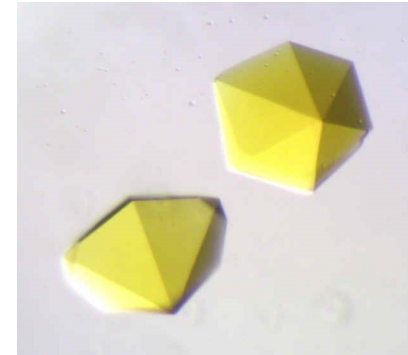
Experimental Determination of Molecular Biological Structures and their Validation

Workshop on Computer-Aided Drug Discovery
Graz, 5.9.2022

Karl Gruber
Institute of Molecular Biosciences, University of Graz

Structure determination of biological macromolecules

- X-ray crystallography
- NMR-spectroscopy
- Cryo-electron microscopy



Macromolecular structure databases

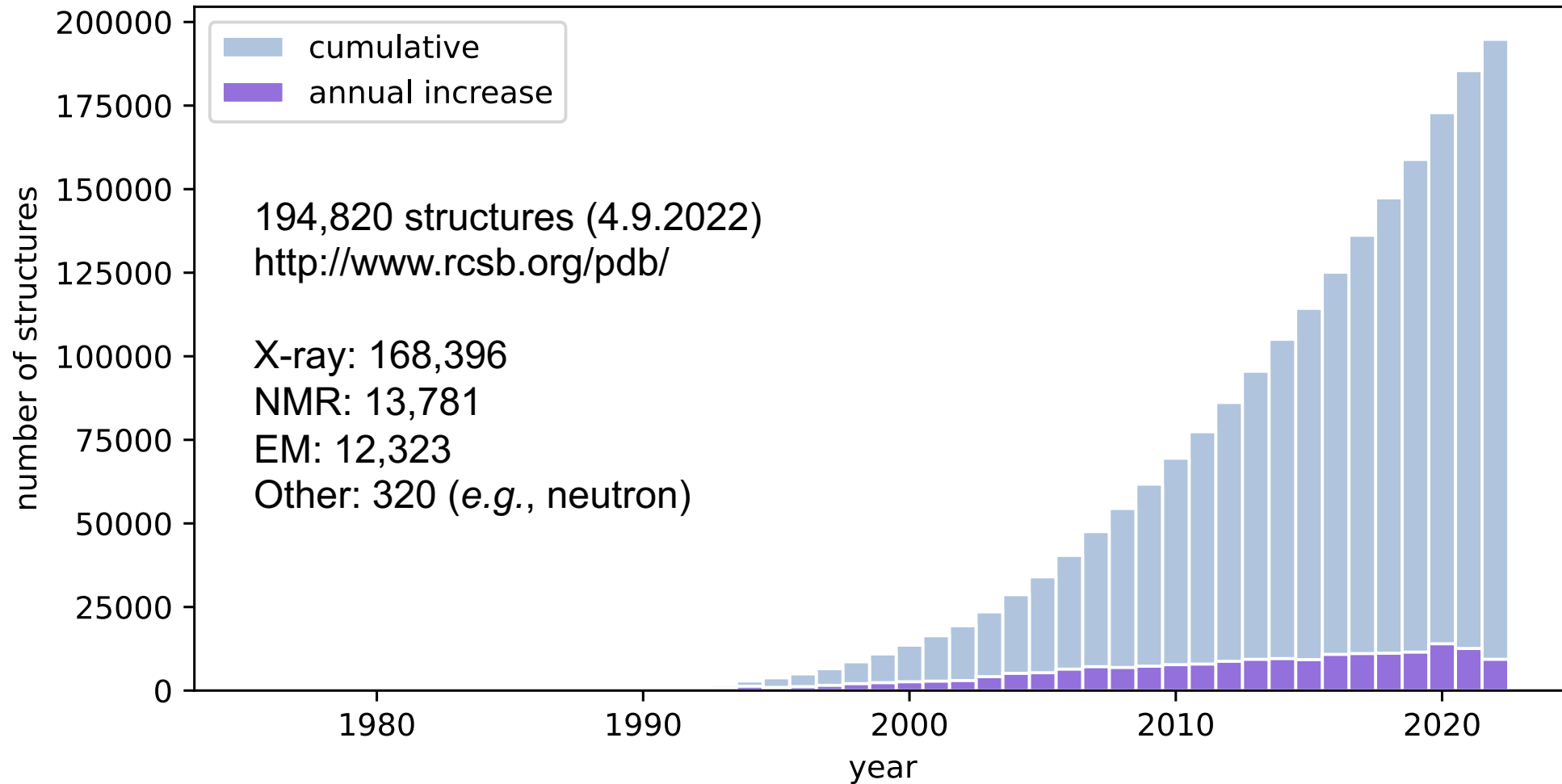


<http://www.rcsb.org>

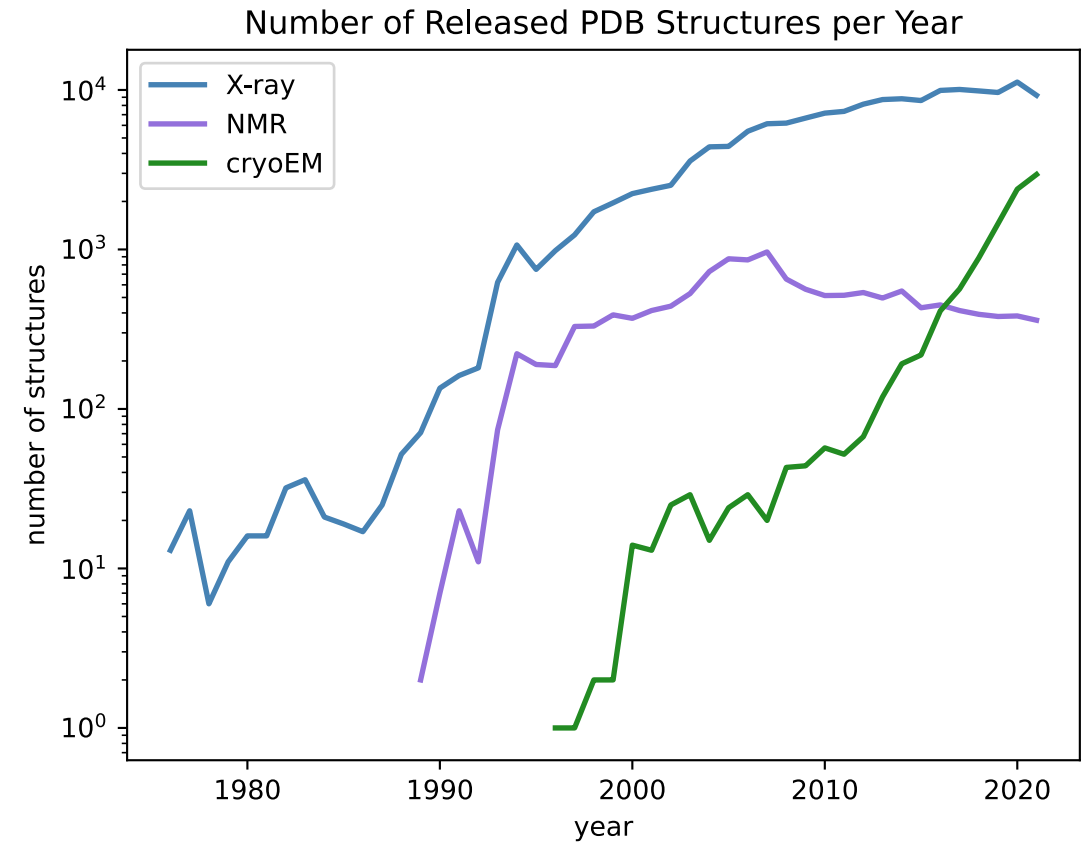
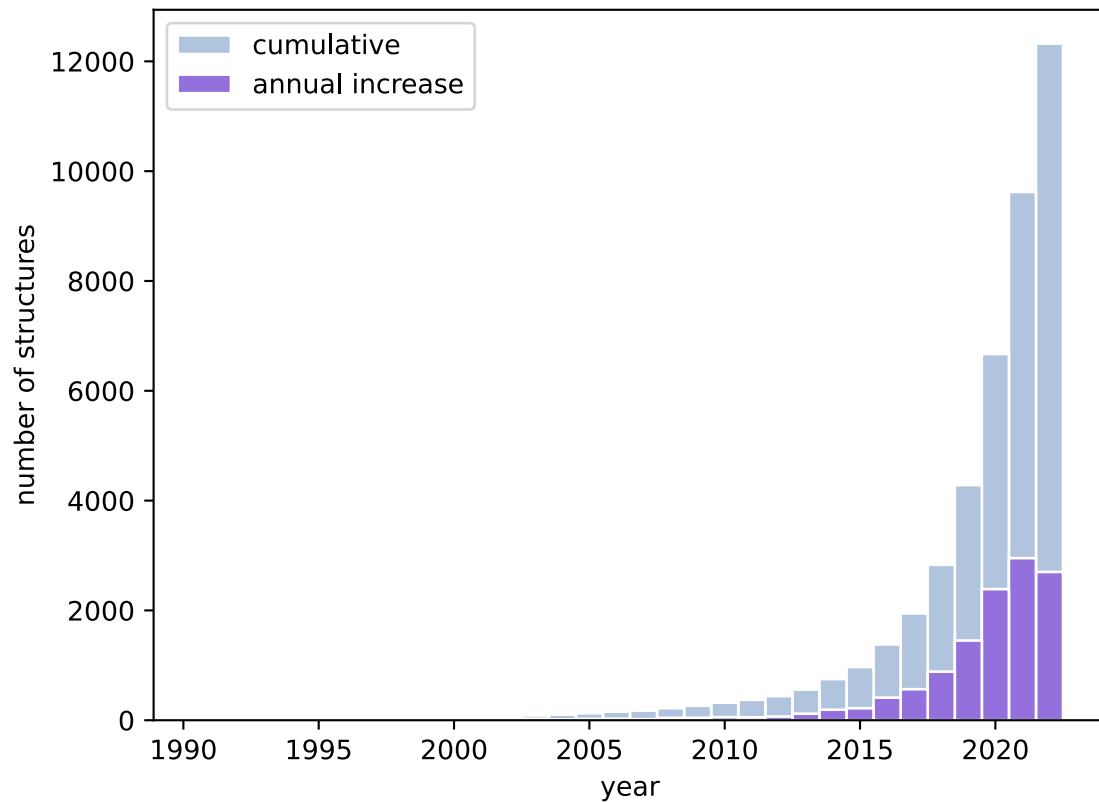


<http://www.ebi.ac.uk/pdbe/>

Protein Data Bank (PDB)



Cryo-electron microscopy (cryoEM)

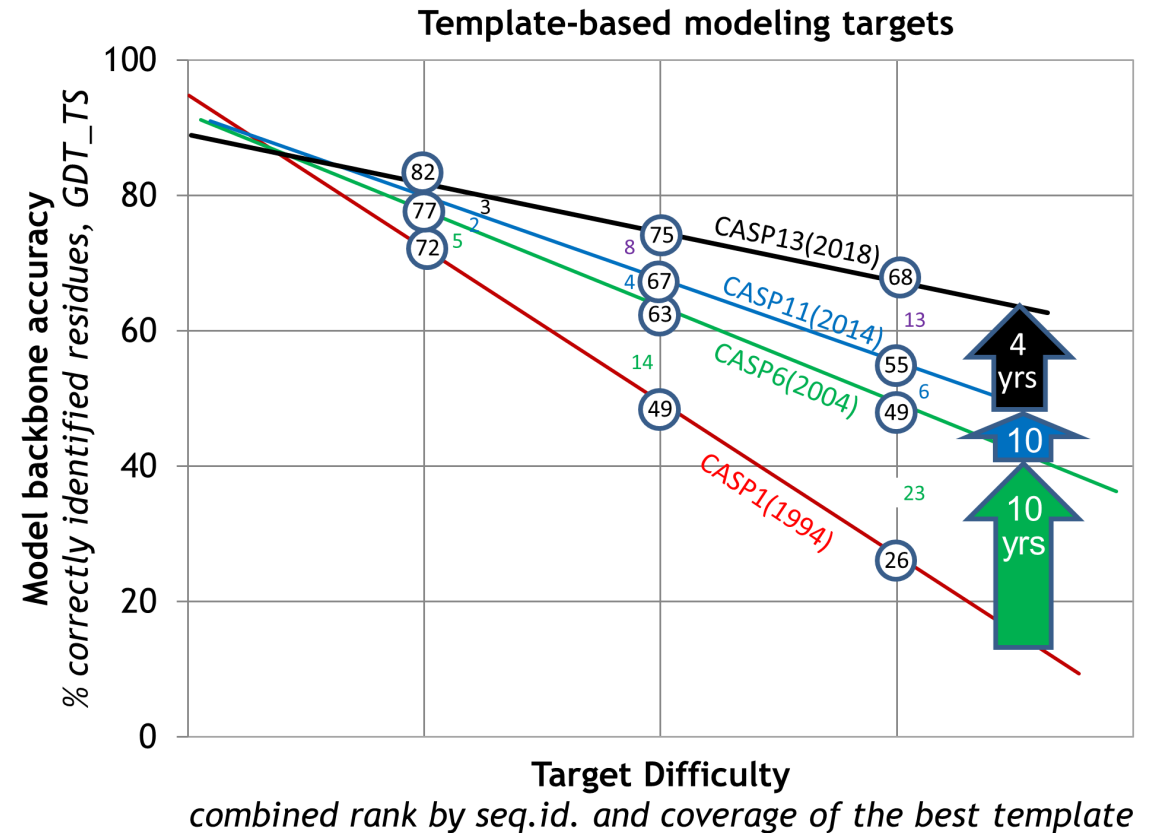


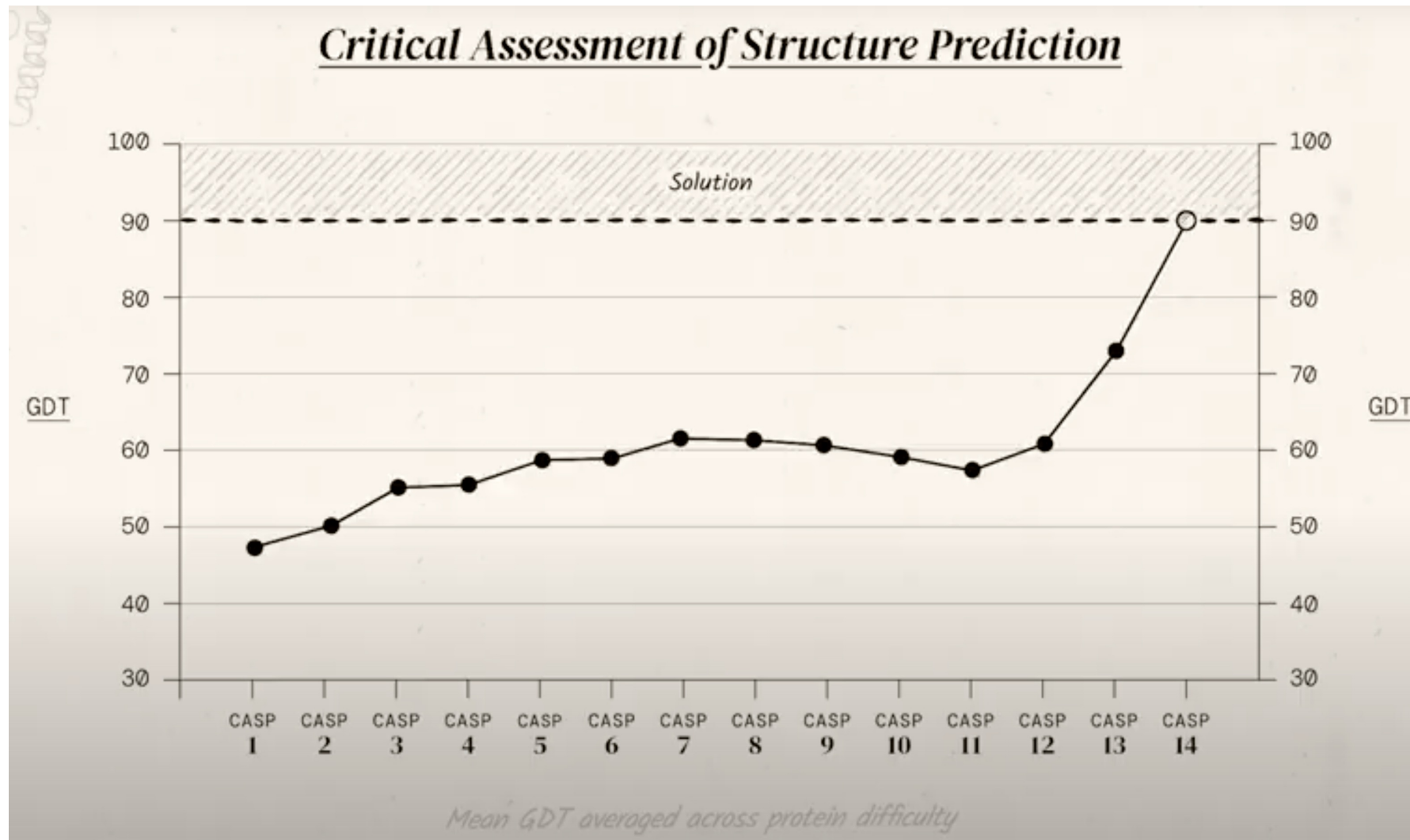


CASP – Critical assessment of protein structure prediction



- biennial competition for protein structure prediction taking place since 1994
- **targets:** structures solved experimentally but not yet published
- **contestants:** expert groups as well as automatic prediction servers
- benchmark of the quality of prediction-algorithms
- upcoming CASP-15, Dec. 2022

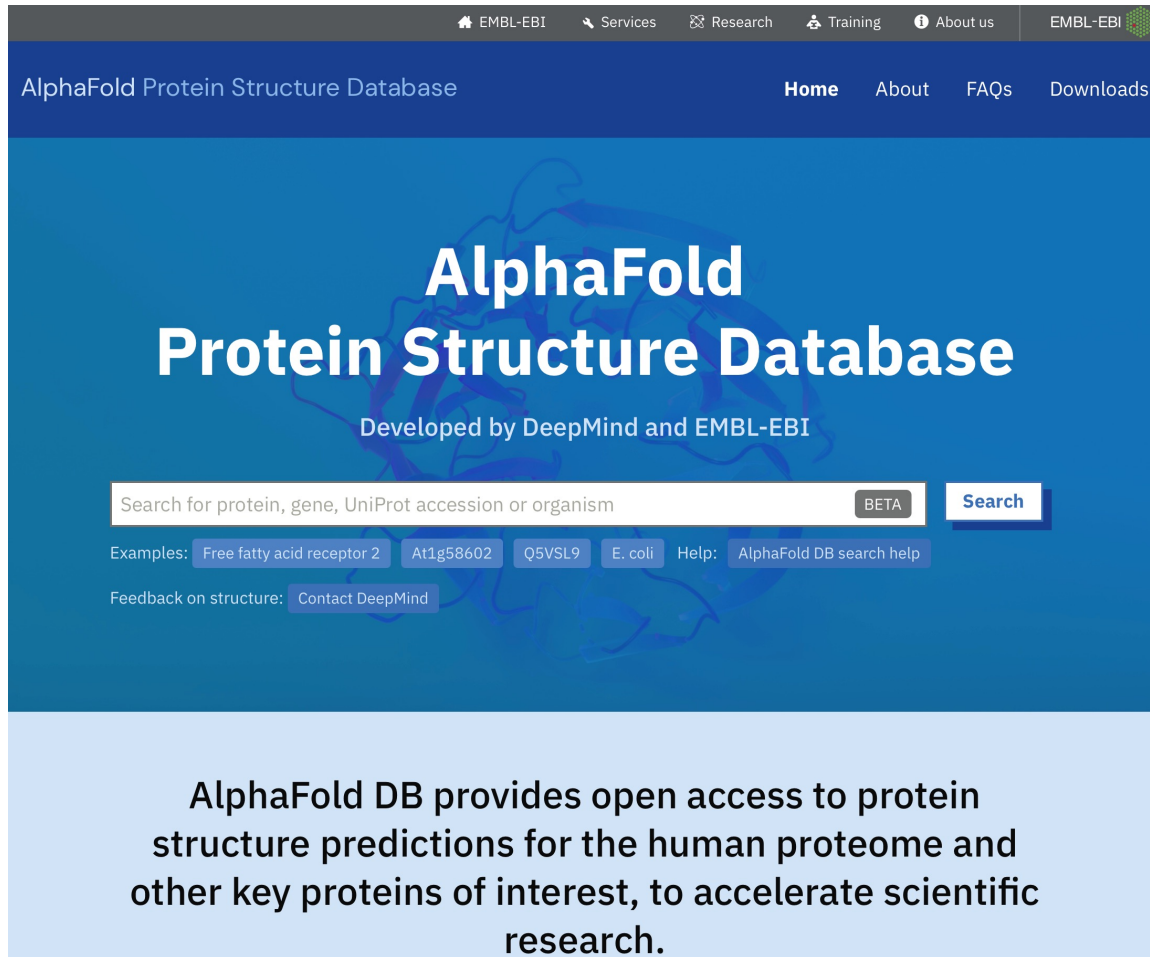




AlphaFold: The making of a scientific breakthrough

<https://www.youtube.com/watch?v=gg7WjuFs8F4>

Database of AlphaFold models



<https://www.alphafold.ebi.ac.uk>

- precomputed models of protein structures
- whole proteomes: human, mouse,...
- complete Swissprot database (curated protein sequences)
- quality parameters useful for validation



Science magazine:
Breakthrough of the Year 2021

<https://www.science.org/content/article/breakthrough-2021>


Run AlphaFold yourself

- GoogleColab notebooks: based on Jupyter Notebooks, Python
- runs on the Google cloud (free service, “restrictions apply”)
- can be installed locally (requires larger GPU)

<https://github.com/sokrypton/ColabFold>

README.md

ColabFold



New Updates

+ 11Mar2022 We use in default AlphaFold-multimer-v2 weights for complex modeling.

+ We also offer the old complex modes "AlphaFold-ptm" or "AlphaFold-multimer-v1"

+ 04Mar2022 ColabFold now uses a much more powerful server for MSAs and searches through the Cola

+ Please let us know if you observe any issues.

+ 26Jan2022 AlphaFold2_mmseqs2, AlphaFold2_batch and colabfold_batch's multimer complexes predict

+ now in default reranked by iptmscore*0.8+ptmscore*0.2 instead of ptmscore

Making Protein folding accessible to all via Google Colab!

Notebooks	monomers	complexes	mmseqs2	jackhmmer	templates
AlphaFold2_mmseqs2	Yes	Yes	Yes	No	Yes
AlphaFold2_batch	Yes	Yes	Yes	No	Yes
RoseTTAFold	Yes	No	Yes	No	No
AlphaFold2 (from Deepmind)	Yes	Yes	No	Yes	No

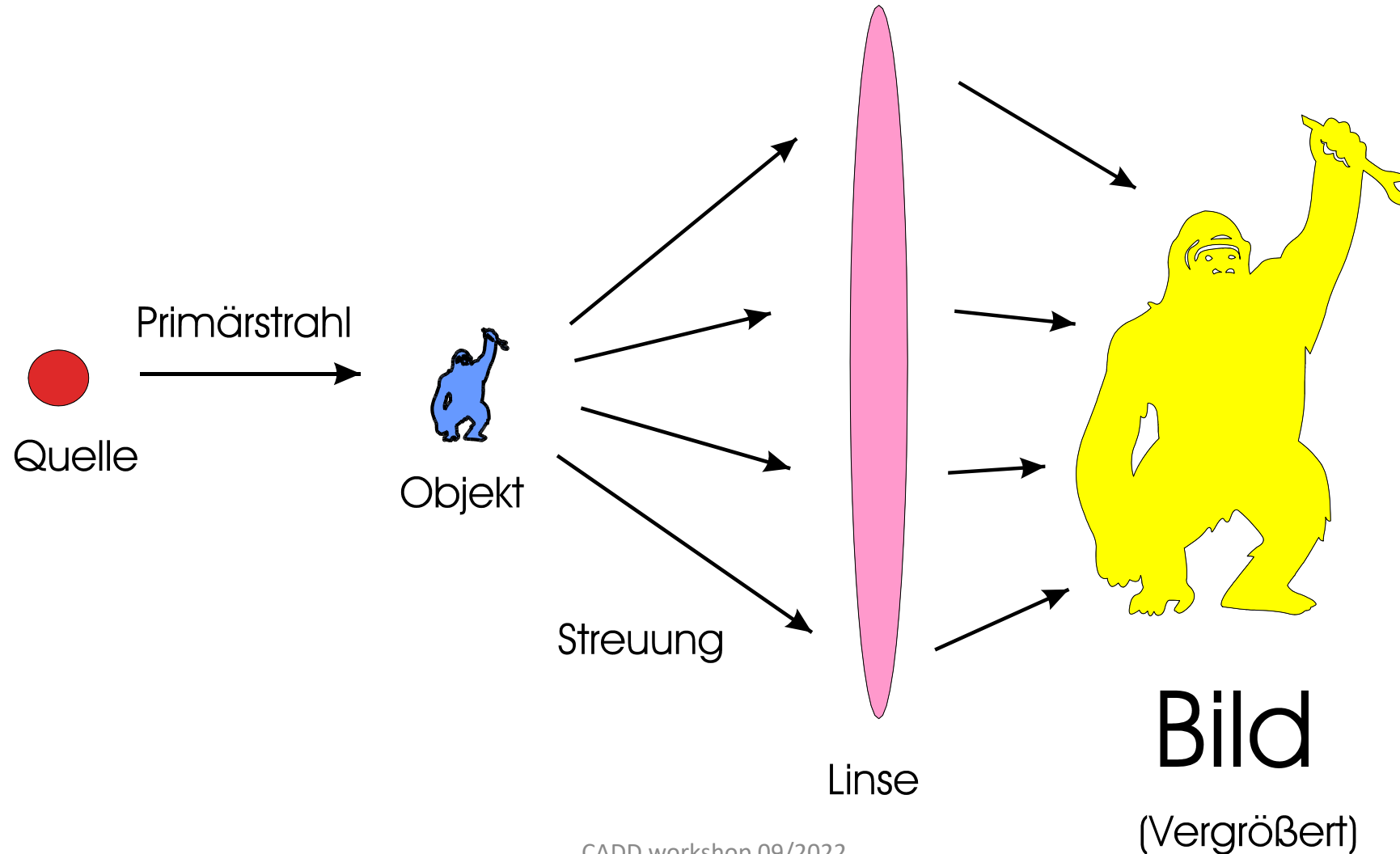
Ultimate aim

Determination of the 3-dimensional structure of a (macro)molecule with **atomic resolution**.

Know the positions of all atoms in the molecule. (At least have a reasonable idea of their average positions.)

“See” the atoms in the molecule as separate entities.

Principle of optical imaging (microscope)



Limits of resolution

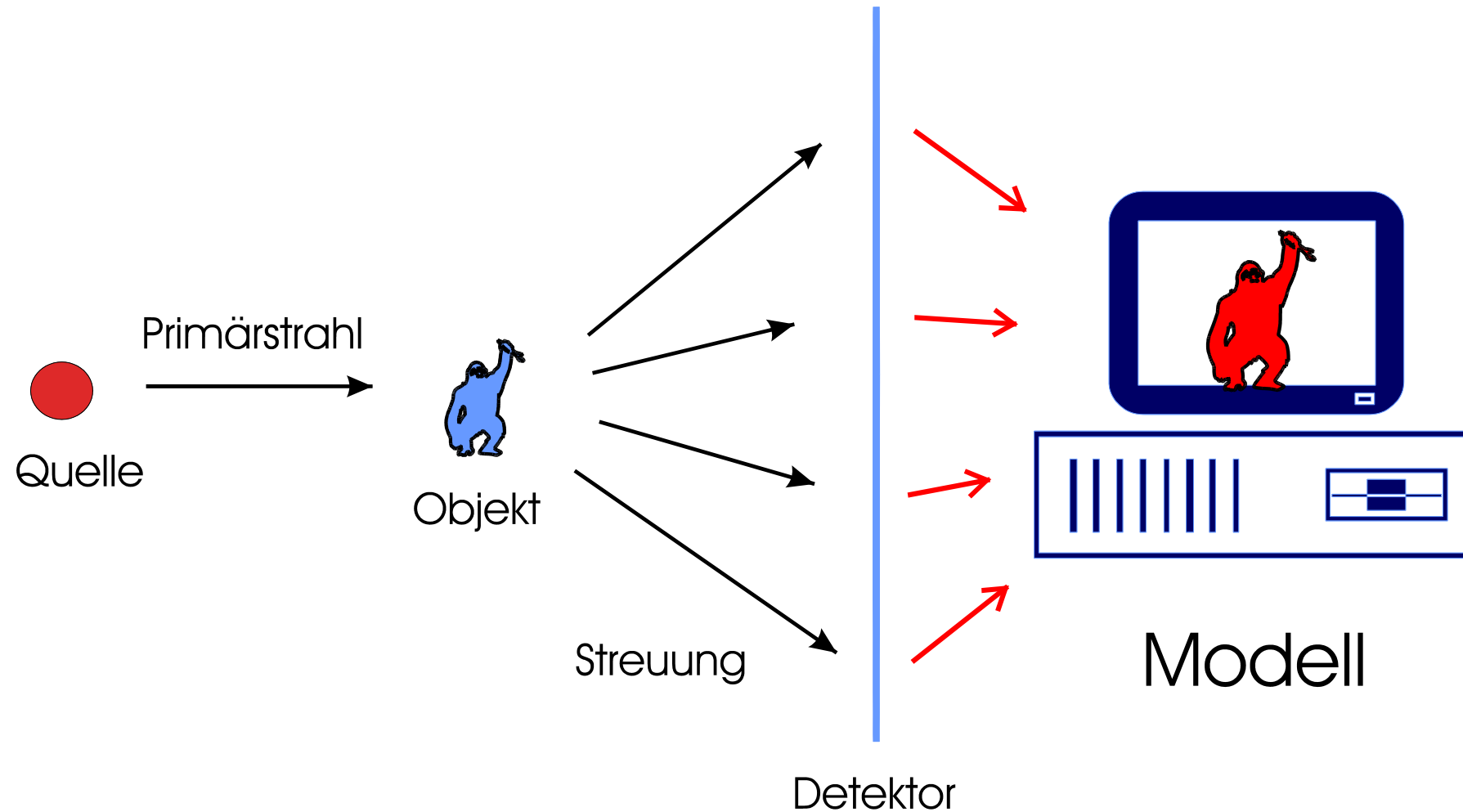
$$d_{\min} = \frac{\lambda}{2 \sin \vartheta_{\max}}$$

d_{\min} ... smallest distance between two points that can be imaged separately

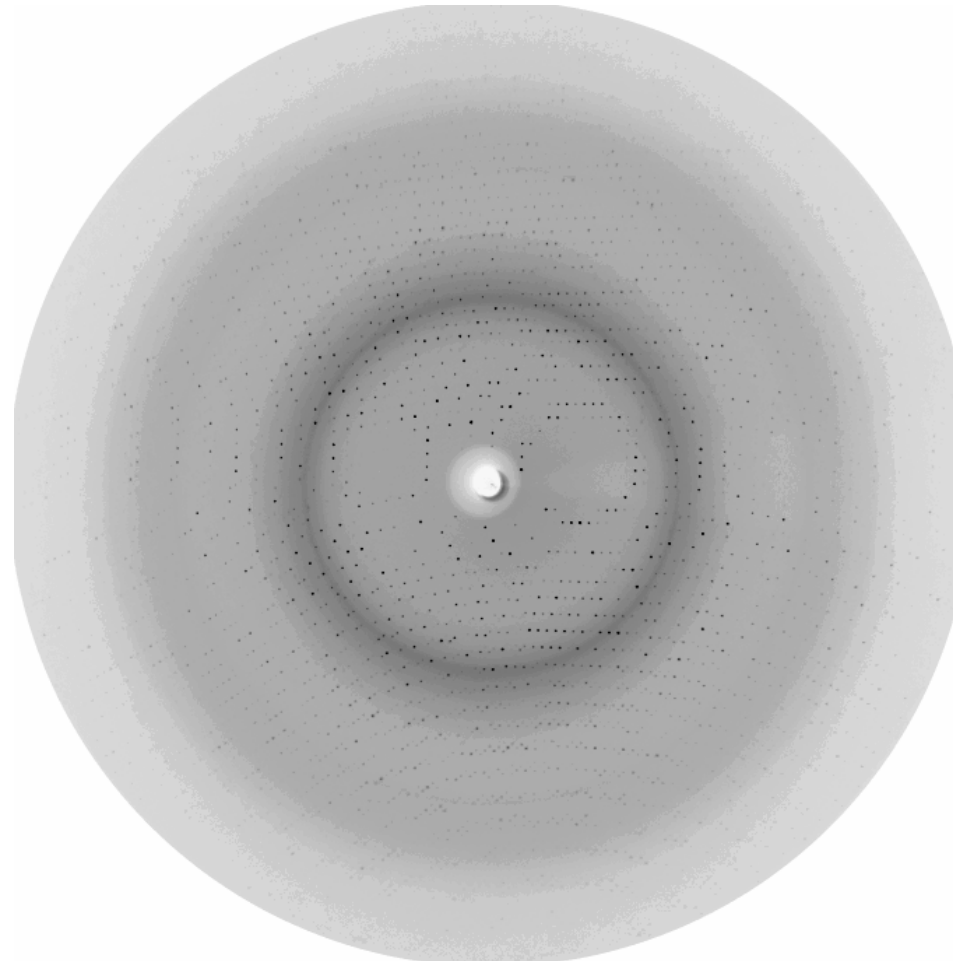
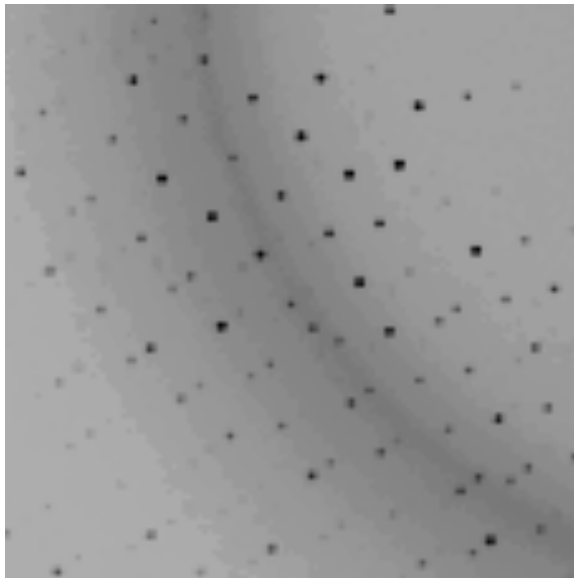
λ ... wavelength

ϑ_{\max} ... maximum diffraction angle

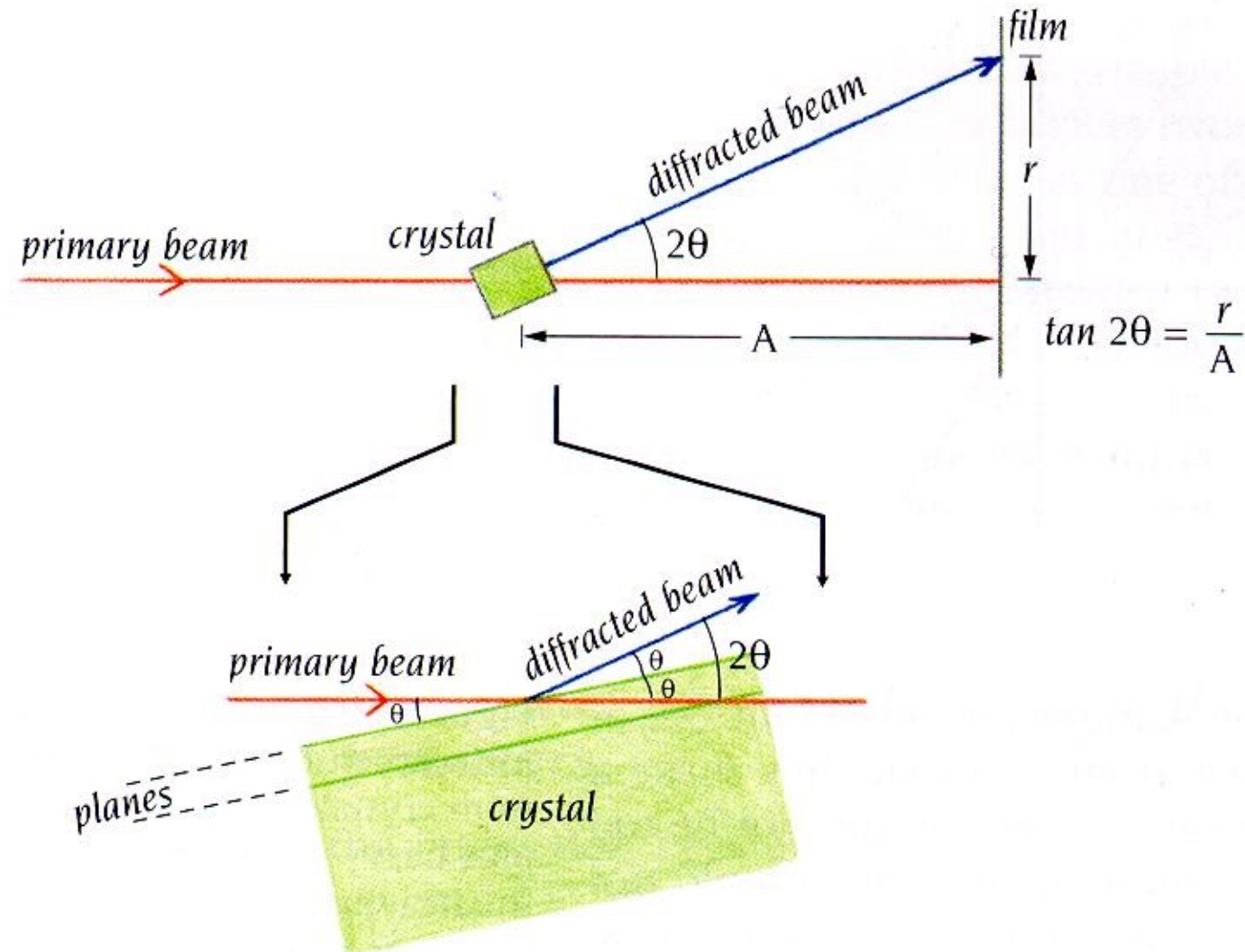
Principle of a diffraction experiment



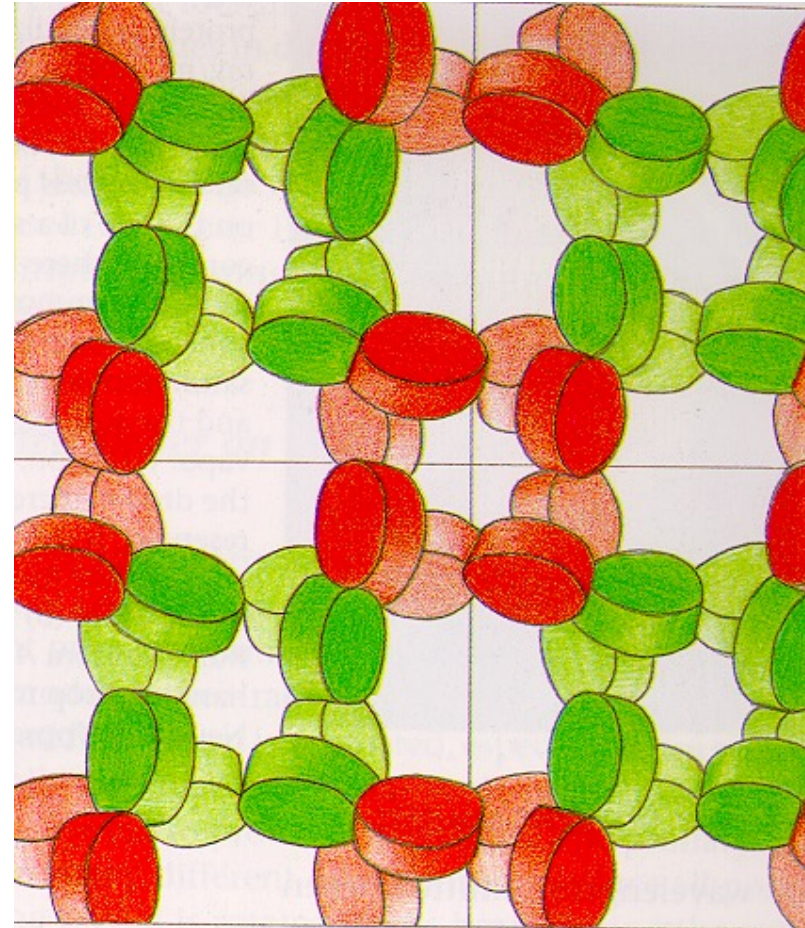
Diffraction image (protein crystal)



Why are the spots called “reflections”?

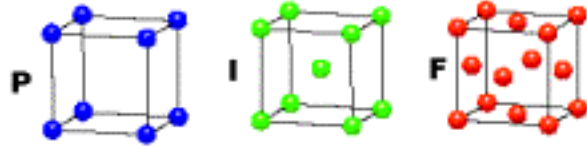


Crystals of biological objects

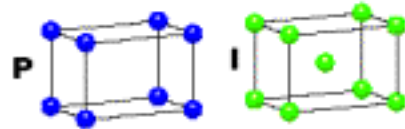


Crystal classes, crystal symmetry

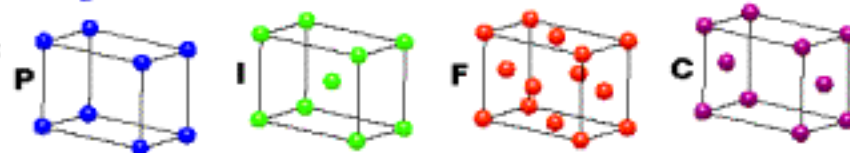
CUBIC
 $a = b = c$
 $\alpha = \beta = \gamma = 90^\circ$



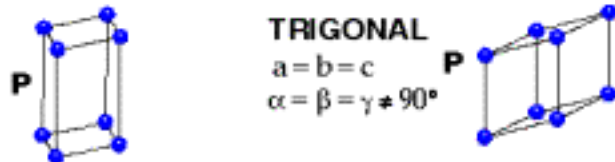
TETRAGONAL
 $a = b \neq c$
 $\alpha = \beta = \gamma = 90^\circ$



ORTHORHOMBIC
 $a \neq b \neq c$
 $\alpha = \beta = \gamma = 90^\circ$

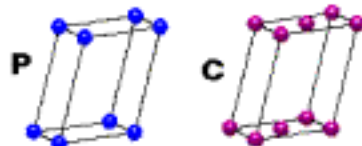


HEXAGONAL
 $a = b \neq c$
 $\alpha = \beta = 90^\circ$
 $\gamma = 120^\circ$

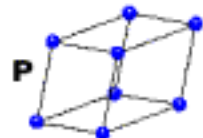


TRIGONAL
 $a = b = c$
 $\alpha = \beta = \gamma \neq 90^\circ$

MONOCLINIC
 $a \neq b \neq c$
 $\alpha = \gamma = 90^\circ$
 $\beta \neq 120^\circ$



TRICLINIC
 $a \neq b \neq c$
 $\alpha \neq \beta \neq \gamma \neq 90^\circ$



4 Types of Unit Cell
P = Primitive
I = Body-Centred
F = Face-Centred
C = Side-Centred
+
7 Crystal Classes
→ 14 Bravais Lattices

Space group symmetry elements:

symmetry axes: only 2-, 3-, 4-, and 6-fold rotations are possible!!

mirror planes

center of inversion

inversion axes: rotation plus center of inversion

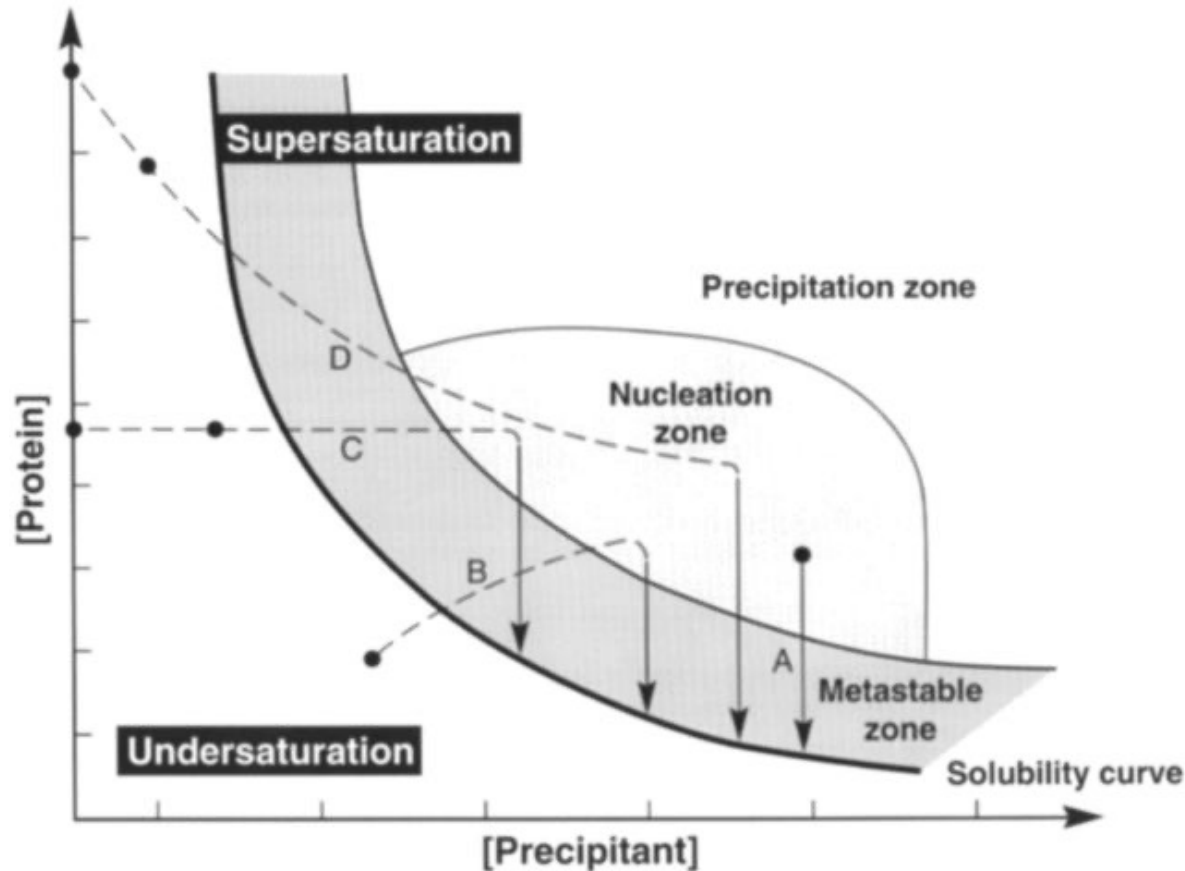
screw axes: rotation plus translation parallel to the axis

glide planes: mirror plane plus translation parallel to the plan

Workflow of an X-ray crystal structure determination

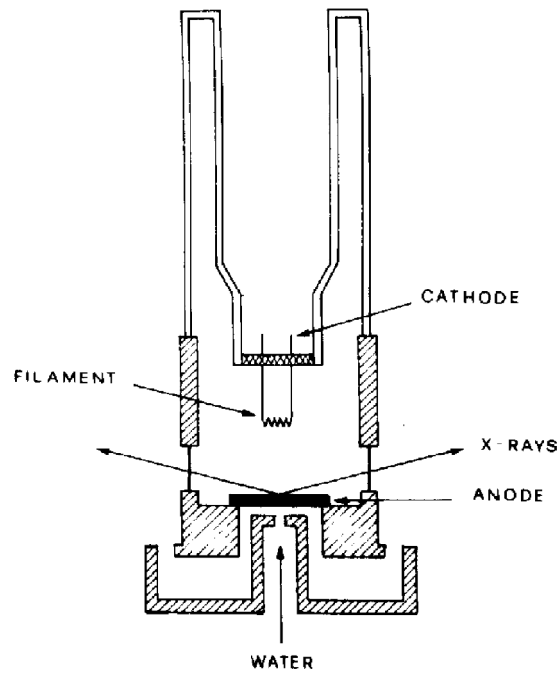
- Preparation and purification (of the protein)
- Crystallization
- Collection of a diffraction dataset
- Solution of the phase problem
- Structure refinement, validation
- Interpretation, publication of the results,...

Crystallization



- **Crystallization = "controlled precipitation"**
- mostly trial and a lot of error/failure
- **important factors:**
 - purity, homogeneity of the sample
 - pH
 - temperature (value, consistency)
 - nature of the precipitating agent (salts, PEG,...)
- ...

Data collection, X-ray sources



Wilhelm Conrad Röntgen,
1895



Synchrotron radiation source (ESRF)

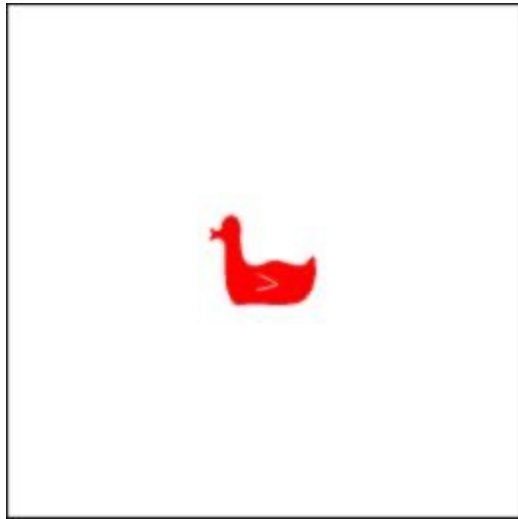
Phase problem

$$\rho(x, y, z) = \sum_{h,k,l} |F_{h,k,l}| e^{2\pi i \varphi_{h,k,l}} e^{-2\pi i (hx + ky + lz)}$$

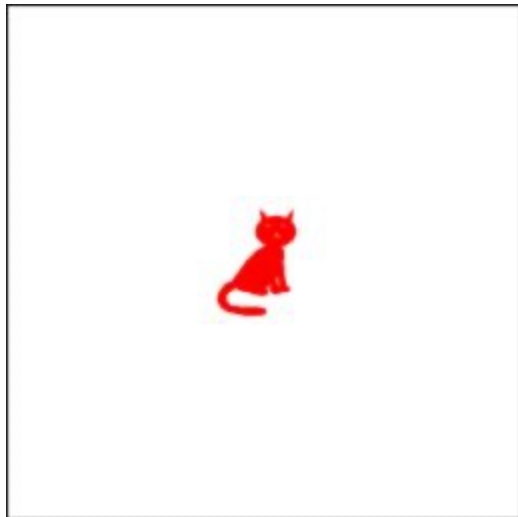
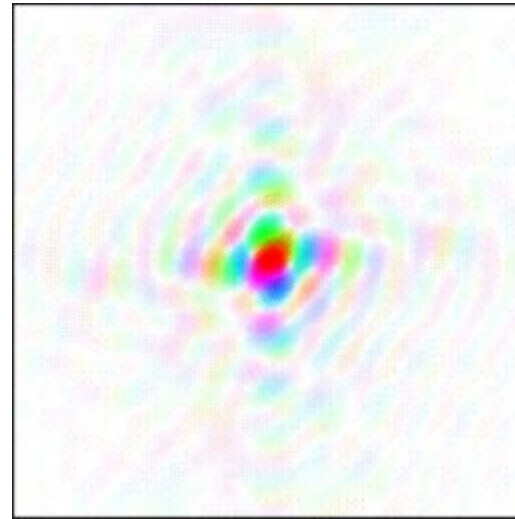
Calculation of the electron density in “real” space (x, y, z) using diffraction data (“reciprocal” space) by a Fourier transform.

Each reflection (h, k, l) is characterized by its amplitude (F) and its phase (φ), and **both are required** for calculating the electron density.

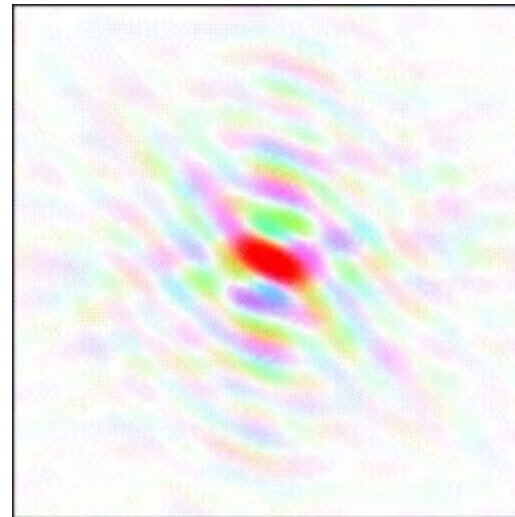
From the diffraction experiment, we only obtain the amplitudes (square roots of the intensities) but **no phase information**. Hence, the **phase problem**.



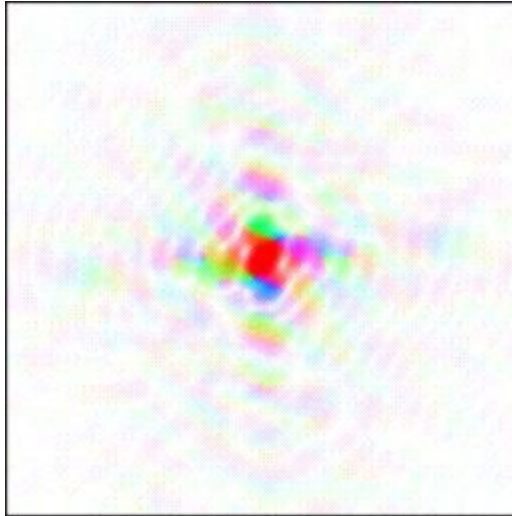
Fourier
transform



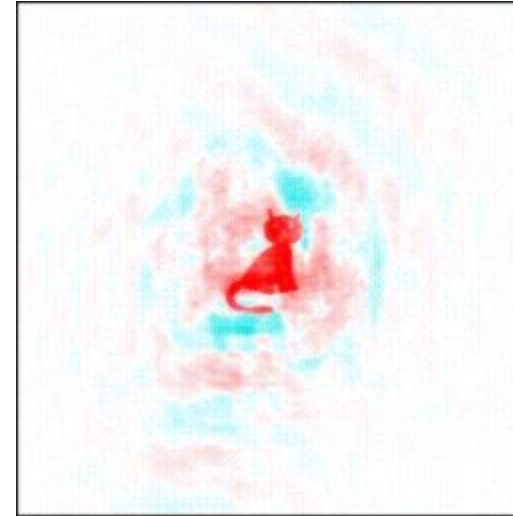
Fourier
transform



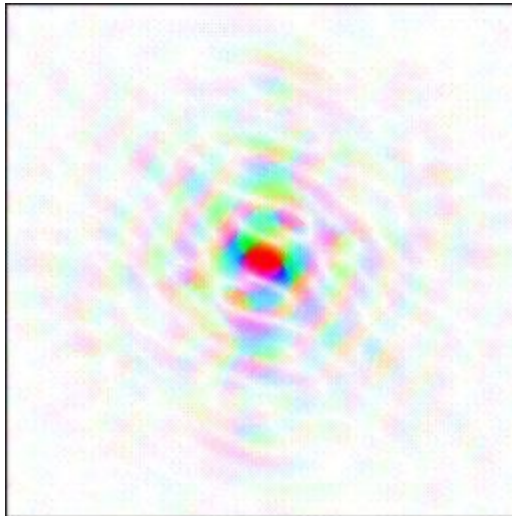
amplitudes: duck
phases: cat



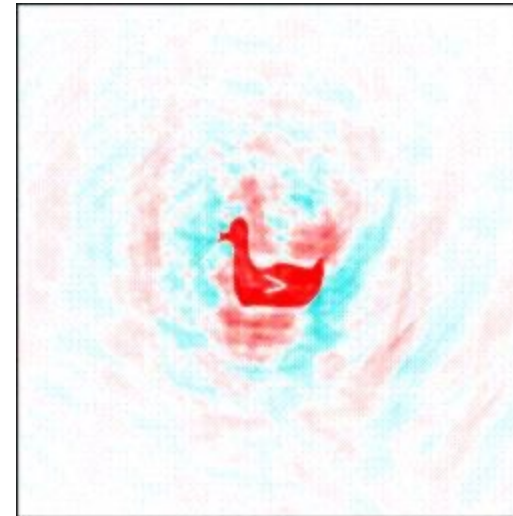
back FT



amplitudes: cat
phases: duck



back FT



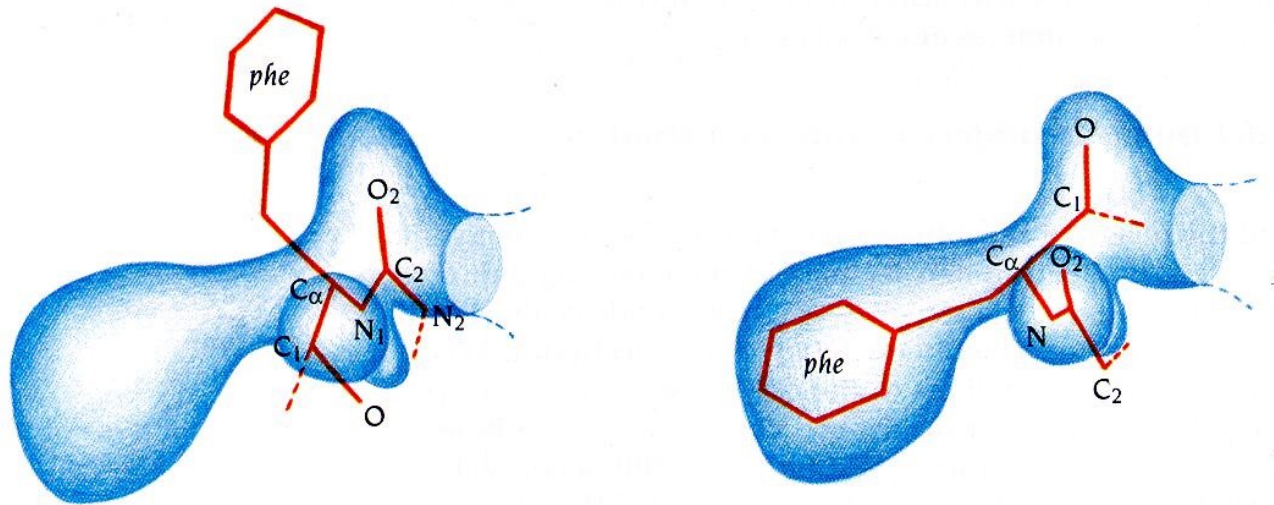
Solving the phase problem

- **MR** **M**olecular **R**eplacement
- **MIR** **M**ultiple isomorphous replacement
- **MAD** **M**ultiple wavelengths anomalous dispersion

$$\Rightarrow \varphi_{h,k,l}^{calc}$$

$$\rho(x,y,z) = \sum_{h,k,l} \left| F_{h,k,l}^{obs} \right| e^{2\pi i \varphi_{h,k,l}^{calc}} e^{-2\pi i(hx+ky+lz)}$$

Model building



We do not observe “molecular structure” directly!

The direct result, after solving the phase problem, is **electron density**. The electron-density map is then interpreted by fitting into it (pieces of) a polypeptide chain.

Individual atoms are not resolved at the resolution typically obtained in protein crystallography. Instead, there are **lumps of density** corresponding to **groups of atoms**.

Structure refinement

Structure factor calculation relies on the “crystallographic structure model”: the unit cell contains **N distinct atoms** at positions (x, y, z) with structure factors f.

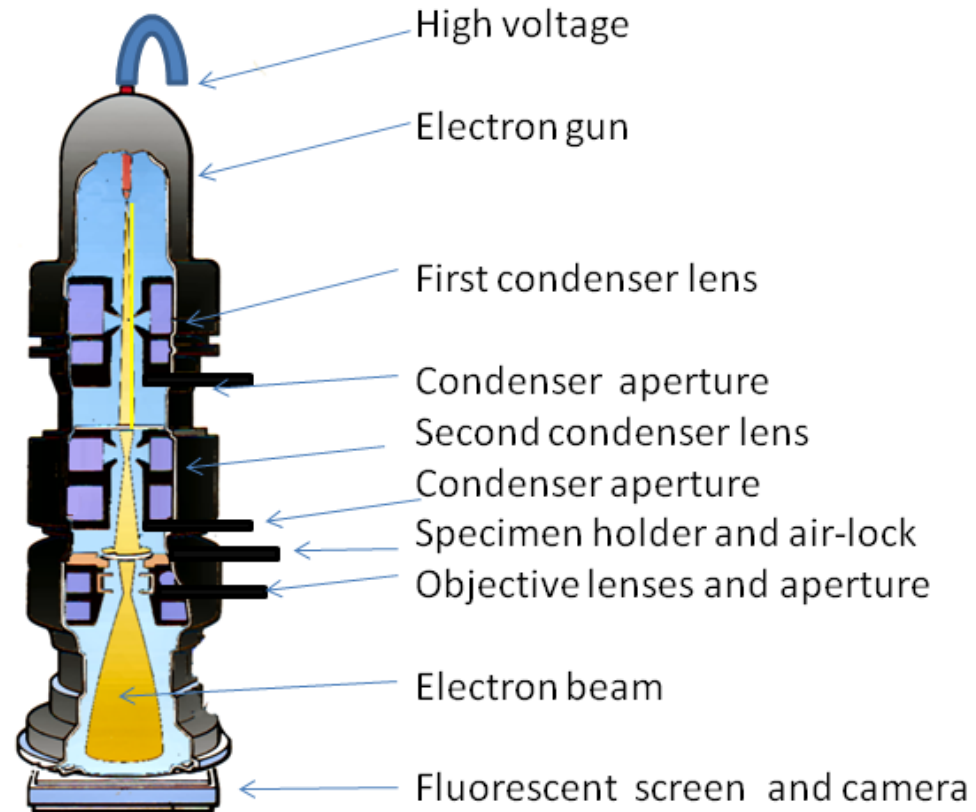
F^{calc} corresponds to the hypothetical diffraction data (**including phases**) we would obtain based on the current model.

The model is continuously adjusted during structure refinement to **minimize the differences** between observed (F^{obs}) and calculated (F^{calc}) reflection amplitudes.

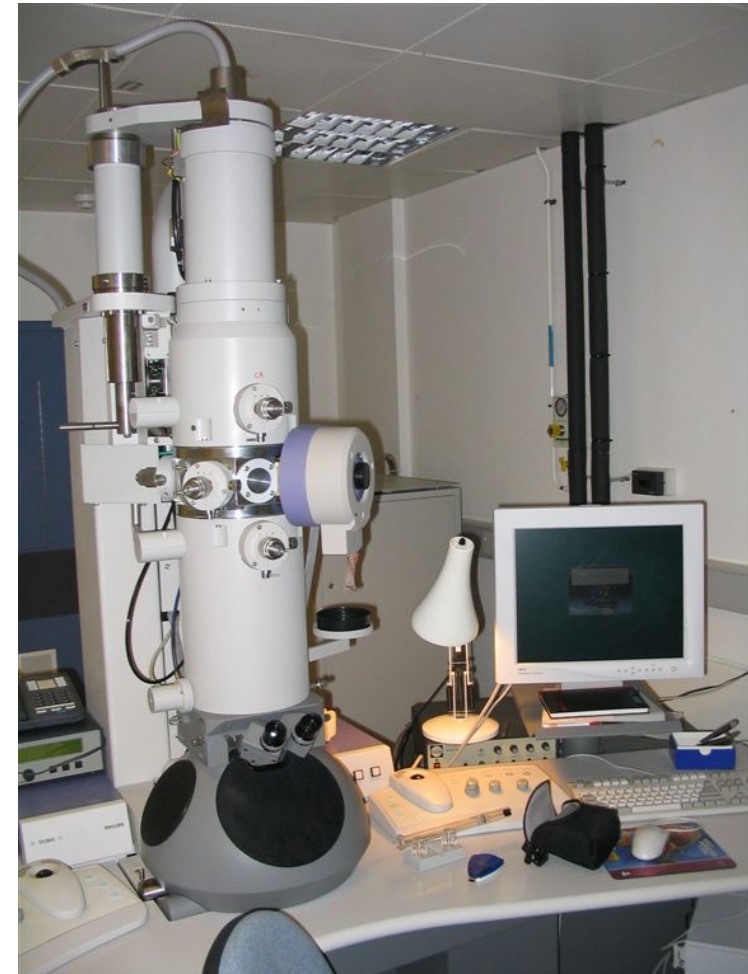
$$F_{h,k,l}^{calc} = \left| F_{h,k,l}^{calc} \right| e^{2\pi i \phi_{h,k,l}} = \sum_{i=1}^N f_i e^{2\pi i (hx_i + ky_i + lz_i)}$$

$$R = \frac{\sum_{h,k,l} \left| |F_{obs}(hkl)| - |F_{calc}(hkl)| \right|}{\sum_{h,k,l} |F_{obs}(hkl)|}$$

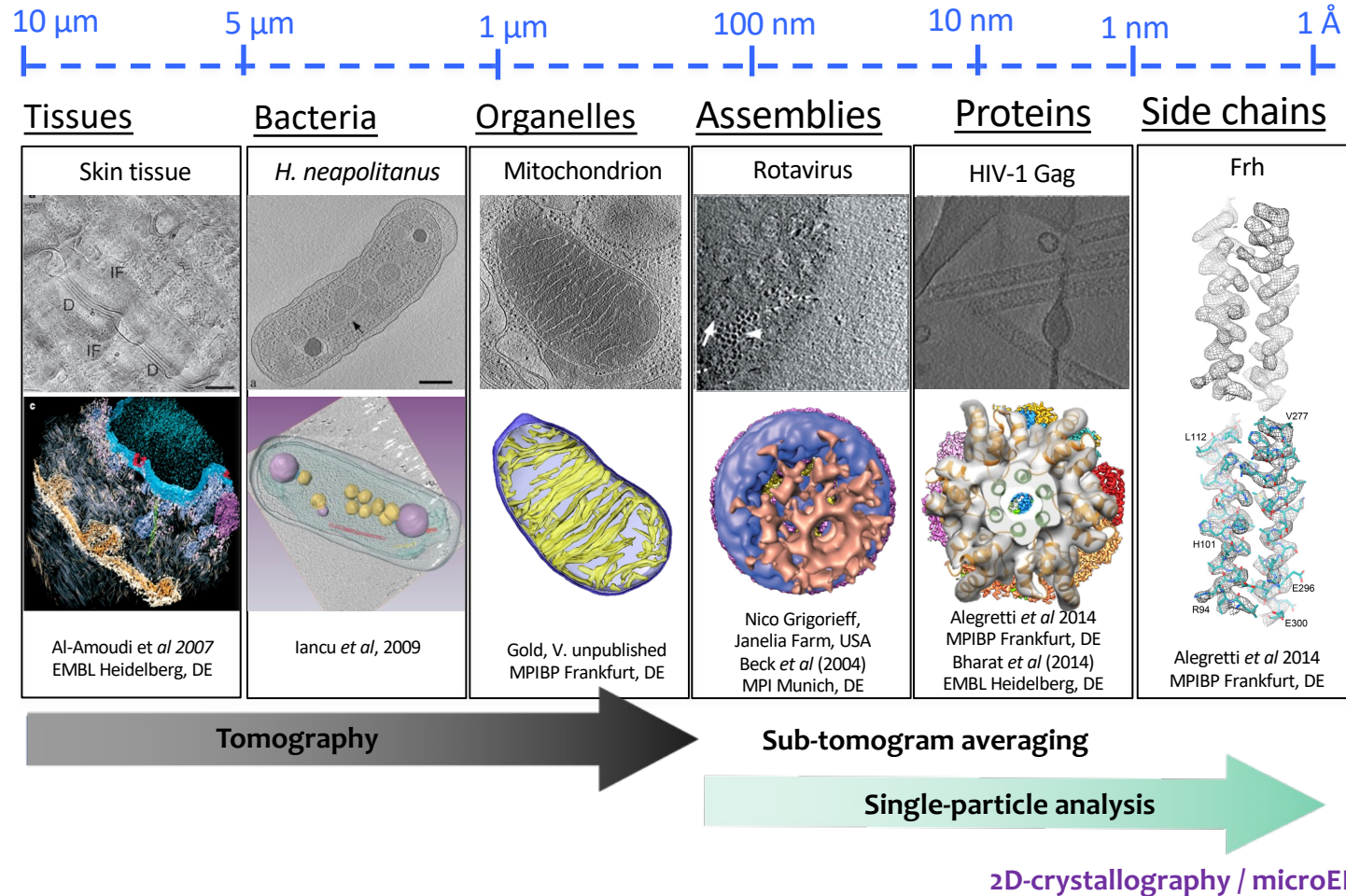
Electron microscopy



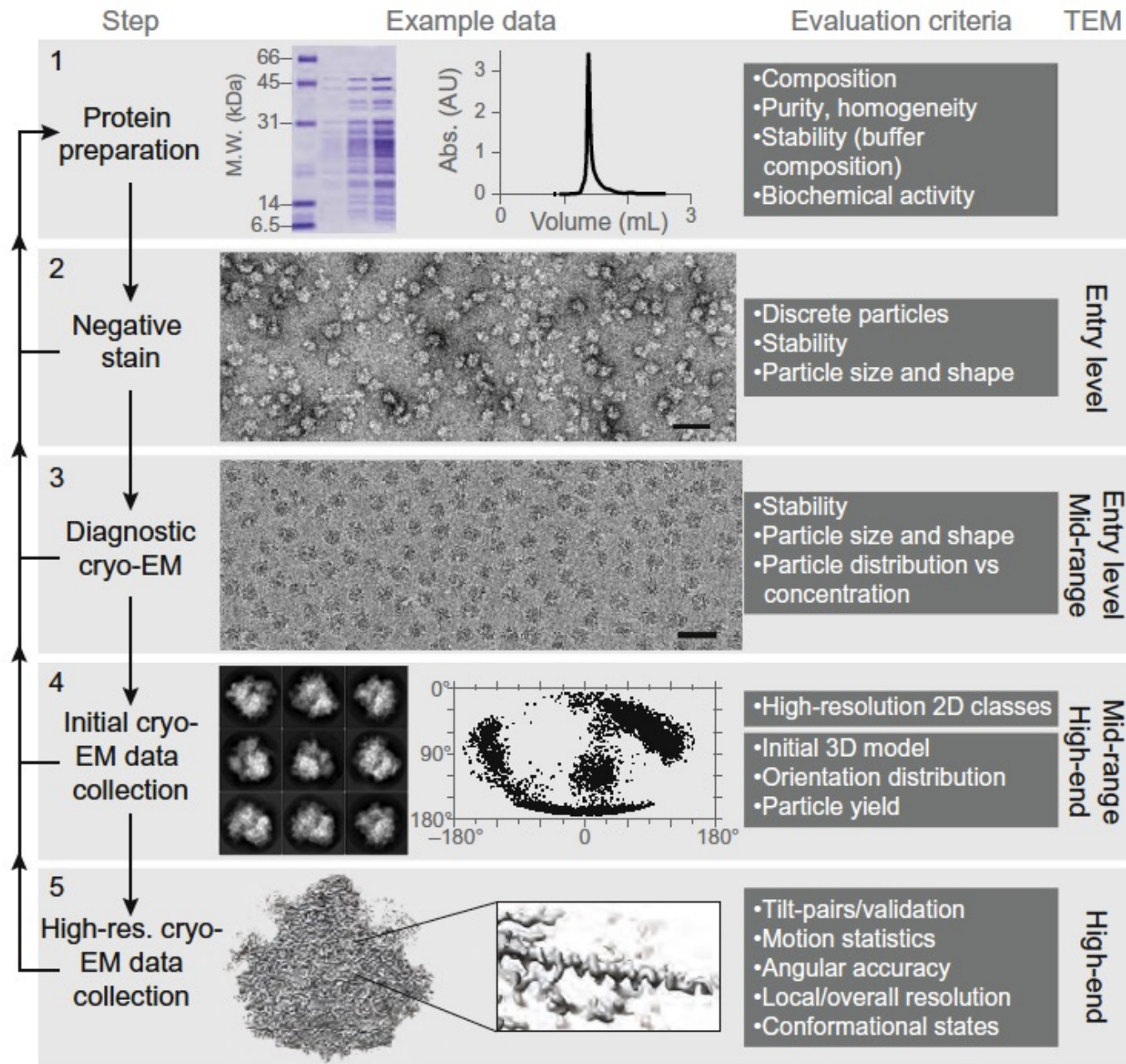
Transmission Electron Microscope



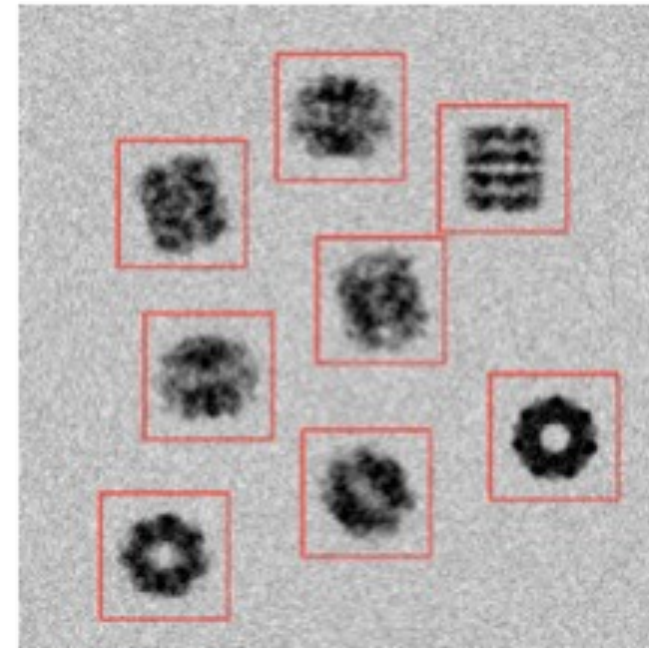
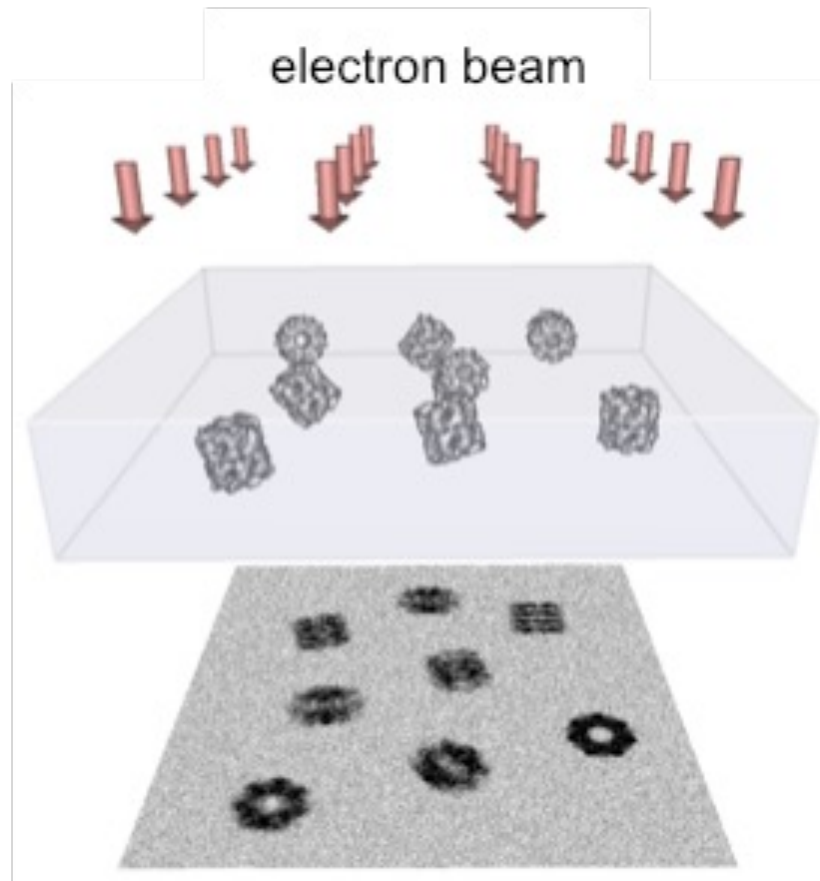
Applications of cryoEM



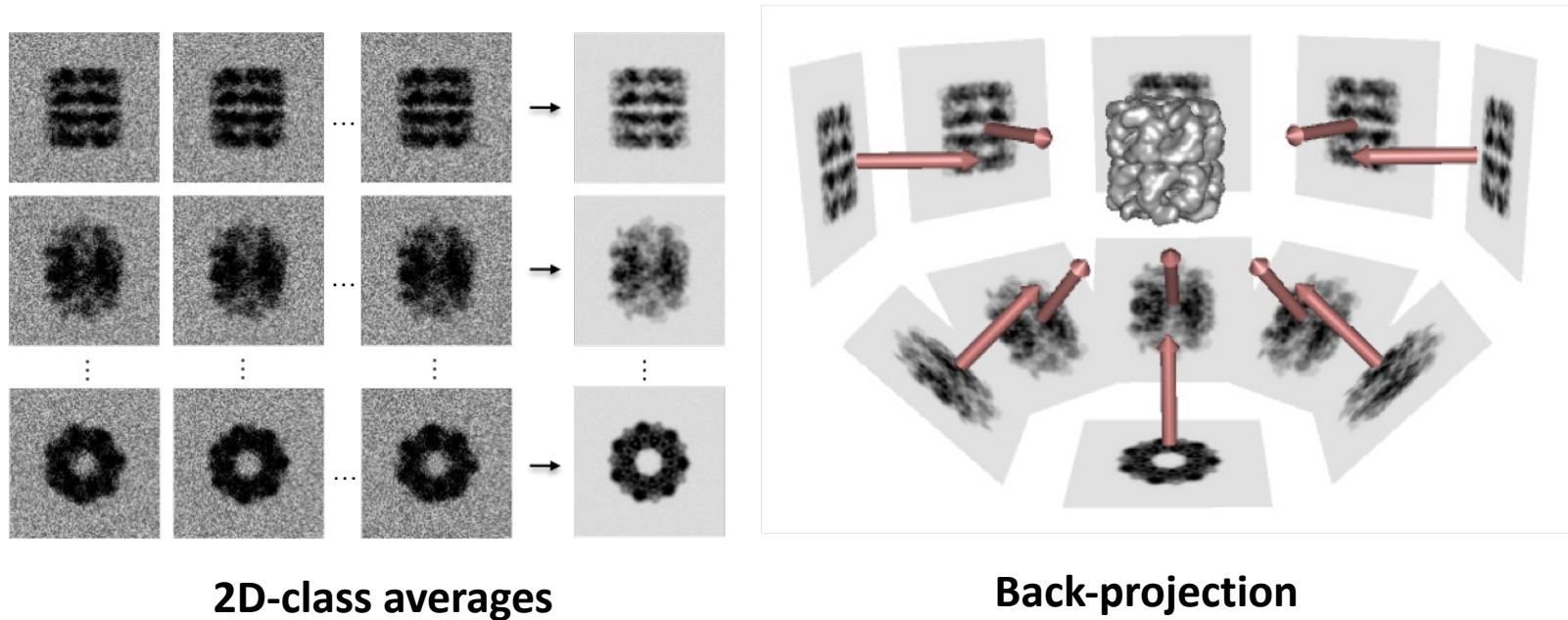
3D structure determination by cryo-EM



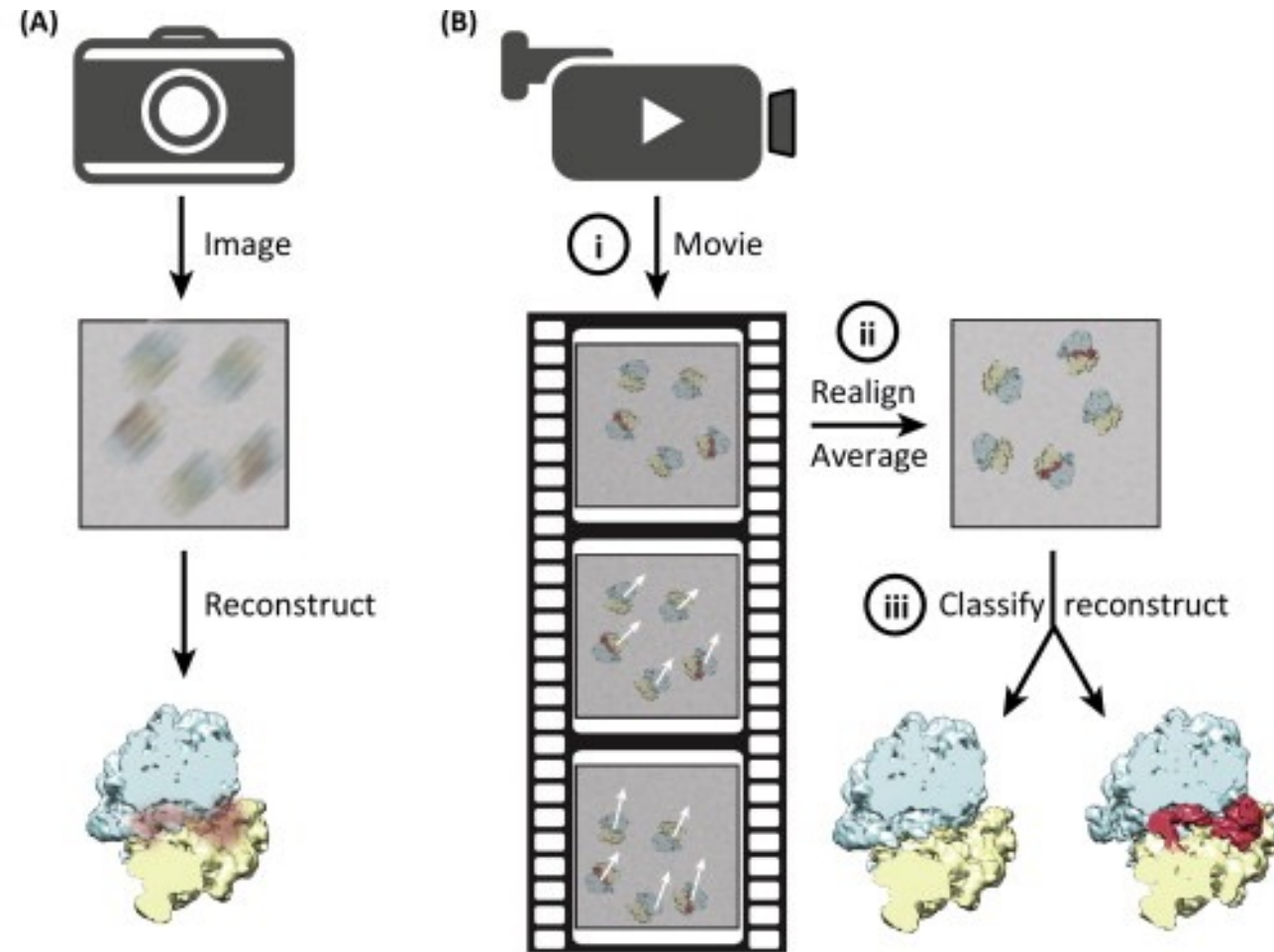
Particle picking



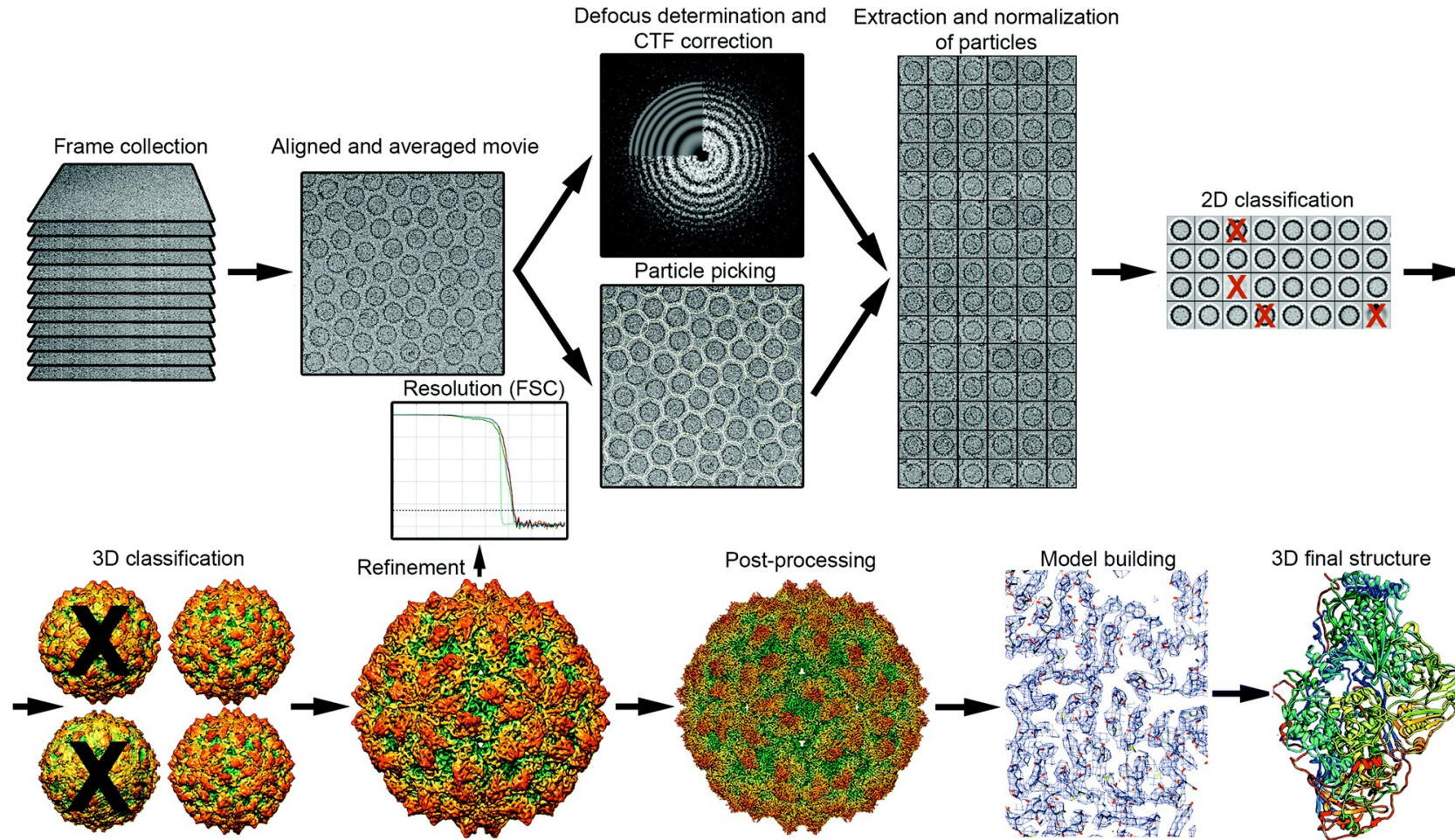
- All the boxed images are then **clustered** - The images in each cluster are **averaged**, to improve the signal-to-noise ratio. The image below shows three clusters, with images from each cluster and the averaged image on the same row.
- The 3D orientation of each of these average images is then found. Using these orientations, the images are **back-projected** to produce a **3D density map**. This 3D density map captures the electron density throughout the macromolecular complex.



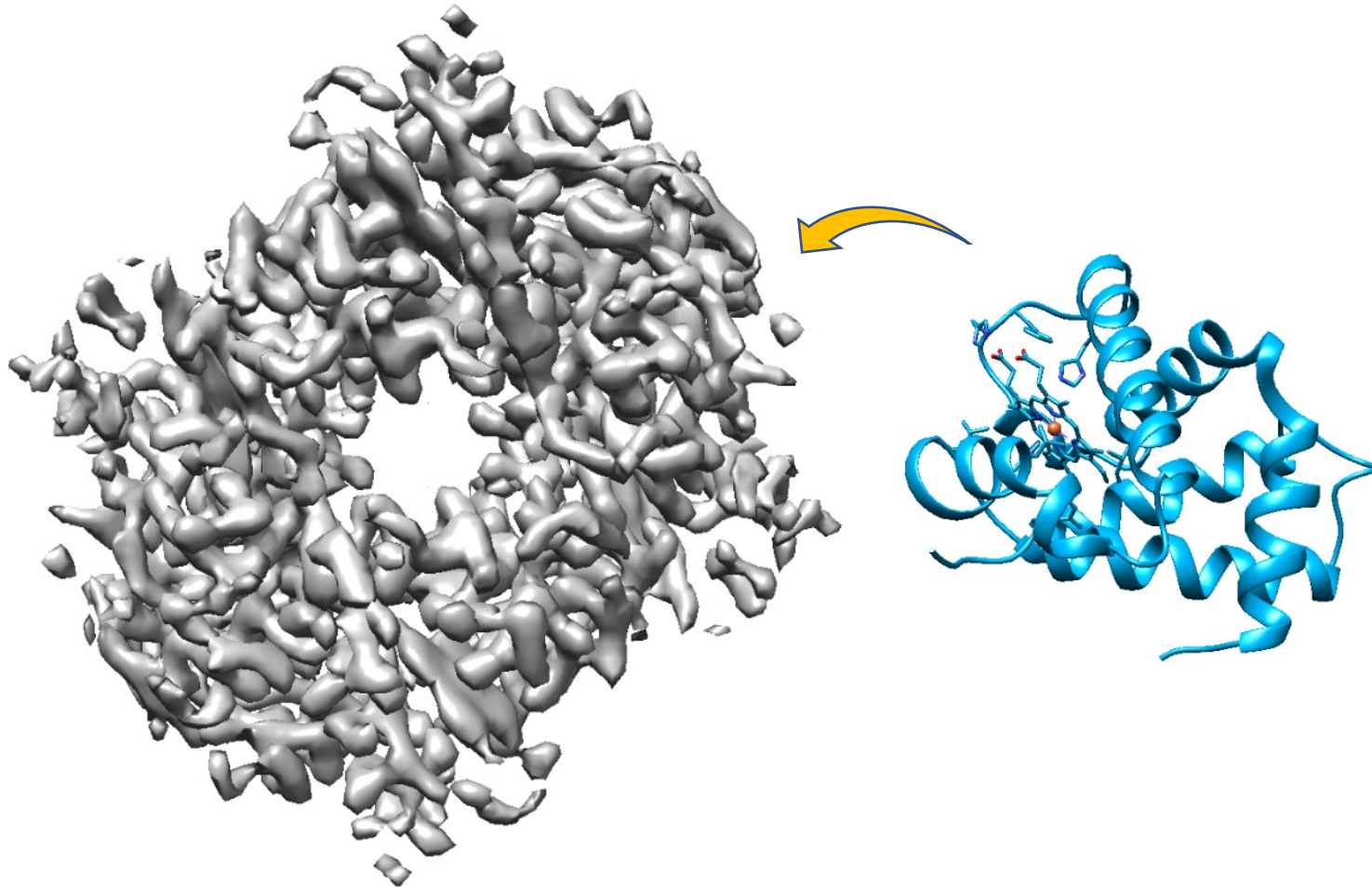
Beam Induced Movement



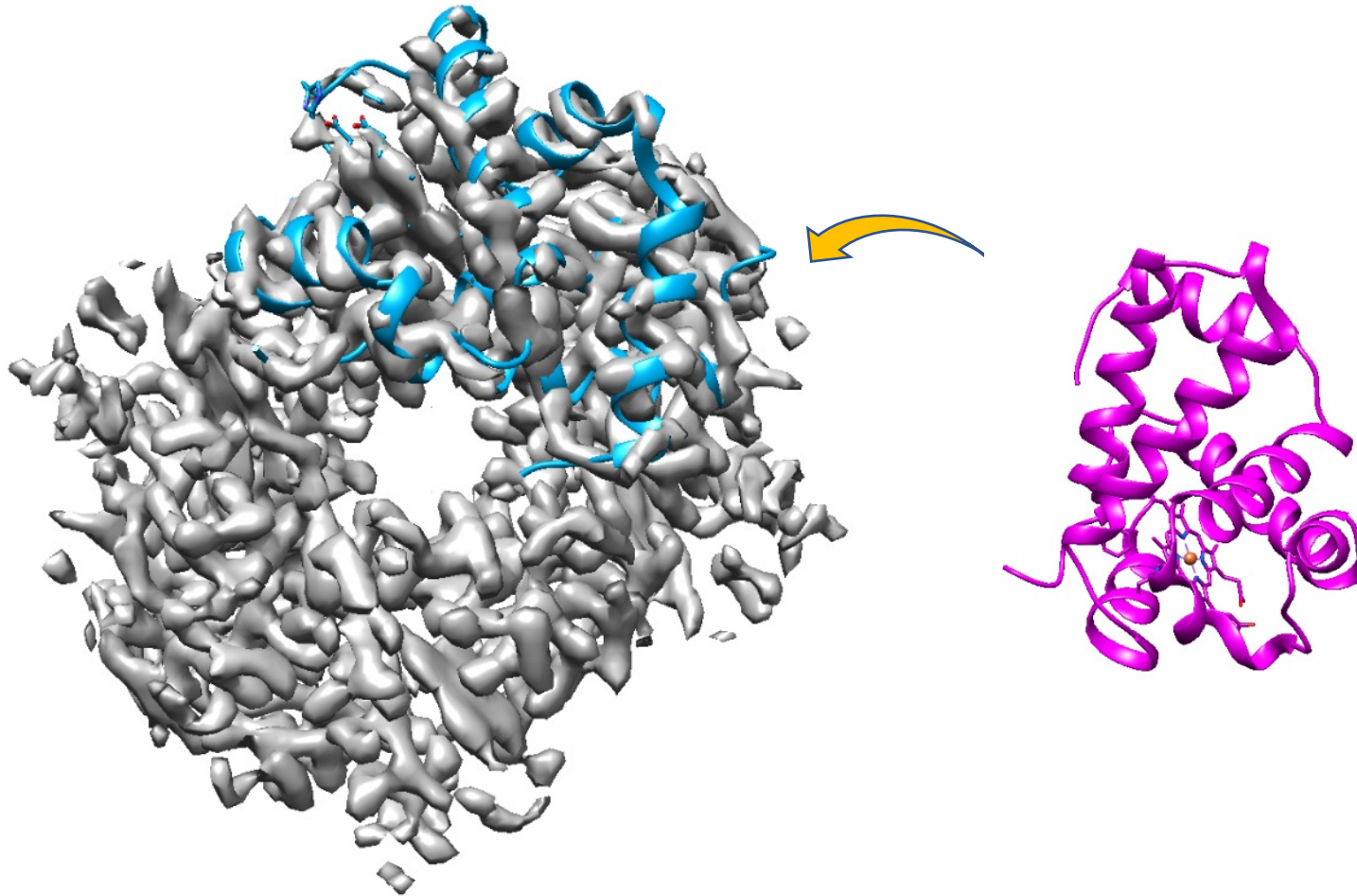
Single particle processing workflow



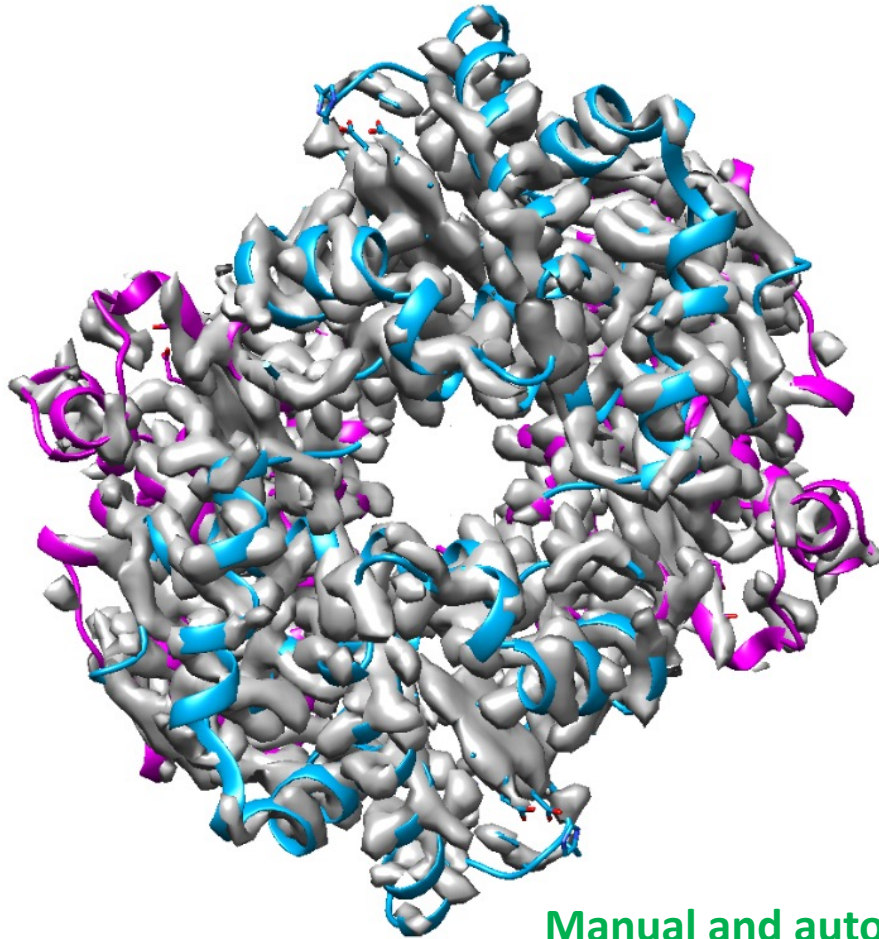
Rigid fitting of the model(s) in the map



Rigid fitting of the model(s) in the map

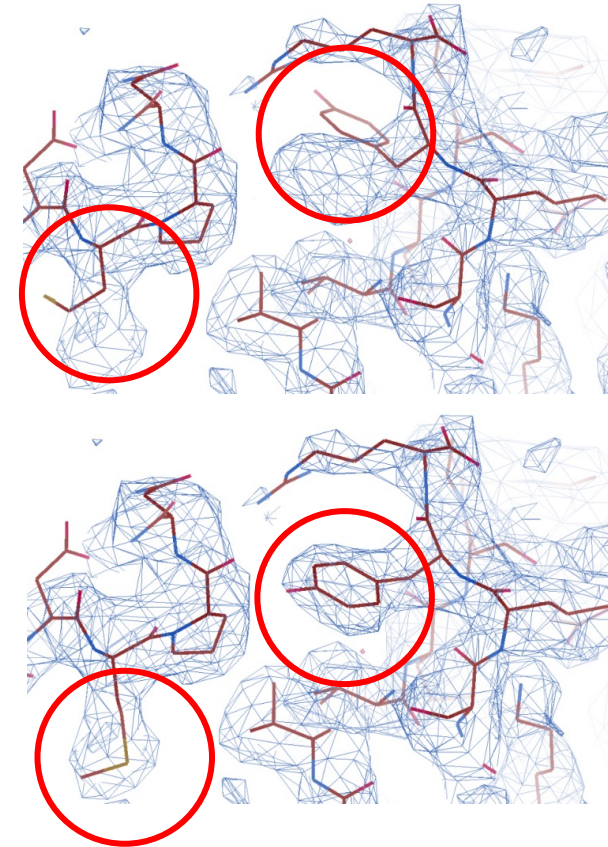


Refinement of models in the map



Manual and automatic

Get the final model(s)



Coffee Break

Structure validation

- Relevance of a (crystal) structure
- Resolution of the diffraction data (X-ray): accuracy of the structure
- R-value and free R-value (cross-validation): difference between observed and calculated scattering data

resolution

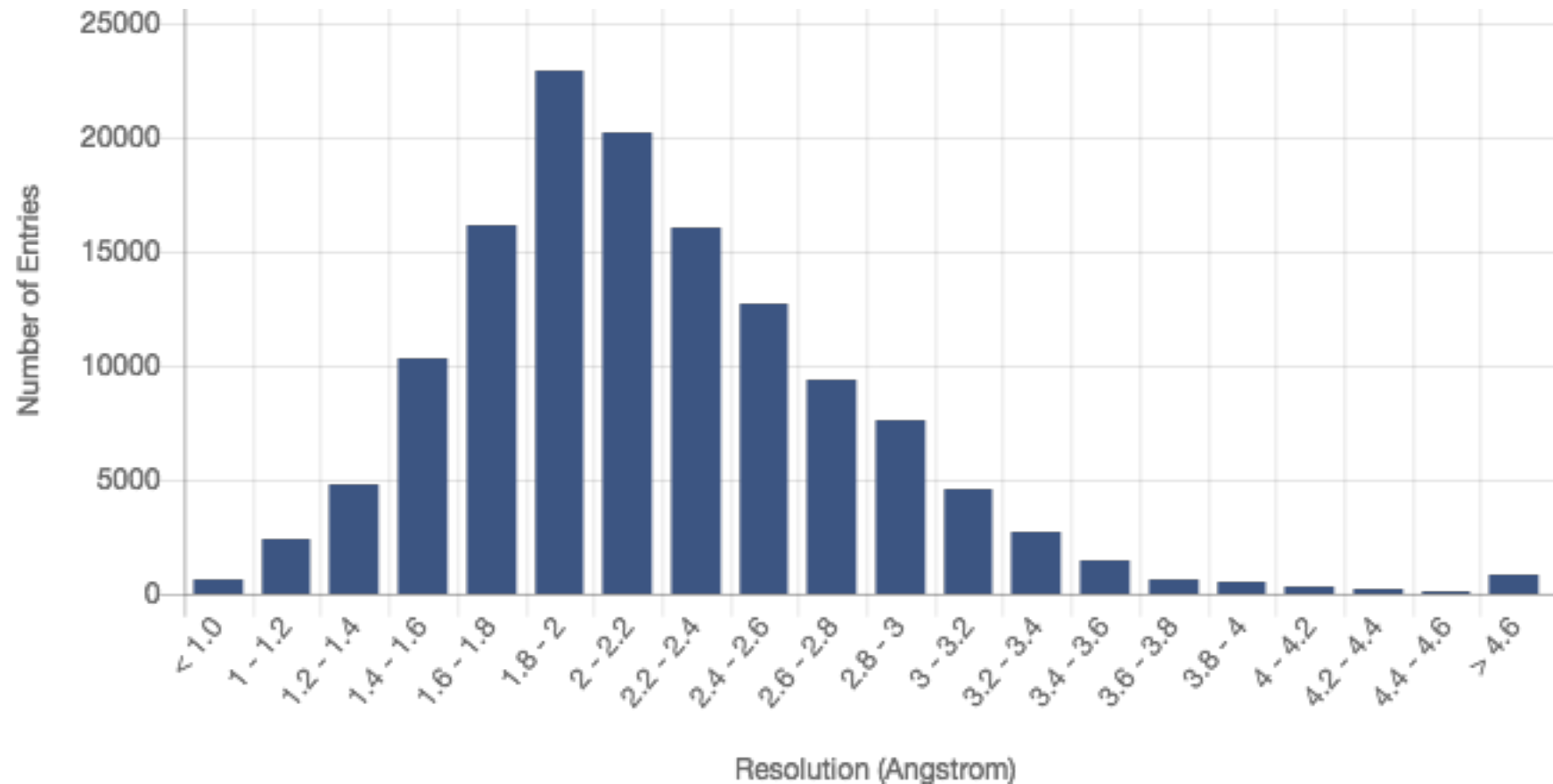
$$n\lambda = 2d \sin \theta$$

wavelength

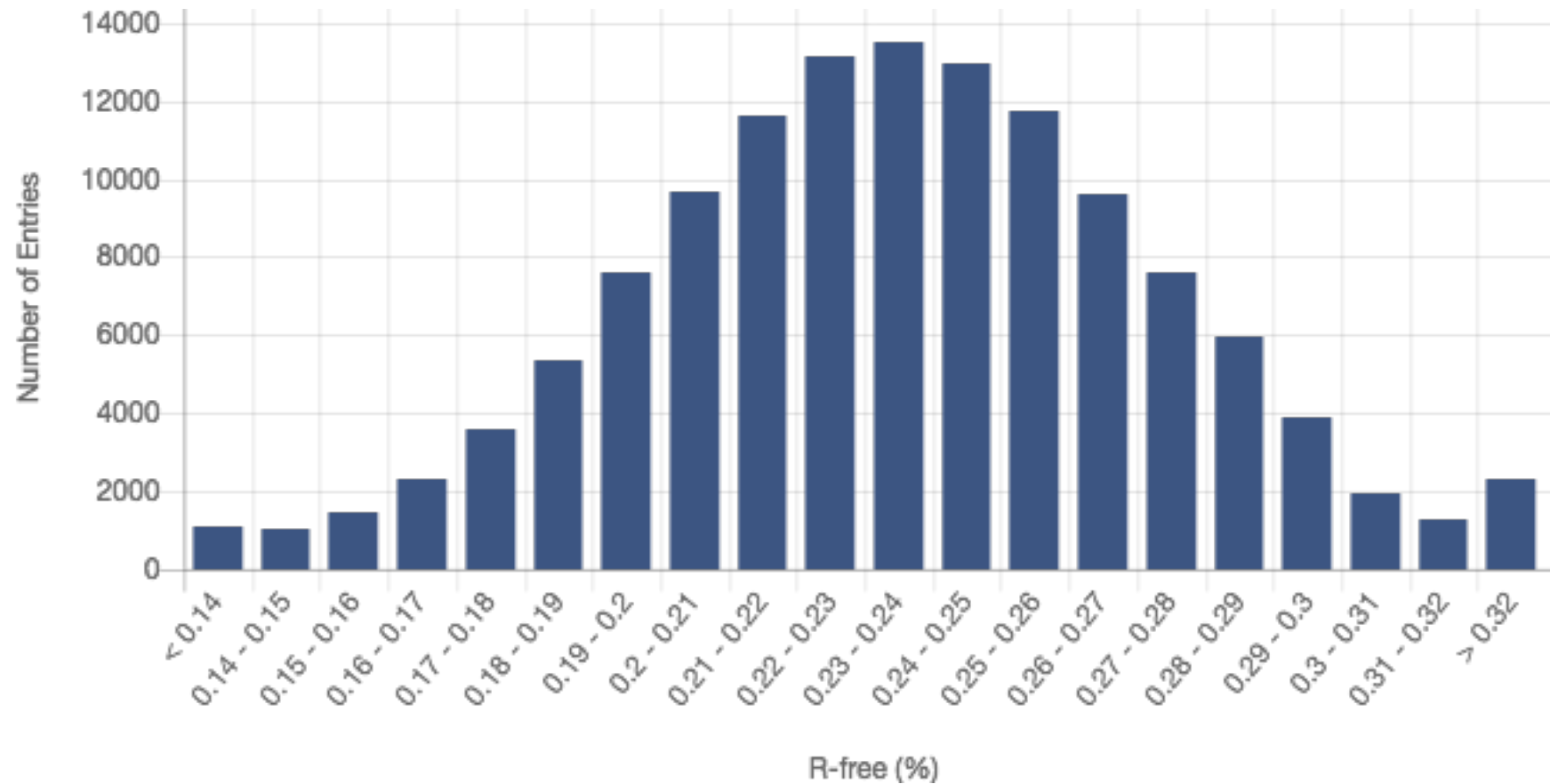
scattering angle



PDB statistics: resolution

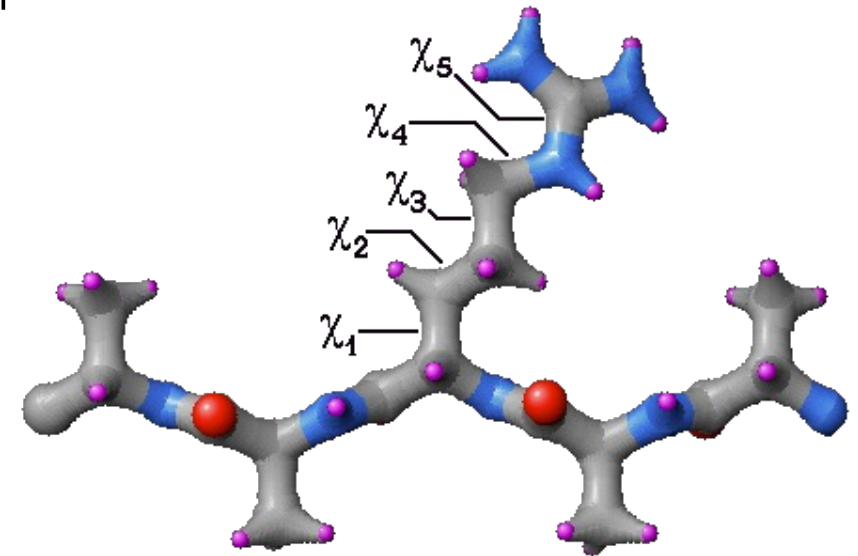
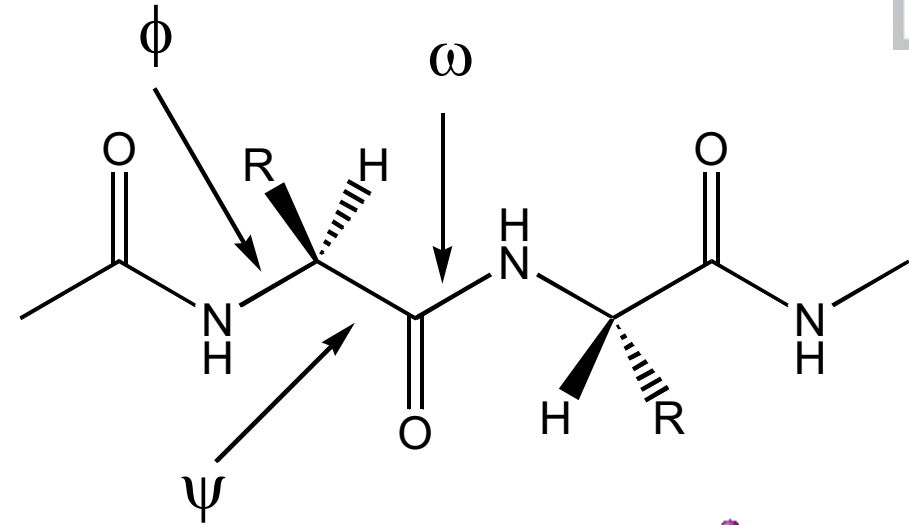


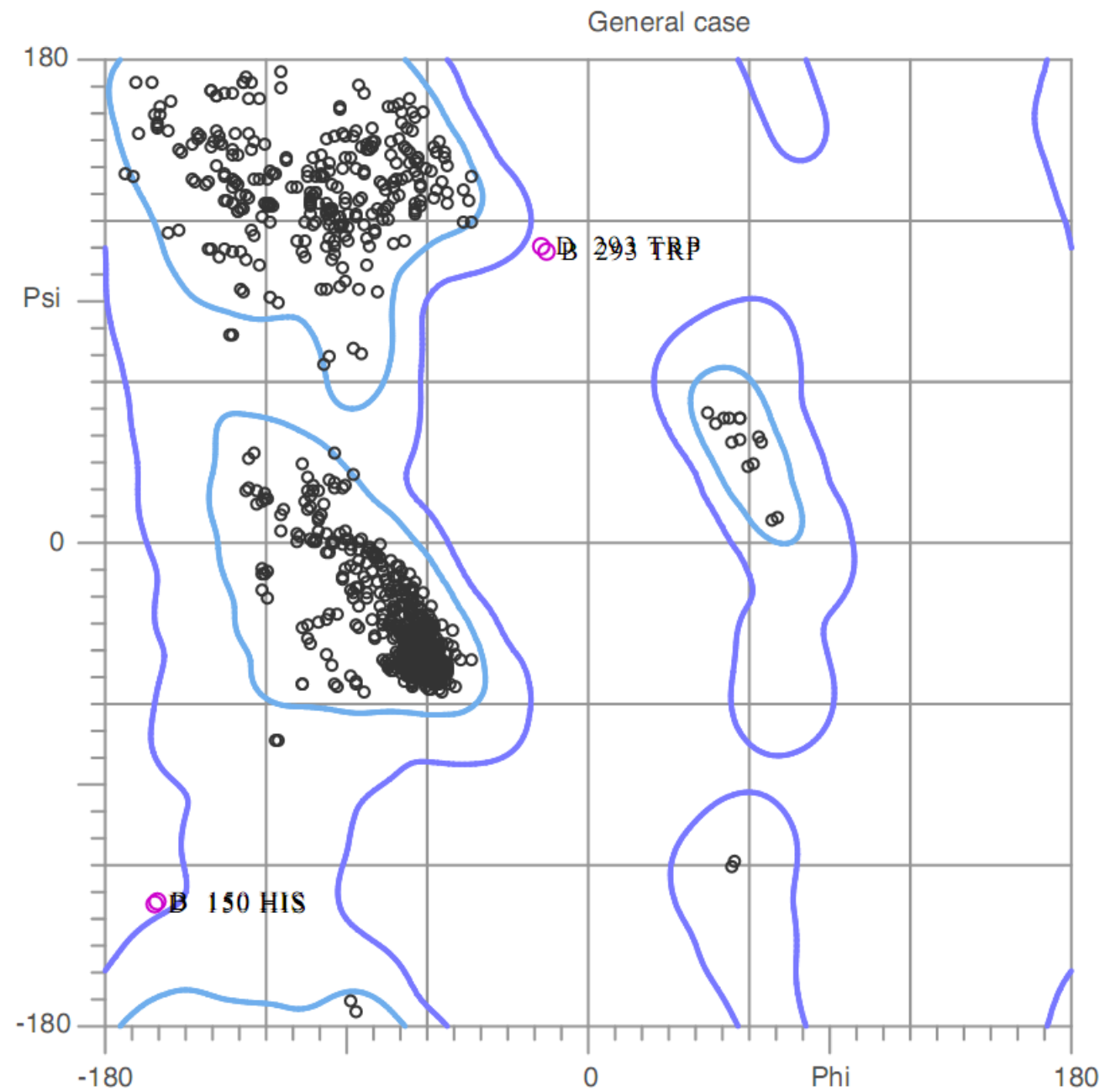
PDB statistics: R_{free} values

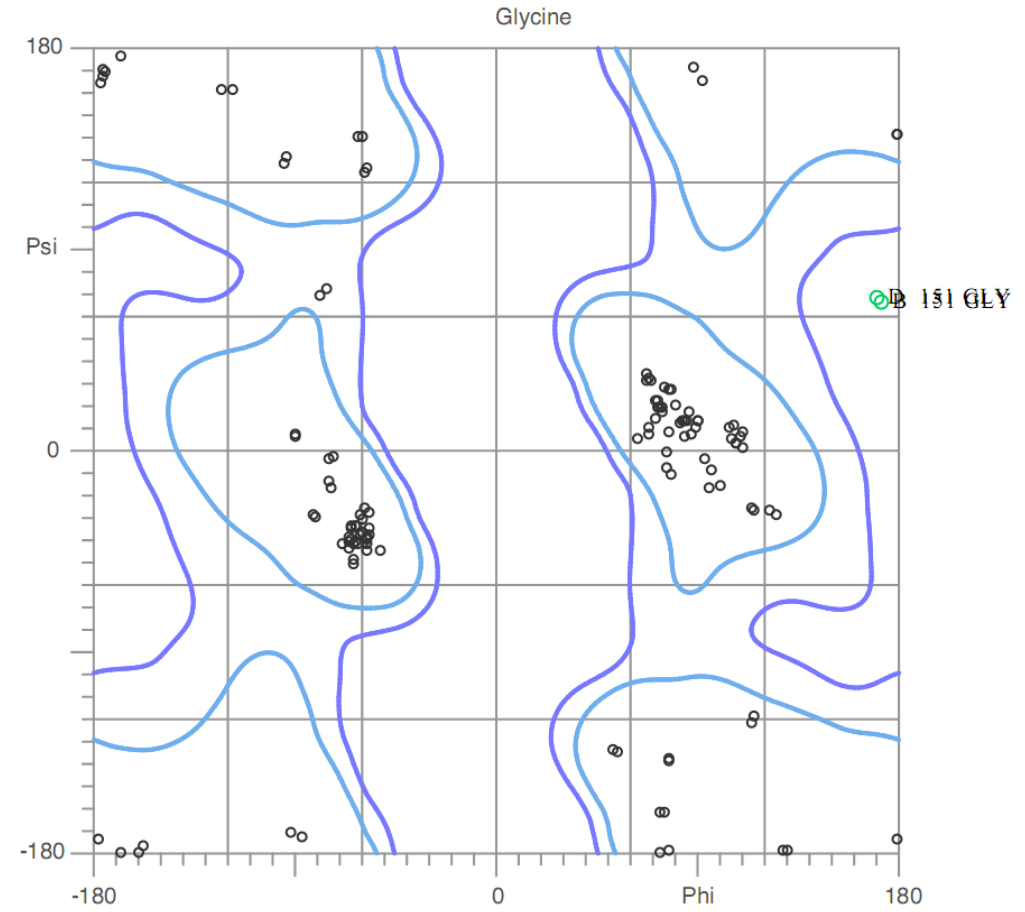
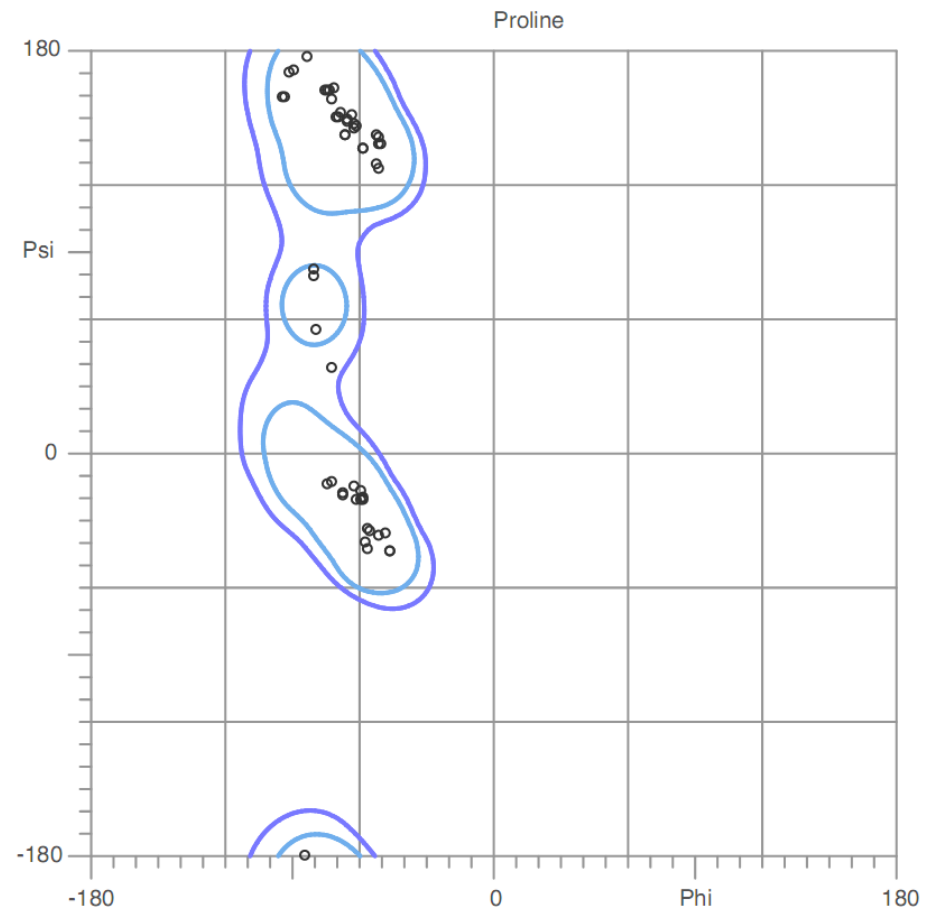


Structure validation

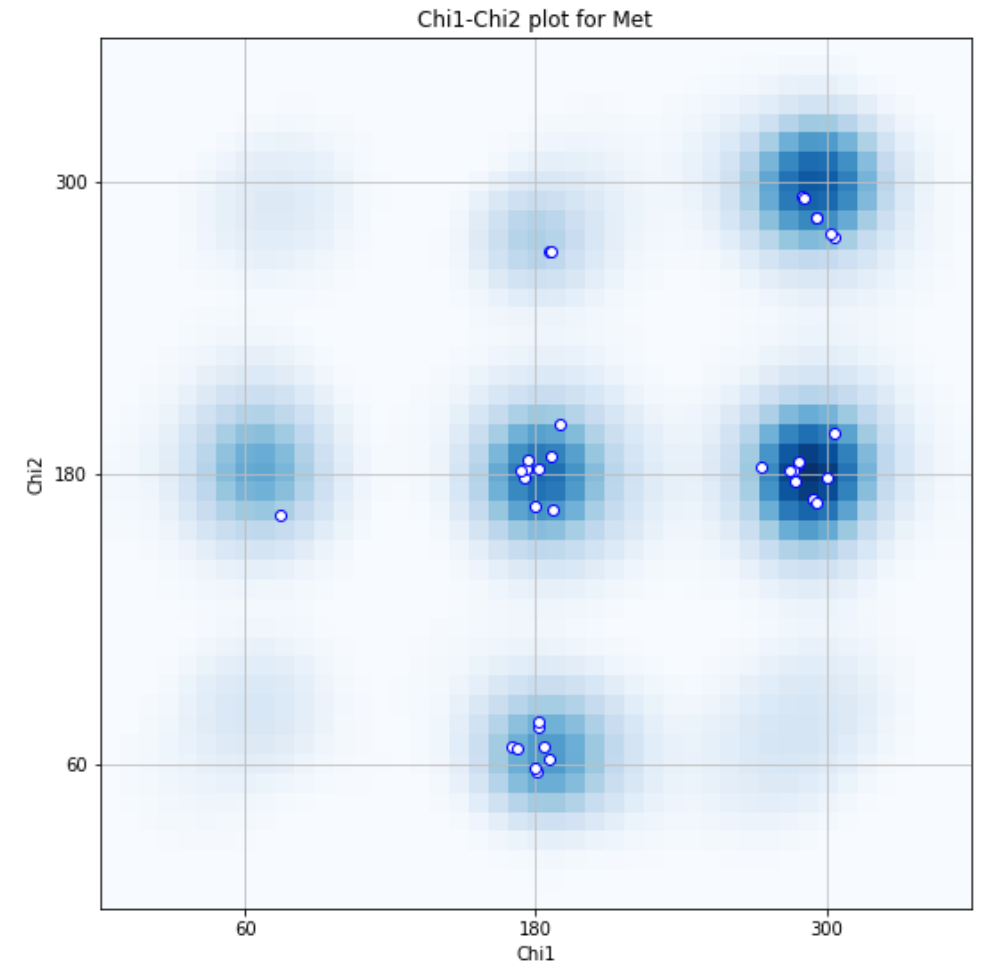
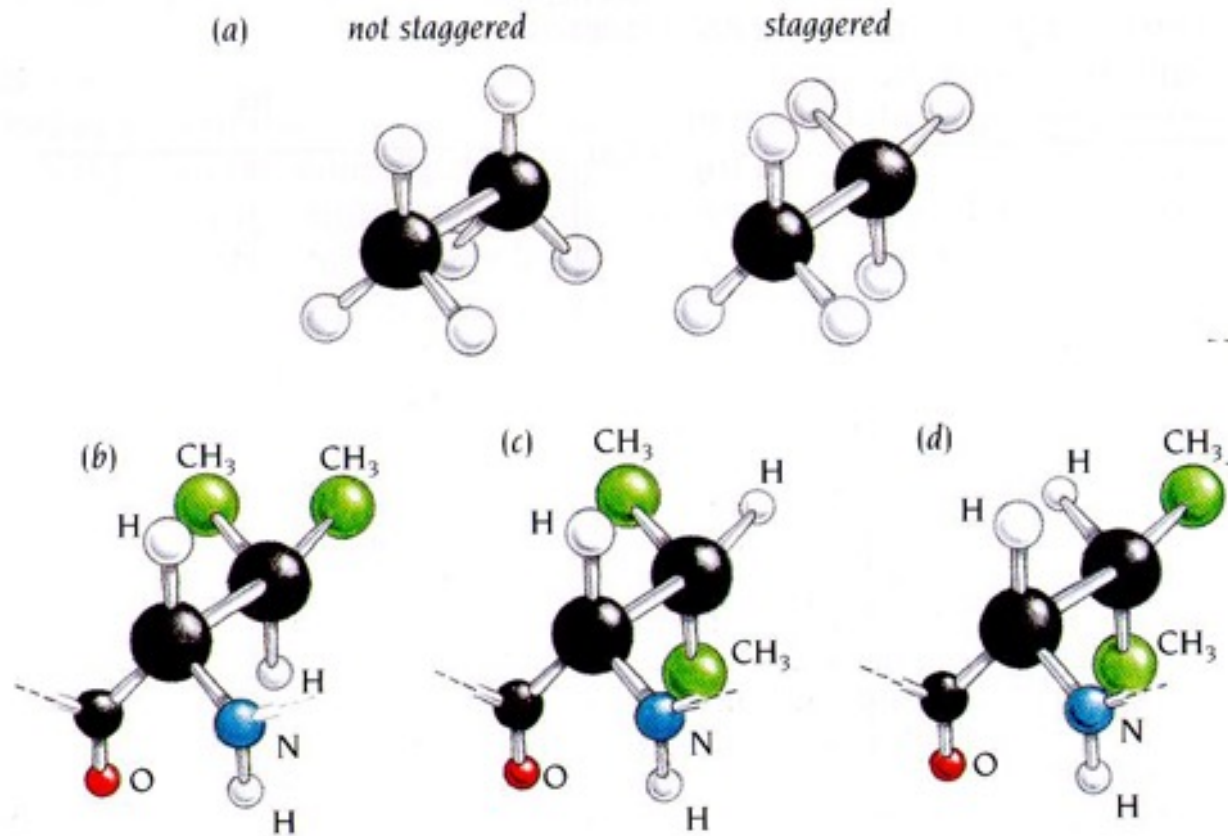
- Geometry:
bond lengths, bond angles
- Ramachandran-Plot:
dihedral angles along the main chain
- Side chain torsion angles





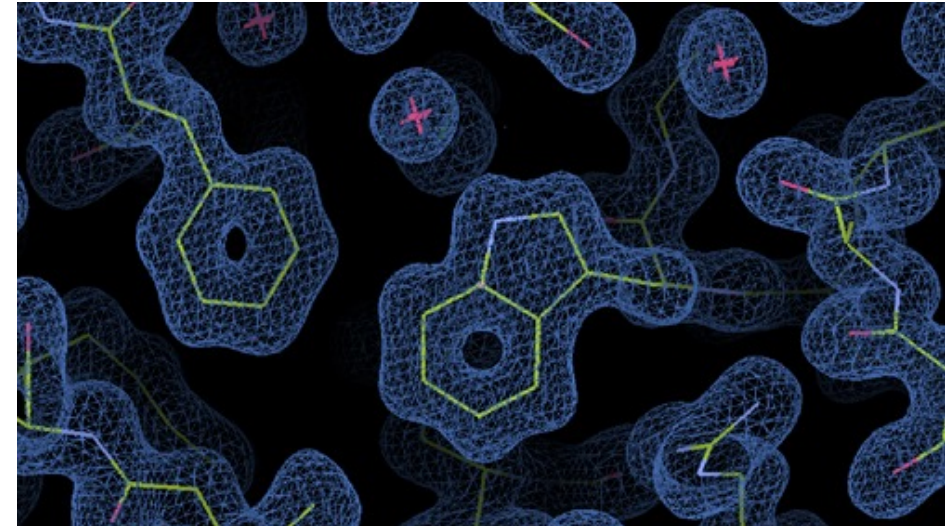


Sidechain conformations



Structure validation

- Biologically active form:
quaternary structure in the crystal
EBI-Pisa server
- Electron density:
primary result of an x-ray
diffraction experiment

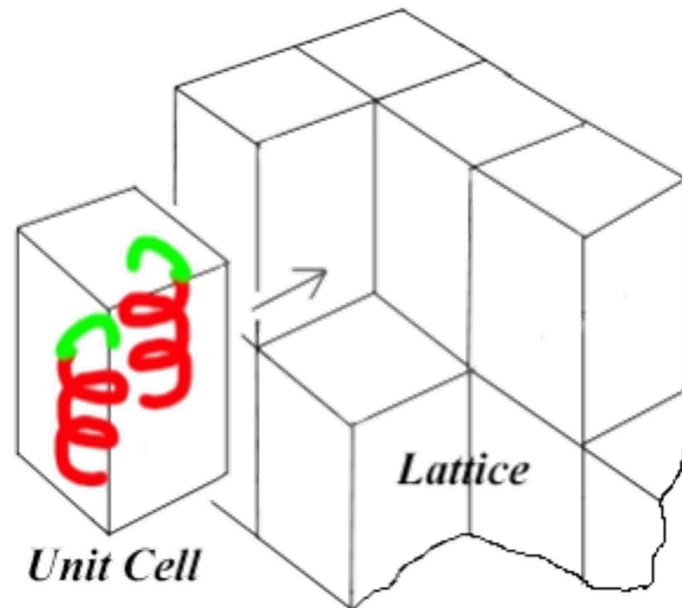


Quarternary structure in crystals

- The coordinates stored in the PDB **may or may not** correspond to the **biologically active oligomer** of a protein!
- *e.g.*: a PDB entry contains **4 chains**, but, in solution, the protein is a **monomer**.
- or *e.g.*: a PDB entry contains **1 chain**, but the biologically active oligomer is a **tetramer**.
- In the last example, the tetramer is generated by **crystal symmetry**.

Crystal symmetry

- Crystals consist of regularly arranged copies of the "unit cell" (translational symmetry).



Crystal symmetry

- **Asymmetric unit:** the portion of the unit cell that has to be considered in a crystallographic experiment
- The content of the entire unit cell results from the application of **symmetry operations** (rotation axes and screw axes).
- The whole crystal is "generated" by copying and translating the unit cell in all three spatial directions.
- **For x-ray crystal structures, only the content of the asymmetric unit is stored in the PDB.**

EBI-Pisa server

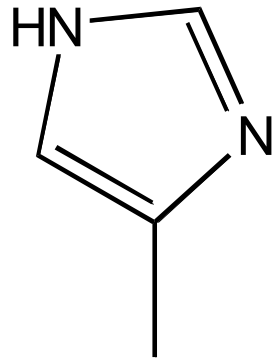
- http://www.ebi.ac.uk/msd-srv/prot_int/pistart.html
- It analyses all **interactions** between molecules in the crystal (within the asymmetric unit and between symmetry-equivalent molecules).
- The interacting molecular surfaces are categorized according to their **size** and **chemical properties**.
- From this data, the program determines those interactions that are likely to be present also in solution.
- **Prediction of the biologically active oligomer (in solution)**

Structure validation – preparation for modeling calculations

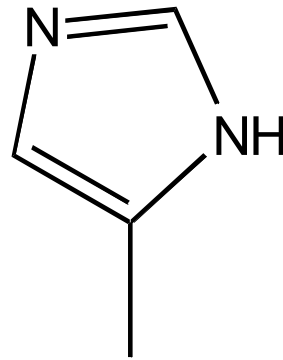
- Missing atoms, missing residues
- Protonation, addition of H-atoms
- Tautomers (His)
- Sidechain conformations (His, Asn, Gln)
- Disulfide bridges

Group Period	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
1	1 H																	2 He
2	3 Li	4 Be											5 B	6 C	7 N	8 O	9 F	10 Ne
3	11 Na	12 Mg											13 Al	14 Si	15 P	16 S	17 Cl	18 Ar
4	19 K	20 Ca	21 Sc	22 Ti	23 V	24 Cr	25 Mn	26 Fe	27 Co	28 Ni	29 Cu	30 Zn	31 Ga	32 Ge	33 As	34 Se	35 Br	36 Kr
5	37 Rb	38 Sr	39 Y	40 Zr	41 Nb	42 Mo	43 Tc	44 Ru	45 Rh	46 Pd	47 Ag	48 Cd	49 In	50 Sn	51 Sb	52 Te	53 I	54 Xe
6	55 Cs	56 Ba	57 La	* 72 Hf	73 Ta	74 W	75 Re	76 Os	77 Ir	78 Pt	79 Au	80 Hg	81 Tl	82 Pb	83 Bi	84 Po	85 At	86 Rn
7	87 Fr	88 Ra	89 Ac	* 104 Rf	105 Db	106 Sg	107 Bh	108 Hs	109 Mt	110 Ds	111 Rg	112 Cn	113 Nh	114 Fl	115 Mc	116 Lv	117 Ts	118 Og
				* 58 Ce	59 Pr	60 Nd	61 Pm	62 Sm	63 Eu	64 Gd	65 Tb	66 Dy	67 Ho	68 Er	69 Tm	70 Yb	71 Lu	
				* 90 Th	91 Pa	92 U	93 Np	94 Pu	95 Am	96 Cm	97 Bk	98 Cf	99 Es	100 Fm	101 Md	102 No	103 Lr	

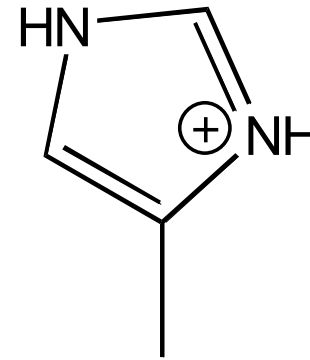
Protonation - Tautomers



His-E
HIE

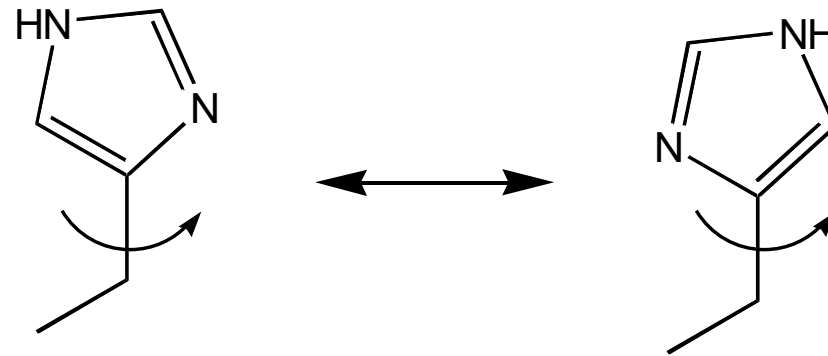


His-D
HID

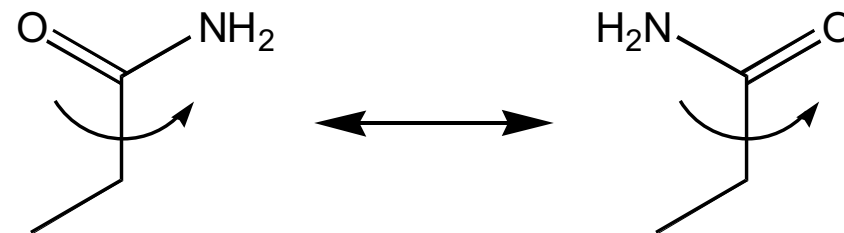


His-P
HIP

Conformations of certain sidechains – problems distinguishing C, N und O



His



Gln oder Asn

Structural databases

8



PDB

>175,000
polypeptides,
nucleotides
& saccharides



CSD

>1.1 million
organic and
metal-organic

ICSD

>230,000
(no C-H and C-C
bonds)

Elements,
minerals,
metals

ICDD

PDF-
4/Organics
>540,000
Includes data
derived from
CSD



FIZ Karlsruhe

Leibniz Institute for Information Infrastructure



CCDC

Cambridge Structural Database (CSD)

The world repository of small molecule crystal structures.

The CSD records bibliographic, chemical and crystallographic information for organic and metal-organic compounds, whose 3D structures have been determined using X-ray or neutron diffraction.

