

Chapter 1

Introduction

1.1 General introduction

1.2 Macromolecular Crystallography

1.2.1 X-ray scattering

X-rays are high energy photons part of the electromagnetic spectrum with a wavelength 0.1-10nm [1]. X-rays can be described as packets of travelling electromagnetic waves, whose electric field vector interacts with the charged electrons of matter [1]. Such interaction, typically termed scattering, results in the diffraction of the incoming wave, which X-ray crystallography relies on.

In its simplest form, scattering of X-ray radiation can be explained in the scenario of exposure to a single free electron. The resulting scattering can be classed as elastic (Thomson scattering) or inelastic (Compton scattering) [1]. The latter — scattering that results in a loss of energy of the emitting photon due to energy transfer onto the electron — does not contribute to discrete scattering, the type of scattering X-ray diffraction relies on. In comparison, Thomson scattering does not result in a loss of energy of the emitting photon. This has significant effects, the incoming photon emits with the same frequency causing the electron to oscillate identically further enhancing the signal. However, it is important to understand that the scattering function of an electron is non-isotropic since the scattered intensity emits strongest in forward and backward directions [1].

If we expand the example to include all electrons in an atom and expose the atom to X-ray radiation, our theory needs to be slightly expanded. Given that one or more

electrons in an atom are not free but orbit around the atom's nucleus in a stable and defined manner, the distribution of these electrons around the nucleus determines the scattering of the incoming X-ray photons. The distribution of scattered photon waves is thus an overall representation of the probability distributions of each electron in the atom and is referred to as electron density $\rho(\vec{r})$. In X-ray scattering, it suffices to approximate the shape of the electron density to a sphere. If we now consider the emitting wave \vec{s}_1 of an X-ray photon scattered by any position \vec{r} in the electron density of an atom, then the phase difference $\Delta\varphi$ to the incoming wave \vec{s}_0 can be described by Eq. (1.7) [1].

$$\Delta\varphi = 2\pi (\vec{s}_1 - \vec{s}_0) \cdot \vec{r} = 2\pi \cdot \vec{S} \cdot \vec{r} \quad (1.1)$$

If more than one electron in an atom's electron density scatter the incoming X-ray wave, then the emitting partial waves can be described by the atomic scattering function f_s (Eq. (1.8)), which describes the interference of all scattered waves [1]. The total scattering power of an atom is proportional to the number of electrons and element-specific with heavier atoms scattering more strongly. Given the approximation of a centrosymmetric electron density, the atomic scattering function is also symmetric.

$$f_s = \int_{\vec{r}}^{V(atoms)} \rho(\vec{r}) \cdot e^{2\pi i \vec{S} \cdot \vec{r}} \cdot d\vec{r} \quad (1.2)$$

With an enhanced understanding of X-ray scattering of electrons orbiting a single atom, it is important to consider X-ray scattering of adjacent atoms, such as it is typically found in molecules. If the electromagnetic wave of a X-ray photon excites all electrons of adjacent atoms, then the resulting partial waves result in constructive or destructive interference. Maximal interference can be obtained when all partial waves are in phase, and maximal destructive interference when out-of-phase. This leads to varying intensities of the emitting X-ray photon at different points in space. To obtain the overall scattering power F_s of all contributing atoms, Eq. (1.8) needs to be modified to include the sum over all atoms j as described in Eq. (1.9).

$$F_s = \sum_{j=1}^{atoms} f_{s,j}^0 \cdot e^{2\pi i \vec{S} \cdot \vec{r}_j} \quad (1.3)$$

If we now translate our hypothetical experiment into a crystal lattice then our un-

derstanding described in Eq. (1.9) needs to be expanded from a 1-dimensional distance vector \vec{r} to the three dimensional lattice translation vectors \vec{a} , \vec{b} and \vec{c} . The Laue equations (Eq. (1.10)) do exactly that and ultimately determine the positions of the diffraction peaks in 3-dimensional space.

$$\vec{S} \cdot \vec{a} = n_1, \quad \vec{S} \cdot \vec{b} = n_2, \quad \vec{S} \cdot \vec{c} = n_3 \quad (1.4)$$

$$n\lambda = 2d_{hkl}\sin\theta \quad (1.5)$$

Such determination is possible through the findings made by Bragg and Bragg [2], who identified the relationship between the scattering vector \vec{S} and the planes in the crystal lattice. Today, this relationship is defined by the Bragg equation (Eq. (1.11)) [2], which allows us to interpret X-ray diffraction as reflections on discrete lattice planes, which relates the diffraction angle θ to the lattice spacing d_{hkl} (Fig. 1.2) [1]. For maximum diffraction n needs to be integer multiples to result in maximum constructive interference of wavelength λ .

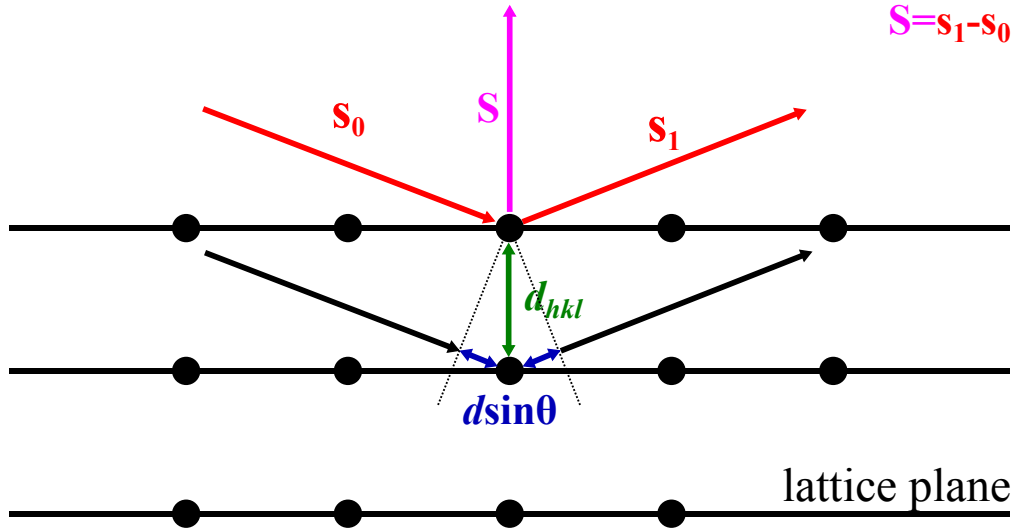


Figure 1.1: Schematic of Bragg scattering.

Lastly, if the hypothetical model is expanded to molecular crystals, then the total scattering from the unit cell is merely a summation of all molecular unit cell scattering contributions in the crystal. Mathematically, this results in Eq. (1.9) being generalised

to Eq. (1.12) through the application of the Laue equations (Eq. (1.10)) to express the scattering vector $\vec{S}\vec{r}_j$ as Miller indices of the reflection planes $\vec{h}\vec{x}_j$.

$$F_h = \sum_{j=1}^{atoms} f_{s,j}^0 \cdot e^{2\pi i \vec{h}\vec{x}_j} \quad (1.6)$$

The structure factor equation defines the scattering power from a crystal in a given reciprocal lattice direction \vec{h} . The scattering is enhanced by the number of repeating units of lattice translation vectors \vec{a} , \vec{b} and \vec{c} , and thus the overall scattering power is proportional to the number of unit cells in the crystal.

It should be noted that Eq. (1.12) is a simplification of the problem at hand. In reality, instrument and experimental corrections need to be applied to the structure factor equation. A correction factor for each experiment-dependent parameter needs to be applied to the structure factor equation. However, in the scope of this work the details of such correction factors do not need to be discussed.

X-rays are high energy photons part of the electromagnetic spectrum with a wavelength 0.1-10nm [1]. X-rays can be described as packets of travelling electromagnetic waves, whose electric field vector interacts with the charged electrons of matter [1]. Such interaction, typically termed scattering, results in the diffraction of the incoming wave, which X-ray crystallography relies on.

In its simplest form, scattering of X-ray radiation can be explained in the scenario of exposure to a single free electron. The resulting scattering can be classed as elastic (Thomson scattering) or inelastic (Crompton scattering) [1]. The latter — scattering that results in a loss of energy of the emitting photon due to energy transfer onto the electron — does not contribute to discrete scattering, the type of scattering X-ray diffraction relies on. In comparison, Thomson scattering does not result in a loss of energy of the emitting photon. This has significant effects, the incoming photon emits with the same frequency causing the electron to oscillate identically further enhancing the signal. However, it is important to understand that the scattering function of an electron is non-isotropic since the scattered intensity emits strongest in forward and backward directions [1].

If we expand the example to include all electrons in an atom and expose the atom to X-ray radiation, our theory needs to be slightly expanded. Given that one or more electrons in an atom are not free but orbit around the atom's nucleus in a stable and

defined manner, the distribution of these electrons around the nucleus determines the scattering of the incoming X-ray photons. The distribution of scattered photon waves is thus an overall representation of the probability distributions of each electron in the atom and is referred to as electron density $\rho(\vec{r})$. In X-ray scattering, it suffices to approximate the shape of the electron density to a sphere. If we now consider the emitting wave \vec{s}_1 of an X-ray photon scattered by any position \vec{r} in the electron density of an atom, then the phase difference $\Delta\varphi$ to the incoming wave \vec{s}_0 can be described by Eq. (1.7) [1].

$$\Delta\varphi = 2\pi (\vec{s}_1 - \vec{s}_0) \cdot \vec{r} = 2\pi \cdot \vec{S} \cdot \vec{r} \quad (1.7)$$

If more than one electron in an atom's electron density scatter the incoming X-ray wave, then the emitting partial waves can be described by the atomic scattering function f_s (Eq. (1.8)), which describes the interference of all scattered waves [1]. The total scattering power of an atom is proportional to the number of electrons and element-specific with heavier atoms scattering more strongly. Given the approximation of a centrosymmetric electron density, the atomic scattering function is also symmetric.

$$f_s = \int_{\vec{r}}^{V(atoms)} \rho(\vec{r}) \cdot e^{2\pi i \vec{S} \cdot \vec{r}} \cdot d\vec{r} \quad (1.8)$$

With an enhanced understanding of X-ray scattering of electrons orbiting a single atom, it is important to consider X-ray scattering of adjacent atoms, such as it is typically found in molecules. If the electromagnetic wave of a X-ray photon excites all electrons of adjacent atoms, then the resulting partial waves result in constructive or destructive interference. Maximal interference can be obtained when all partial waves are in phase, and maximal destructive interference when out-of-phase. This leads to varying intensities of the emitting X-ray photon at different points in space. To obtain the overall scattering power F_s of all contributing atoms, Eq. (1.8) needs to be modified to include the sum over all atoms j as described in Eq. (1.9).

$$F_s = \sum_{j=1}^{atoms} f_{s,j}^0 \cdot e^{2\pi i \vec{S} \cdot \vec{r}_j} \quad (1.9)$$

If we now translate our hypothetical experiment into a crystal lattice then our understanding described in Eq. (1.9) needs to be expanded from a 1-dimensional distance

vector \vec{r} to the three dimensional lattice translation vectors \vec{a} , \vec{b} and \vec{c} . The Laue equations (Eq. (1.10)) do exactly that and ultimately determine the positions of the diffraction peaks in 3-dimensional space.

$$\vec{S} \cdot \vec{a} = n_1, \quad \vec{S} \cdot \vec{b} = n_2, \quad \vec{S} \cdot \vec{c} = n_3 \quad (1.10)$$

$$n\lambda = 2d_{hkl}\sin\theta \quad (1.11)$$

Such determination is possible through the findings made by Bragg and Bragg [2], who identified the relationship between the scattering vector \vec{S} and the planes in the crystal lattice. Today, this relationship is defined by the Bragg equation (Eq. (1.11)) [2], which allows us to interpret X-ray diffraction as reflections on discrete lattice planes, which relates the diffraction angle θ to the lattice spacing d_{hkl} (Fig. 1.2) [1]. For maximum diffraction n needs to be integer multiples to result in maximum constructive interference of wavelength λ .

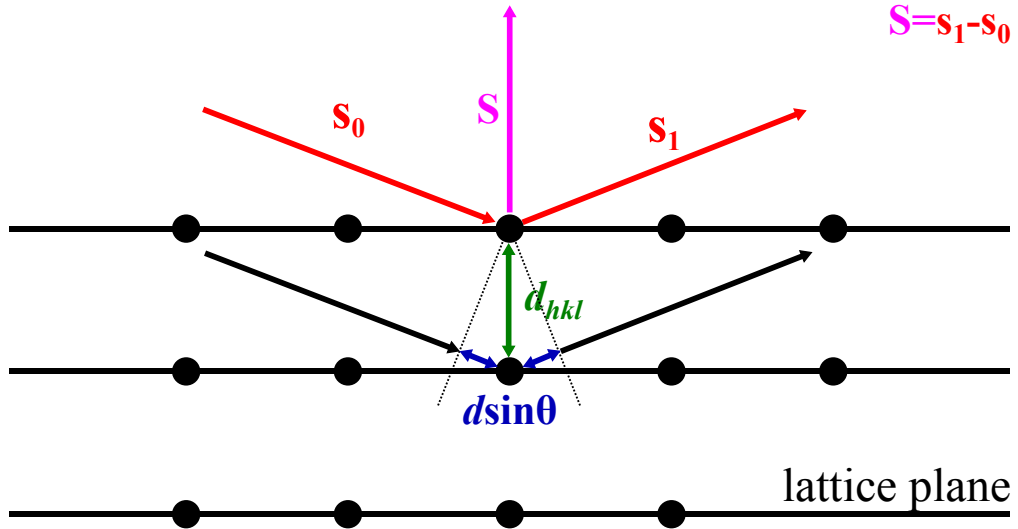


Figure 1.2: Schematic of Bragg scattering.

Lastly, if the hypothetical model is expanded to molecular crystals, then the total scattering from the unit cell is merely a summation of all molecular unit cell scattering contributions (all scattered partial waves emanating from all atoms in all unit cells) in the crystal. Mathematically, this results in Eq. (1.9) being generalised to Eq. (1.12) through

the application of the Laue equations (Eq. (1.10)) to express the scattering vector $\vec{S}\vec{r}_j$ as Miller indices of the reflection planes $\vec{h}\vec{x}_j$.

$$F_h = \sum_{j=1}^{atoms} f_{s,j}^0 \cdot e^{2\pi i \vec{h}\vec{x}_j} \quad (1.12)$$

The structure factor equation defines the scattering power from a crystal in a given reciprocal lattice direction \vec{h} . The scattering is enhanced by the number of repeating units of lattice translation vectors \vec{a} , \vec{b} and \vec{c} , and thus the overall scattering power is proportional to the number of unit cells in the crystal.

It should be noted that Eq. (1.12) is a simplification of the problem at hand. In reality, instrument and experimental corrections need to be applied to the structure factor equation. A correction factor for each experiment-dependent parameter needs to be applied to the structure factor equation. However, in the scope of this work the details of such correction factors do not need to be discussed.

1.2.2 Structure determination

Chapter 2

General methodology

2.1 Dataset creation

2.1.1 FLUME dataset

A test set of 21 globular protein targets was manually selected to include a range of chain lengths, fold architectures, X-ray diffraction data resolutions and Multiple Sequence Alignment (MSA) depths for contact prediction (Table 2.1). The test set covered the three fold classes (α -helical, mixed α - β and β -sheet) and targets were grouped using their DSSP [3] secondary-structure assignment. Target chain lengths fell in the range of [62, 221] residues. Each crystal structure contained one molecule per asymmetric unit and the resolutions of the experimental data was in range from 1.0 to 2.3Å.

2.1.2 KEENO dataset

An unbiased selection of 27 non-redundant protein targets was selected using the following protocol (Table 2.2).

The Pfam v29.0 [24] database was filtered for all protein families with at least one representative structure in the RCSB PDB [25] database. Each representative had to have monomeric protein stoichiometry and its fold classified in the SCOPe v2.05 database [26]. Targets with fold assignments other than "a" (all- α), "b" (all- β), "c" (mixed α + β) or "d" (mixed α / β) were excluded to exclusively focus on regular globular protein folds. Each resulting protein target was screened against the RESTful API of the RCSB PDB (www.rcsb.org) webserver to identify targets meeting the following criteria: experimental technique is X-ray crystallography; chain length is ≥ 100 residues and ≤ 250 residues;

Table 2.1: Summary of the FLUME dataset.

PDB ID	Molecule	Resolution (Å)	Space Group	Chain ID	Chain Length	Molecules per ASU	Matthew's Coeffi- cient	Solvent Content (%)	Fold	Citation
1a6m	Oxy-myoglobin	1.00	P2 ₁	A	151	1	1.90	36.00	all- α	[4]
1aba	T4 glutaredoxin	1.45	P2 ₁ 2 ₁ 2 ₁	A	87	1	2.22	44.62	mixed α/β	[5]
1bdo	Biotinyl domain of acetyl-coenzyme A carboxylase	1.80	P2 ₁ 2 ₁ 2	A	80	1	2.48	49.00	all- β	[6]
1bkr	Calponin Homology (CH) domain from β -spectrin	1.10	P2 ₁	A	109	1	2.04	39.80	all- α	[7]
1chd	CheB methyltransferase domain	1.75	P3 ₂ 21	A	203	1	2.35	47.65	mixed α/β	[8]
1e0s	G-protein Arf6-GDP	2.28	P6 ₁ 22	A	174	1	2.18	37.00	mixed α/β	[9]
1eaz	Phosphoinositol (3,4)-bisphosphate PH domain	1.40	C222 ₁	A	125	1	2.48	48.00	mixed $\alpha+\beta$	[10]
1hh8	N-terminal region of P67Phox	1.80	P3 ₁	A	213	1	2.71	45.00	all- α	[11]
1kjl	Galectin-3 domain	1.40	P2 ₁ 2 ₁ 2 ₁	A	146	1	2.15	42.68	all- β	[12]
1kw4	Polyhomeotic SAM domain	1.75	P6 ₅	A	89	1	2.25	45.27	all- α	[13]
1lo7	4-hydroxybenzoyl CoA thioesterase	1.50	I222	A	141	1	2.06	40.22	mixed $\alpha+\beta$	[14]
1npu	Extracellular domain of murine PD-1	2.00	P2 ₁ 2 ₁ 2 ₁	A	117	1	1.67	25.80	all- β	[15]
1pnc	Poplar plastocyanin	1.60	P2 ₁ 2 ₁ 2 ₁	A	99	1	1.82	32.48	all- β	[16]
1tjx	Synaptotagmin I C2B domain	1.04	P3 ₂ 21	A	159	1	2.40	48.00	mixed $\alpha+\beta$	[17]
1tlv	LicT PRD	1.95	P3 ₂ 21	A	221	1	2.80	50.00	all- α	[18]
2nuz	α -spectrin SH3 domain	1.85	P2 ₁ 2 ₁ 2 ₁	A	62	1	2.57	52.16	all- β	[19]
2qyj	Ankyrin	2.05	P6 ₁	A	166	1	2.28	45.99	all- α	[20]
3w56	C2 domain	1.60	I2	A	131	1	2.05	40.10	all- β	[21]
4cl9	N-terminal bromodomain of Brd4	1.40	P2 ₁ 2 ₁ 2 ₁	A	127	1	2.21	44.37	all- α	[22]
4u3h	FN3con	1.98	P4 ₁ 32	A	100	1	2.47	50.27	all- β	[22]
4w97	KstR2	1.60	C2	A	200	1	2.75	55.25	all- α	[23]

resolution is between 1.3 and 2.3Å; structure factor amplitudes are deposited in the Protein Data Bank [25] database; and there is only a single molecule in the asymmetric unit. The resulting protein structures were cross-validated against the Protein Data Bank of Transmembrane Proteins (PDBTM) [27] to exclude any possible matches. Subsequently, one representative entry was randomly selected for each Pfam family.

The final set of 27 non-redundant targets was determined using further target characterisation and grouping of Pfam families. All targets were grouped using three criteria: domain fold, target chain length and alignment depth. The former consisted of the three fold classes all- α , all- β , and mixed α - β (α + β and α/β) and targets were group using the SCOPe assignment. The target chain lengths were obtained from the deposited information via the RESTful API of the RCSB PDB web server and split into three bins, using 150 and 200 residues as bin edges. Furthermore, the alignment depth was calculated for the sequence alignment of each Pfam family and three bins established with bin edges of 100 and 200 sequences. Thus, all targets were classed in three bins for each of the three features.

The final selection of the 27 targets was performed by randomly selecting one target for each feature combination. To ensure even sampling across the three different fold categories, a target function was employed to identify roughly even target characteristics in each group. The alignment depth and chain length were used as metrics, and had to be within ± 15 units to the values of the other fold classes. This created two conditions that had to be met for a randomly chosen sample to be accepted.

2.1.3 ETHERWOOD dataset

The selection of this dataset was done by [51]. In summary, 14 non-redundant transmembrane protein targets were selected from the PDBTM [27], with a chain length of < 250 residues and resolution of $< 2.5\text{\AA}$. The final selection is summarised in Table 2.3.

2.2 Enhancement of β -sheet restraints

Structure prediction of β -strand containing protein targets *ab initio* is a notoriously challenging task. β -strands, potentially far in sequence space, form a β -sheet in 3-dimensions. Since fragment-assembly algorithms work on the basis of randomly inserting one fragment

Table 2.2: Summary of the KEENO dataset.

PDB ID	Molecule	Resolution (Å)	Space Group	Chain ID	Chain Length	Molecules per ASU	Matthew's Coeffi- cient	Solvent Content (%)	Fold	Citation
1fcy	Retinoic acid nuclear receptor HRAR	1.30	P4 ₁ 2 ₁ 2	A	236	1	2.25	45.50	all-α	[28]
1fvq	Peptide methionine sulfoxide reductase	1.60	C121	A	199	1	2.10	41.55	mixed α+β	[29]
1gm4	Cytochrome C3	2.05	P6 ₁ 22	A	107	1	2.48	50.43	all-α	[30]
1gv8	N-II domain of ovotransferrin	1.95	P3 ₁	A	159	1	2.24	45.00	mixed α/β	[31]
1k40	FAT domain of focal adhesion kinase	2.25	C121	A	126	1	2.21	44.40	all-α	[32]
1oee	Hypothetical protein YodA	2.10	C121	A	193	1	2.30	46.20	all-β	[33]
1oz9	Hypothetical protein AQ-1354	1.89	P4 ₃ 2 ₁ 2	A	150	1	2.76	55.07	mixed α+β	[34]
1q8c	Hypothetical protein MG027	2.00	P4 ₁	A	151	1	2.42	49.25	all-α	[35]
1rlh	Conserved hypothetical protein	1.80	P6 ₃	A	173	1	2.12	41.98	mixed α+β	[36]
1s2x	Cag-Z	1.90	P2 ₁ 2 ₁ 2 ₁	A	206	1	2.74	54.70	all-α	[36]
1u61	Putative Ribonuclease III	2.15	I4 ₁ 32	A	138	1	6.50	80.80	all-α	[36]
1zxu	At5g01750 protein	1.70	P2 ₁ 2 ₁ 2 ₁	A	217	1	2.50	50.20	mixed α+β	[37]
2eum	Glycolipid transfer protein	2.30	C121	A	209	1	2.25	45.39	all-α	[38]
2ol8	Outer surface protein A	1.90	P12 ₁ 1	O	249	1	2.19	43.87	all-β	[39]
2oqz	Sortase B	1.60	P12 ₁ 1	A	223	1	2.07	40.71	all-β	[40]
2x6u	T-Box transcription factor TBX5	1.90	P2 ₁ 2 ₁ 2 ₁	A	203	1	2.20	44.21	all-β	[41]
2y64	Xylanase	1.40	P2 ₁ 2 ₁ 2 ₁	A	167	1	2.15	43.00	all-β	[42]
2yjm	TtrD	1.84	C121	A	176	1	2.08	40.80	all-α	[43]
2yq9	2, 3-cyclic-nucleotide phosphodiesterase	1.90	P2 ₁ 2 ₁ 2 ₁	A	221	1	2.10	41.70	mixed α+β	[44]
3dju	Protein BTG2	2.26	P2 ₁ 2 ₁ 2 ₁	B	122	1	1.98	37.73	mixed α+β	[45]
3g0m	Cysteine desulfuration protein sufe	1.76	P12 ₁ 1	A	141	1	1.88	34.58	mixed α+β	[46]
3qzl	Iron-regulated surface determinant protein A	1.30	P2 ₁ 2 ₁ 2	A	127	1	2.42	49.12	all-β	[47]
4aaJ	N-(5-phosphoribosyl)anthranilate isomerase	1.75	P6 ₁	A	228	1	2.38	48.30	mixed α/β	[48]
4dbb	Amyloid-β A4 precursor protein-binding family A1	1.90	P4 ₁ 2 ₁ 2	A	162	1	3.25	62.10	all-β	[49]
4e9e	Methyl-CpG-binding domain protein 4	1.90	H3	A	161	1	2.42	49.23	all-α	[50]
4lbj	Galectin-3	1.80	P2 ₁ 2 ₁ 2 ₁	A	138	1	2.09	41.01	all-β	[50]
4pgo	Hypothetical protein PF0907	2.30	P6 ₅ 22	A	116	1	3.25	62.10	all-β	[50]

Table 2.3: Summary of the ETHERWOOD dataset.

PDB ID	Molecule	Resolution (Å)	ResolutionSpace Group	Chain ID	Chain Length	Molecules per ASU	Matthew's Coeffi- cient	Solvent Content (%)	Fold	Citation
1gu8	Sensory rhodopsin II	2.27	C222 ₁	A	239	1	2.75	53.00	all-α	[52]
2bhw	Chlorophyll A-B binding protein	2.50	C121	A	232	3	4.10	69.00	all-α	[53]
2evu	Aquaporin aqpM	2.30	I4	A	246	1	3.38	63.57	all-α	[54]
2o9g	Aquaporin Z	1.90	I4	A	234	1	3.34	63.19	all-α	[55]
2wie	ATP synthase C chain	2.13	P6 ₃ 22	A	82	5	3.41	68.00	all-α	[56]
2xov	Rhomboid protease GLPG	1.65	H32	A	181	1	3.50	64.92	all-α	[57]
3gd8	Aquaporin 4	1.80	P42 ₁ 2	A	223	1	2.73	54.97	all-α	[58]
3hap	Bacteriorhodopsin	1.60	C222 ₁	A	249	1	2.73	54.99	all-α	[59]
3ldc	Calcium-gated potassium channel	1.45	P42 ₁ 2	A	82	1	2.48	50.44	all-α	[60]
3ouf	PTHK	1.55	I2	A	97	2	2.40	48.76	all-α	[61]
3pcv	Potassium channel protein	1.90	F23	A	156	1	4.91	74.77	all-α	[62]
3rlb	Leukotriene C4 synthase	2.00	C121	A	192	2	3.89	68.39	all-α	[63]
3u2f	ThiT	2.00	P4 ₂ 22	K	76	5	2.32	46.92	all-α	[64]
4dve	ATP synthase subunit C	2.09	C121	A	198	3	3.27	62.40	all-α	[65]
	Biotin transporter BioY									

at the time, the probability of β -strand formation is much lower compared to α -helices.

Recent advances in *ab initio* structure prediction have seen great improvements in structure prediction quality through the use of predicted residue-residue contacts as distance restraints (see [Introduction](#)). However, only a single approach specifically focused on improvements to the structure prediction of β -sheet formation [66]. To enhance the probability of β -sheet formation in *ab initio* structure prediction, part of this thesis focused on a more general model to enrich restraints between β -strands to attempt better super-secondary quality in the final decoys.

A more general approach, compared to [66] focusing on β -barrel proteins, was developed combining a starting set of contact pairs with a specifically-prepared set obtained from BBCONTACTS [67]. A HHBLITS [68] MSA was constructed using two sequence-search iterations with an E-value cutoff of 10^{-3} against the UniProt20 database [69]. Redundant sequences were removed from the MSA to 90% sequence identity using HHFILTER [68]. Subsequently, the MSA was subjected to CCMPRED [70] for co-evolution based contact prediction. The BBCONTACTS algorithm also requires a secondary-structure prediction, which was obtained using the ADDSS.PL script [68] distributed with the HHSUITE [71]. Both input files were subjected to BBCONTACTS to obtain a final set of β -strand specific contact pairs.

The BBCONTACTS contact pairs were added to a base set of contact pairs usually obtained from a separate (meta-)predictor. The combination of the two sets of contact pairs was done by simple union of the lists; however, if a contact pair was in the intersection, a contact-pair related weight was doubled to allow subsequent modifications of the energy term in distance restraint creation. Furthermore, additional contact pairs were inferred if not present in the base set of contact pairs. The inference worked on the basis that any neighbouring contacts (i.e. $i, j \pm 1$; $i, j \pm 2$; $i \pm 1, j$; $i \pm 2, j$) to contact i, j must be present, and thus any missing were automatically added to the final list. Again, any already present contact pair was assigned double the weight compared to the rest.

2.3 Evaluation of data

This section defines and describes concepts used throughout this thesis to assess and/or validate various data.

2.3.1 Sequence alignment data

2.3.1.1 Sequence alignment depth

Co-evolution based residue-residue contact prediction is dependent on an input MSA ideally containing all homologous sequences found in the queried database. However, the MSA needs a certain level of sequence diversity amongst the homologs to accurately capture the co-evolution signal. The alignment depth — often also referred to as Number of Effective Sequences (M_{eff}) — captures this diversity by computing the number of non-redundant sequences in the MSA.

$$M_{eff} = \sum_i \frac{1}{\sum_j S_{i,j}} \quad (2.1)$$

Various approaches exist for computing M_{eff} [72–74] yielding similar results [75]. In this thesis, the approach defined by Morcos et al. [72] is used. Morcos et al. [72] first described the approach by which sequence weights are computed by means of Hamming distances between all possible sequence combinations in the MSA (Eq. (2.1)). If a Hamming distance was < 0.2 (sequence identity of 80%), the binary value $S_{i,j}$ was assigned 1 and otherwise a 0. The sum of fractional weights of the similarity of each sequence compared to all others ultimately describes the alignment depth.

2.3.2 Contact prediction data

2.3.2.1 Contact map coverage

The fraction of residues covered by a set of contact pairs (N_{map}) out of the total number of residues in the target sequence ($N_{sequence}$) (Eq. (2.2)).

$$Cov = \frac{N_{map}}{N_{sequence}} \quad (2.2)$$

2.3.2.2 Contact map precision

The precision of a set of contact pairs is equivalent to the the proportion of True Positive (TP) contact pairs in the overall set (Eq. (2.3)). A contact pair was defined as TP if the equivalent C β (C α in case of Gly) atoms in the native crystal structure were $< 8\text{\AA}$ apart. The precision value is in range $[0, 1]$, whereby a value of 1 means all contact pairs are

TPs.

$$Prec = \frac{TP}{TP - FP} \quad (2.3)$$

If contacts were unmatched between the target sequence and reference structure, they were not taken into account in the calculation of the precision score.

2.3.2.3 Contact map Jaccard index

The Jaccard index quantifies the similarity between two sets of contact pairs. It describes the proportion of contact pairs in the intersection compared to the union between the two sets [76] (Eq. (2.4)).

$$J_{x,y} = \frac{|x \cap y|}{|x \cup y|} \quad (2.4)$$

The variables x and y are two sets of contact pairs. $|x \cap y|$ is the number of elements in the intersection of x and y , and the $|x \cup y|$ represents the number of elements in the union of x and y . The Jaccard index falls in the range $[0,1]$, with a value of 1 corresponding to identical sets of contact pairs and 0 to non-identical ones. It is worth noting that only exact matches are considered and the neighbourhood of a single contact ignored.

2.3.2.4 Contact map singleton content

Almost all sliced sets of residue-residue contact pairs contain a fraction of contact pairs not co-localising with others. These contact pairs — referred to as singleton contact pairs from here onwards — typically show a high False Positive (FP) rate and could be considered noise (although sometimes they encode TP contacts in an oligomeric interface). To quantify this fraction, a distance-based clustering analysis was defined to identify singleton contact pairs, and thus describe the level of noise in the prediction, or alternatively how well contact pairs co-localise typically between secondary structure features.

To identify singleton contact pairs in a set of contacts, the neighbourhood of each pair was searched for the presence of other contacts. The search radius was defined by ± 2 residues in a 2D-representation of the contact map. If no other contact pair was identified under such constraint, the contact pair was classified as singleton.

2.3.3 Structure prediction data

2.3.3.1 Root Mean Squared Deviation

The Root Mean Square Deviation (RMSD) is a measure to quantify the average atomic distance between two protein structures (Eq. (2.5)). The RMSD is sequence-independent, and measures the distance between Ca atoms.

$$RMSD = \sqrt{\frac{1}{n} \sum_{i,j} (x_i - x_j)^2 + (y_i - y_j)^2 + (z_i - z_j)^2} \quad (2.5)$$

2.3.3.2 Template-Modelling score

The Template-Modelling score (TM-score) is a more accurate measure of structure similarity between two protein structures than the RMSD [77]. Unlike the RMSD, the TM-score score assigns a length-dependent weight to the distances between atoms, with shorter distances getting assigned stronger weights [77]. The TM-score has widely been accepted as a standard for assessing the similarity between two structures, particularly in the field of *ab initio* structure prediction.

$$TMscore = \max \left[\frac{1}{L_{target}} \sum_i^{L_{aligned}} \frac{1}{1 + \left(\frac{d_i}{d_0}\right)^2} \right] \quad (2.6)$$

d_i describes the distance between the i th pair of residues. The distance scale d_0 to normalise the distances is defined by the equation $1.24\sqrt[3]{L_{target} - 15} - 1.8$. The TM-score value falls in the range $(0, 1]$. A TM-score value of < 0.2 indicates two random unrelated structures, and a value > 0.5 roughly the same fold [78]

2.3.3.3 Long-range contact precision

The long-range contact precision score is computed identically to the precision of sets of contact pairs (Section 2.3.2.2). However, the precision score is computed solely for long-range contacts (> 23 residues sequence separation).

2.3.4 Molecular Replacement data

2.3.4.1 Register-Independent Overlap

The Residue-Independent Overlap (RIO) score [79] is a measure of structural similarity between two protein structures considering the total number of atoms within $< 1.5\text{\AA}$. The RIO can be separated into the in- (RIO_{in}) and out-of-register (RIO_{out}) score considering the sequence register between the model and the target. The RIO score is primarily a measure for post-Molecular Replacement (MR) search models to assess the placement of search model atoms with respect to the previously solved crystal structure. To avoid the addition of single atoms place correctly purely by chance, the RIO metric requires at least three consecutive C α atoms to be within the 1.5\AA threshold.

2.3.4.2 Structure solution

MR structure solutions were assessed throughout this thesis always by the SHELXE Correlation Coefficient (CC) and Average Chain Length (ACL) scores. SHELXE performs density modification and main-chain tracing of the refined MR solution [80]. Thorn and Sheldrick [80] highlighted in their work that a CC of $\geq 25\%$ indicates a successful structure solution. Additionally, previous research with Ab initio Modelling of Proteins for moLEcular replacement (AMPLE) [79] has shown that an ACL of the trace needs to be ≥ 10 residues.

In most studies in this thesis, additionally to the SHELXE metrics the post-SHELXE auto-built structures needed R values of ≤ 0.45 . The R values had to be acquired by at least one of the Buccaneer [81] or ARP/wARP [82] solutions.

Lastly, the PHASER Translation Function Z-score (TFZ) and Log-Likelihood Gain (LLG) metrics were also considered when automatically judging a MR solution. Values of > 8 and > 120 were required, respectively. However, the PHASER metrics do not always indicate a structure solution — particularly for smaller fragments — and thus was not considered an essential metric to pass to be considered a successful solution.

Bibliography

- [1] B. Rupp, *Biomolecular crystallography : principles, practice, and application to structural biology*, English, Garland Science, New York, **2010**.
- [2] W. H. Bragg, W. L. Bragg, *Proceedings of the Royal Society A: Mathematical Physical and Engineering Sciences* **July 1913**, *88*, 428–438.
- [3] W. Kabsch, C. Sander, en, *Biopolymers* **Dec. 1983**, *22*, 2577–2637.
- [4] J. Vojtěchovský, K. Chu, J. Berendzen, R. M. Sweet, I. Schlichting, en, *Biophys. J.* **Oct. 1999**, *77*, 2153–2174.
- [5] H. Eklund, M. Ingelman, B. O. Söderberg, T. Uhlin, P. Nordlund, M. Nikkola, U. Sonnerstam, T. Joelson, K. Petratos, en, *J. Mol. Biol.* **Nov. 1992**, *228*, 596–618.
- [6] F. K. Athappilly, W. A. Hendrickson, en, *Structure* **Dec. 1995**, *3*, 1407–1419.
- [7] S. Bañuelos, M. Saraste, K. D. Carugo, en, *Structure* **Nov. 1998**, *6*, 1419–1431.
- [8] A. H. West, E. Martinez-Hackert, A. M. Stock, en, *J. Mol. Biol.* **July 1995**, *250*, 276–290.
- [9] J. Ménétrey, E. Macia, S. Pasqualato, M. Franco, J. Cherfils, en, *Nat. Struct. Biol.* **June 2000**, *7*, 466–469.
- [10] C. C. Thomas, S Dowler, M Deak, D. R. Alessi, D. M. van Aalten, en, *Biochem. J.* **Sept. 2001**, *358*, 287–294.
- [11] S. Grizot, F. Fieschi, M. C. Dagher, E. Pebay-Peyroula, en, *J. Biol. Chem.* **June 2001**, *276*, 21627–21631.
- [12] P. Sörme, P. Arnoux, B. Kahl-Knutsson, H. Leffler, J. M. Rini, U. J. Nilsson, en, *J. Am. Chem. Soc.* **Feb. 2005**, *127*, 1737–1743.
- [13] C. A. Kim, M. Gingery, R. M. Pilpa, J. U. Bowie, en, *Nat. Struct. Biol.* **June 2002**, *9*, 453–457.

- [14] J. B. Thoden, H. M. Holden, Z. Zhuang, D. Dunaway-Mariano, en, *J. Biol. Chem.* **July 2002**, *277*, 27468–27476.
- [15] X. Zhang, J.-C. D. Schwartz, X. Guo, S. Bhatia, E. Cao, M. Lorenz, M. Cammer, L. Chen, Z.-Y. Zhang, M. A. Edidin, S. G. Nathenson, S. C. Almo, en, *Immunity* **Mar. 2004**, *20*, 337–347.
- [16] B. A. Fields, H. H. Bartsch, H. D. Bartunik, F Cordes, J. M. Guss, H. C. Freeman, en, *Acta Crystallogr. D Biol. Crystallogr.* **Sept. 1994**, *50*, 709–730.
- [17] Y. Cheng, S. M. Sequeira, L. Malinina, V. Tereshko, T. H. Söllner, D. J. Patel, en, *Protein Sci.* **Oct. 2004**, *13*, 2665–2672.
- [18] M. Graille, C. Z. Zhou, V. Receveur-Bréchet, B. Collinet, N. Declerck, H. Van Tilbeurgh, en, *J. Biol. Chem.* **Apr. 2005**, *280*, 14780–14789.
- [19] T. Merz, S. K. Wetzel, S. Firbank, A. Plückthun, M. G. Grütter, P. R. E. Mittl, en, *J. Mol. Biol.* **Feb. 2008**, *376*, 232–240.
- [20] D. A. K. Traore, A. J. Brennan, R. H. P. Law, C. Dogovski, M. A. Perugini, N. Lukoyanova, E. W. W. Leung, R. S. Norton, J. A. Lopez, K. A. Browne, H. Yagita, G. J. Lloyd, A. Ciccone, S. Verschoor, J. A. Trapani, J. C. Whisstock, I. Voskoboinik, en, *Biochem. J.* **Dec. 2013**, *456*, 323–335.
- [21] S. J. Atkinson, P. E. Soden, D. C. Angell, M. Bantscheff, C.-W. Chung, K. A. Giblin, N. Smithers, R. C. Furze, L. Gordon, G. Drewes, I. Rioja, J. Witherington, N. J. Parr, R. K. Prinjha, en, *Med. Chem. Commun.* **Feb. 2014**, *5*, 342–351.
- [22] B. T. Porebski, A. A. Nickson, D. E. Hoke, M. R. Hunter, L. Zhu, S. McGowan, G. I. Webb, A. M. Buckle, en, *Protein Eng. Des. Sel.* **Mar. 2015**, *28*, 67–78.
- [23] A. M. Crowe, P. J. Stogios, I. Casabon, E. Evdokimova, A. Savchenko, L. D. Eltis, en, *J. Biol. Chem.* **Jan. 2015**, *290*, 872–882.
- [24] R. D. Finn, P. Coggill, R. Y. Eberhardt, S. R. Eddy, J. Mistry, A. L. Mitchell, S. C. Potter, M. Punta, M. Qureshi, A. Sangrador-Vegas, G. A. Salazar, J. Tate, A. Bateman, en, *Nucleic Acids Res.* **Jan. 2016**, *44*, D279–D285.
- [25] H. M. Berman, J. Westbrook, Z. Feng, G. Gilliland, T. N. Bhat, H. Weissig, I. N. Shindyalov, P. E. Bourne, *Nucleic Acids Res.* **Jan. 2000**, *28*, 235–242.

- [26] J. M. Chandonia, N. K. Fox, S. E. Brenner, en, *J. Mol. Biol.* **Feb. 2017**, *429*, 348–355.
- [27] G. E. Tusnady, Z. Dosztanyi, I. Simon, en, *Nucleic Acids Res.* **Jan. 2005**, *33*, D275–8.
- [28] B. P. Klaholz, A. Mitschler, D. Moras, en, *J. Mol. Biol.* **Sept. 2000**, *302*, 155–170.
- [29] W. T. Lowther, N. Brot, H. Weissbach, B. W. Matthews, en, *Biochemistry* **Nov. 2000**, *39*, 13307–13312.
- [30] R. O. Louro, I. Bento, P. M. Matias, T. Catarino, A. M. Baptista, C. M. Soares, M. A. Carrondo, D. L. Turner, A. V. Xavier, en, *J. Biol. Chem.* **Nov. 2001**, *276*, 44044–44051.
- [31] P. Kuser, D. R. Hall, L. H. Mei, M. Neu, R. W. Evans, P. F. Lindley, en, *Acta Crystallogr. D Biol. Crystallogr.* **May 2002**, *58*, 777–783.
- [32] I. Hayashi, K. Vuori, R. C. Liddington, en, *Nat. Struct. Biol.* **Feb. 2002**, *9*, 101–106.
- [33] G. David, K. Blondeau, M. Schiltz, S. Penel, A. Lewit-Bentley, en, *J. Biol. Chem.* **Oct. 2003**, *278*, 43728–43735.
- [34] V. Oganessian, D. Busso, J. Brandsen, S. Chen, J. Jancarik, R. Kim, S. H. Kim, en, *Acta Crystallographica - Section D Biological Crystallography* **July 2003**, *59*, 1219–1223.
- [35] J. Liu, H. Yokota, R. Kim, S. H. Kim, en, *Proteins: Structure Function and Genetics* **June 2004**, *55*, 1082–1086.
- [36] L. Cendron, A. Seydel, A. Angelini, R. Battistutta, G. Zanotti, en, *J. Mol. Biol.* **July 2004**, *340*, 881–889.
- [37] L. Malinina, M. L. Malakhova, A. T. Kanack, M. Lu, R. Abagyan, R. E. Brown, D. J. Patel, en, *PLoS Biol.* **Nov. 2006**, *4*, 1996–2011.
- [38] K. Makabe, S. Yan, V. Tereshko, G. Gawlak, S. Koide, en, *J. Am. Chem. Soc.* **Nov. 2007**, *129*, 14661–14669.
- [39] A. W. Maresso, R. Wu, J. W. Kern, R. Zhang, D. Janik, D. M. Missiakas, M. E. Duban, A. Joachimiak, O. Schneewind, en, *J. Biol. Chem.* **Aug. 2007**, *282*, 23129–23139.

- [40] C. U. Stirnimann, D. Ptchelkine, C. Grimm, C. W. Müller, en, *J. Mol. Biol.* **July 2010**, *400*, 71–81.
- [41] L. Von Schantz, M. Håkansson, D. T. Logan, B. Walse, J. Osterlin, E. Nordberg-Karlsson, M. Ohlin, M. Hkansson, D. T. Logan, B. Walse, J. Österlin, E. Nordberg-Karlsson, M. Ohlin, en, *Glycobiology* **July 2012**, *22*, 948–961.
- [42] S. J. Coulthurst, A. Dawson, W. N. Hunter, F. Sargent, en, *Biochemistry* **Feb. 2012**, *51*, 1678–1686.
- [43] M. Myllykoski, A. Raasakka, M. Lehtimäki, H. Han, I. Kursula, P. Kursula, en, *J. Mol. Biol.* **Nov. 2013**, *425*, 4307–4322.
- [44] X. Yang, M. Morita, H. Wang, T. Suzuki, W. Yang, Y. Luo, C. Zhao, Y. Yu, M. Bartlam, T. Yamamoto, Z. Rao, en, *Nucleic Acids Res.* **Dec. 2008**, *36*, 6872–6881.
- [45] J. C. Grigg, C. X. Mao, M. E. P. Murphy, en, *J. Mol. Biol.* **Oct. 2011**, *413*, 684–698.
- [46] H. Repo, J. S. Oeemig, J. Djupsjöbacka, H. Iwai, P. Heikinheimo, en, *Acta Crystallogr. D Biol. Crystallogr.* **Nov. 2012**, *68*, 1479–1487.
- [47] M. F. Matos, Y. Xu, I. Dulubova, Z. Otwinowski, J. M. Richardson, D. R. Tomchick, J. Rizo, A. Ho, en, *Proc. Natl. Acad. Sci. U. S. A.* **Mar. 2012**, *109*, 3802–3807.
- [48] S. Moréra, I. Grin, A. Vigouroux, S. Couvé, V. Henriot, M. Saparbaev, A. A. Ishchenko, en, *Nucleic Acids Res.* **Oct. 2012**, *40*, 9917–9926.
- [49] P. M. Collins, K. Bum-Erdene, X. Yu, H. Blanchard, en, *J. Mol. Biol.* **Apr. 2014**, *426*, 1439–1451.
- [50] T. Weinert, V. Olieric, S. Waltersperger, E. Panepucci, L. Chen, H. Zhang, D. Zhou, J. Rose, A. Ebihara, S. Kuramitsu, D. Li, N. Howe, G. Schnapp, A. Pautsch, K. Bargsten, A. E. Protá, P. Surana, J. Kottur, D. T. Nair, F. Basilico, V. Cecatiello, S. Pasqualato, A. Boland, O. Weichenrieder, B. C. Wang, M. O. Steinmetz, M. Caffrey, M. Wang, en, *Nat. Methods* **Feb. 2015**, *12*, 131–133.
- [51] J. M. H. Thomas, F. Simkovic, R. M. Keegan, O. Mayans, Y. Zhang, D. J. Rigden, C. Zhang, Y. Zhang, D. J. Rigden, *Acta Crystallographica Section D Structural Biology* **Dec. 2017**, *73*, 985–996.
- [52] K. Edman, A. Royant, P. Nollert, C. A. Maxwell, E. Pebay-Peyroula, J. Navarro, R. Neutze, E. M. Landau, en, *Structure* **Apr. 2002**, *10*, 473–482.

- [53] J. Standfuss, A. C. T. Van Scheltinga, M. Lamborghini, W. Kühlbrandt, en, *EMBO J.* **Mar. 2005**, *24*, 919–928.
- [54] J. K. Lee, D. Kozono, J. Remis, Y. Kitagawa, P. Agre, R. M. Stroud, en, *Proc. Natl. Acad. Sci. U. S. A.* **Dec. 2005**, *102*, 18932–18937.
- [55] D. F. Savage, R. M. Stroud, en, *J. Mol. Biol.* **May 2007**, *368*, 607–617.
- [56] D. Pogoryelov, Ö. Yildiz, J. D. Faraldo-Gómez, T. Meier, en, *Nat. Struct. Mol. Biol.* **Oct. 2009**, *16*, 1068–1073.
- [57] K. R. Vinothkumar, K. Strisovsky, A. Andreeva, Y. Christova, S. Verhelst, M. Freeman, en, *EMBO J.* **Nov. 2010**, *29*, 3797–3809.
- [58] J. D. Ho, R. Yeh, A. Sandstrom, I. Chorny, W. E. C. Harries, R. A. Robbins, L. J. W. Miercke, R. M. Stroud, en, *Proceedings of the National Academy of Sciences* **May 2009**, *106*, 7437–7442.
- [59] N. H. Joh, A. Oberai, D. Yang, J. P. Whitelegge, J. U. Bowie, en, *J. Am. Chem. Soc.* **Aug. 2009**, *131*, 10846–10847.
- [60] S. Ye, Y. Li, Y. Jiang, en, *Nat. Struct. Mol. Biol.* **Aug. 2010**, *17*, 1019–1023.
- [61] M. G. Derebe, D. B. Sauer, W. Zeng, A. Alam, N. Shi, Y. Jiang, en, *Proceedings of the National Academy of Sciences* **Jan. 2011**, *108*, 598–602.
- [62] H. Saino, Y. Ukita, H. Ago, D. Irikura, A. Nisawa, G. Ueno, M. Yamamoto, Y. Kanaoka, B. K. Lam, K. F. Austen, M. Miyano, en, *J. Biol. Chem.* **May 2011**, *286*, 16392–16401.
- [63] G. B. Erkens, R. P. A. Berntsson, F. Fulyani, M. Majsnerowska, A. Vujičić-Žagar, J. Ter Beek, B. Poolman, D. J. Slotboom, en, *Nat. Struct. Mol. Biol.* **June 2011**, *18*, 755–760.
- [64] J. Symersky, V. Pagadala, D. Osowski, A. Krah, T. Meier, J. D. Faraldo-Gómez, D. M. Mueller, en, *Nat. Struct. Mol. Biol.* **Apr. 2012**, *19*, 485–91, S1.
- [65] R. P.-A. Berntsson, J. ter Beek, M. Majsnerowska, R. H. Duurkens, P. Puri, B. Poolman, D.-J. Slotboom, en, *Proceedings of the National Academy of Sciences* **Aug. 2012**, *109*, 13990–13995.
- [66] S. Hayat, C. Sander, D. S. Marks, A. Elofsson, en, *Proc. Natl. Acad. Sci. U. S. A.* **Apr. 2015**, *112*, 5413–5418.

- [67] J. Andreani, J. Söding, en, *Bioinformatics* **June 2015**, *31*, 1729–1737.
- [68] M. Remmert, A. Biegert, A. Hauser, J. Söding, en, *Nat. Methods* **Dec. 2011**, *9*, 173–175.
- [69] A. Bateman, M. J. Martin, C. O'Donovan, M. Magrane, E. Alpi, R. Antunes, B. Bely, M. Bingley, C. Bonilla, R. Britto, B. Bursteinas, H. Bye-AJee, A. Cowley, A. Da Silva, M. De Giorgi, T. Dogan, F. Fazzini, L. G. Castro, L. Figueira, P. Garmiri, G. Georghiou, D. Gonzalez, E. Hatton-Ellis, W. Li, W. Liu, R. Lopez, J. Luo, Y. Lussi, A. MacDougall, A. Nightingale, B. Palka, K. Pichler, D. Poggioli, S. Pundir, L. Pureza, G. Qi, S. Rosanoff, R. Saidi, T. Sawford, A. Shypitsyna, E. Speretta, E. Turner, N. Tyagi, V. Volynkin, T. Wardell, K. Warner, X. Watkins, R. Zaru, H. Zellner, I. Xenarios, L. Bougueleret, A. Bridge, S. Poux, N. Redaschi, L. Aimò, G. ArgoudPuy, A. Auchincloss, K. Axelsen, P. Bansal, D. Baratin, M. C. Blatter, B. Boeckmann, J. Bolleman, E. Boutet, L. Breuza, C. Casal-Casas, E. De Castro, E. Coudert, B. CuChe, M. Doche, D. Dornevil, S. Duvaud, A. Estreicher, L. Famiglietti, M. Feuermann, E. Gasteiger, S. Gehant, V. Gerritsen, A. Gos, N. Gruaz-Gumowski, U. Hinz, C. Hulo, F. Jungo, G. Keller, V. Lara, P. Lemercier, D. Lieberherr, T. Lombardot, X. Martin, P. Masson, A. Morgat, T. Neto, N. Nospikel, S. Paesano, I. Pedruzzi, S. Pilbout, M. Pozzato, M. Pruess, C. Rivoire, B. Roechert, M. Schneider, C. Sigrist, K. Sonesson, S. Staehli, A. Stutz, S. Sundaram, M. Tognolli, L. Verbregue, A. L. Veuthey, C. H. Wu, C. N. Arighi, L. Arminski, C. Chen, Y. Chen, J. S. Garavelli, H. Huang, K. Laiho, P. McGarvey, D. A. Natale, K. Ross, C. R. Vinayaka, Q. Wang, Y. Wang, L. S. Yeh, J. Zhang, en, *Nucleic Acids Res.* **Jan. 2017**, *45*, D158–D169.
- [70] S. Seemayer, M. Gruber, J. Söding, en, *Bioinformatics* **Nov. 2014**, *30*, 3128–3130.
- [71] J. Söding, en, *Bioinformatics* **Apr. 2005**, *21*, 951–960.
- [72] F. Morcos, A. Pagnani, B. Lunt, A. Bertolino, D. S. Marks, C. Sander, R. Zecchina, J. N. Onuchic, T. Hwa, M. Weigt, en, *Proceedings of the National Academy of Sciences* **Dec. 2011**, *108*, E1293–E1301.
- [73] D. T. Jones, D. W. A. Buchan, D. Cozzetto, M. Pontil, *Bioinformatics* **Jan. 2012**, *28*, 184–190.

- [74] D. T. Jones, T. Singh, T. Kosciulek, S. Tetchner, en, *Bioinformatics* **Apr. 2015**, *31*, 999–1006.
- [75] M. J. Skwark, D. Raimondi, M. Michel, A. Elofsson, en, *PLoS Comput. Biol.* **Nov. 2014**, *10*, e1003889.
- [76] Q. Wuyun, W. Zheng, Z. Peng, J. Yang, *Brief. Bioinform.* **Oct. 2016**, bbw106.
- [77] Y. Zhang, J. Skolnick, en, *Proteins: Structure Function and Genetics* **Dec. 2004**, *57*, 702–710.
- [78] J. Xu, Y. Zhang, en, *Bioinformatics* **Apr. 2010**, *26*, 889–895.
- [79] J. M. H. Thomas, R. M. Keegan, J. Bibby, M. D. Winn, O. Mayans, D. J. Rigden, en, *IUCrJ* **Mar. 2015**, *2*, 198–206.
- [80] A. Thorn, G. M. Sheldrick, en, *Acta Crystallogr. D Biol. Crystallogr.* **Nov. 2013**, *69*, 2251–2256.
- [81] K. Cowtan, en, *Acta Crystallogr. D Biol. Crystallogr.* **Sept. 2006**, *62*, 1002–1011.
- [82] S. X. Cohen, M. Ben Jelloul, F. Long, A. Vagin, P. Knipscheer, J. Lebbink, T. K. Sixma, V. S. Lamzin, G. N. Murshudov, A. Perrakis, en, *Acta Crystallogr. D Biol. Crystallogr.* **Jan. 2007**, *64*, 49–60.