

AlphaFold2-Assisted Domain Determination of the  
SARS-CoV-2 Protein NSP3  
&  
Exploration of the Solvent Region of Macromolecular  
Crystals with Experimental Phasing

Maximilian Edich

Dissertation  
zur Erlangung des Doktorgrades  
an der Fakultät für Mathematik, Informatik und Naturwissenschaften  
Fachbereich Chemie der Universität Hamburg

Angefertigt am Institut für Nanostruktur- und Festkörperphysik

vorgelegt von

**Maximilian Edich**

an der Universität Hamburg eingereichte Dissertation



**Universität Hamburg**  
DER FORSCHUNG | DER LEHRE | DER BILDUNG

Hamburg, Deutschland  
2024

Gutachter/-innen der Dissertation:  
Prof. Dr. Markus Perbandt  
Dr. Andrea Thorn

Prüfungskommission:  
Prof. Dr. Andrew Torda  
Prof. Dr. Gabriel Bester  
Prof. Dr. Alf Mews

Betreuer/-innen:  
Dr. Andrea Thorn  
Prof. Dr. Andrew Torda

Datum der Disputation:  
12.07.2024

Datum der Druckfreigabe:  
02.09.2024



Diese Dissertation wurde angefertigt am Hamburg Advanced Research Centre for  
Bioorganic Chemistry in Hamburg im Zeitraum vom Mai 2021 bis Mai 2024.

## I. List of Publications

### Published during PhD

1. Kandler, L., Kippes, O., Edich, M., Stüb, S., Santoni, G., Thorn, A. (2023). SARS-CoV-2 envelope protein and its relationship to the membrane protein
2. Edich, M., Briggs, D. C., Kippes, O., Gao, Y., & Thorn, A. (2022). The impact of AlphaFold2 on experimental structure solution. *Faraday Discussions*, 240, 184-195.
3. von Soosten, L. C., Edich, M., Nolte, K., Kaub, J., Santoni, G., & Thorn, A. (2022). The Swiss army knife of SARS-CoV-2: the structures and functions of NSP3. *Crystallography Reviews*, 28(1), 39-61.

### Published before PhD

1. Whitford, C. M., Dymek, S., Kerkhoff, D., März, C., Schmidt, O., Edich, M., ... & Kalinowski, J. (2018). Auxotrophy to Xeno-DNA: an exploration of combinatorial mechanisms for a high-fidelity biosafety system for synthetic biology applications. *Journal of Biological Engineering*, 12, 1-28.

## II. Table of Contents

<b>List of Publications</b>	<b>III</b>
<b>List of Abbreviations</b>	<b>VIII</b>
<b>Zusammenfassung</b>	<b>IX</b>
<b>Abstract</b>	<b>X</b>
<b>Part I</b>	
<b>Unlocking the Secrets of SARS-CoV-2: AlphaFold2-assisted Domain Determination of NSP3</b>	<b>1</b>
<b>1 Introduction and Objectives</b>	<b>2</b>
1.1 Introduction . . . . .	2
1.2 Objectives . . . . .	3
<b>2 Theoretical Background</b>	<b>4</b>
2.1 The coronaviral multi domain protein NSP3 . . . . .	4
2.1.1 Coronavirus biology . . . . .	4
2.1.2 Domains and functions of NSP3 . . . . .	6
2.1.3 Double membrane vesicles and virion assembly . . . . .	8
2.2 Protein analysis techniques . . . . .	11
2.2.1 Sequence analysis . . . . .	11
2.2.2 Experimental structure determination . . . . .	13
2.2.3 Prediction of protein structures . . . . .	14
2.2.4 Integrative modelling . . . . .	16
2.2.5 Structure similarity search . . . . .	17
<b>3 Results</b>	<b>18</b>
3.1 Utilizing AlphaFold2 for domain boundary determination and construct design . . . . .	18
3.1.1 Preliminary domain ranges . . . . .	18
3.1.2 Sequence alignments of preliminary domain ranges . . . . .	20
3.1.3 Comparison between AlphaFold2 predictions and experimentally determined structures . . . . .	20
3.1.4 Classification into regions of order and disorder . . . . .	23
3.1.5 Final determination of domain boundaries . . . . .	25
3.1.6 Structural similarity of final domains . . . . .	26
3.1.7 Structure similarity search . . . . .	31
3.2 Experimental validation of the Betacoronavirus-specific linker domain ( $\beta$ SLD)	31
3.2.1 Structure prediction of linker domain between PL2 <sup>pro</sup> and NAB domain . . . . .	31
3.2.2 Conservation of the linker domain . . . . .	32
3.2.3 Experimental validation . . . . .	34
3.2.4 Nucleic acid binding domain and Betacoronavirus-specific marker domain . . . . .	36
3.3 Hexameric pore complex . . . . .	37
3.3.1 Predicted arrangement of domains . . . . .	38
3.3.2 Structure prediction of Y1 and CoV-Y . . . . .	39
3.3.3 Conservation of Y1 among nidoviruses . . . . .	43

3.3.4	Multimer prediction of Y1 . . . . .	44
3.3.5	Multimer prediction of Y1+CoV-Y . . . . .	45
3.3.6	Integrative modeling . . . . .	47
3.4	Automated domain determination . . . . .	50
3.4.1	Implementation . . . . .	50
3.4.2	Benchmarking on database of multidomain proteins . . . . .	52
<b>4</b>	<b>Discussion</b>	<b>53</b>
4.1	Utilizing AlphaFold2 for domain boundary determination and construct design . . . . .	53
4.1.1	Classification of regions as ordered or disordered . . . . .	53
4.1.2	Domain boundary determination of NSP3 . . . . .	56
4.1.3	Automated domain determination . . . . .	57
4.2	Experimental validation of the Betacoronavirus-specific linker domain ( $\beta$ SLD) . . . . .	57
4.2.1	Structure analysis . . . . .	57
4.2.2	Potential functions of the linker domain, NAB, and $\beta$ SM . . . . .	59
4.3	C-terminal domains and the hexameric pore complex . . . . .	61
4.3.1	Nomenclature of Y1 and CoV-Y . . . . .	61
4.3.2	Structure prediction of Y1 . . . . .	61
4.3.3	Conservation of Y1 among <i>Nidovirales</i> . . . . .	62
4.3.4	Multimer prediction of Y1 and Y1+CoV-Y . . . . .	63
4.3.5	Experimental validation . . . . .	65
4.3.6	Integrative modeling . . . . .	66
4.3.7	Y1 as potential drug target . . . . .	70
<b>5</b>	<b>Outlook</b>	<b>72</b>
<b>6</b>	<b>Materials and Methods</b>	<b>73</b>
6.1	Utilizing AlphaFold2 for domain boundary determination and construct design . . . . .	73
6.1.1	Sequence information and preliminary domain ranges . . . . .	73
6.1.2	AlphaFold2 predictions . . . . .	74
6.1.3	Classification into ordered or disordered regions . . . . .	75
6.1.4	Structure Similarity Search . . . . .	75
6.2	Experimental validation of the Betacoronavirus-specific linker domain ( $\beta$ SLD) . . . . .	76
6.2.1	Structure prediction analysis . . . . .	76
6.2.2	Conservation of domain sequence . . . . .	76
6.2.3	Expression of Betacoronavirus-specific linker domain construct . . . . .	77
6.2.4	Purification of Betacoronavirus-specific linker domain construct . . . . .	78
6.2.5	Crystallization of Betacoronavirus-specific linker domain construct . . . . .	78
6.2.6	SAXS analysis of Betacoronavirus-specific linker domain construct . . . . .	79
6.3	Hexameric pore complex . . . . .	80
6.3.1	Prediction of multidomain segments . . . . .	80
6.3.2	Conservation of Y1 . . . . .	81
6.3.3	Multimer predictions of Y1 and Y1+CoV-Y . . . . .	81
6.3.4	Expression of Y1 . . . . .	82
6.3.5	Integrative modelling . . . . .	82

**Part II**

<b>Beyond the Surface: Exploring the Solvent Region of Macromolecular Crystals with Experimental Phasing</b>	<b>85</b>
<b>7 Introduction and Objectives</b>	<b>86</b>
7.1 Introduction . . . . .	86
7.2 Objectives . . . . .	86
<b>8 Theoretical Background</b>	<b>88</b>
8.1 Macromolecular Crystallography . . . . .	88
8.2 Experimental phasing . . . . .	88
8.2.1 Anomalous differences . . . . .	88
8.2.2 MAD phasing . . . . .	89
8.3 Crystallographic data . . . . .	90
8.3.1 Electron density maps . . . . .	91
8.3.2 Structure factor tables and contribution of separate components . .	93
<b>9 Results</b>	<b>95</b>
9.1 MAD phasing . . . . .	95
9.1.1 Experimental phasing via SHELXE . . . . .	95
9.1.2 Experimental phasing via SHARP . . . . .	97
9.2 Comparison of electron density maps . . . . .	99
9.2.1 Surface water comparison . . . . .	100
9.2.2 Similarity of partial maps . . . . .	101
9.2.3 Molecular dynamics simulation . . . . .	103
9.3 Refinement against electron density map from SHARP . . . . .	105
9.4 Analysis of solvent region contribution . . . . .	106
9.4.1 Map separation error . . . . .	106
9.4.2 Solvent contribution to structure factors . . . . .	108
9.4.3 Fifty percent anomaly . . . . .	110
<b>10 Discussion</b>	<b>112</b>
10.1 MAD phasing . . . . .	112
10.2 Comparison of electron density maps . . . . .	114
10.3 Refinement against electron density map from SHARP . . . . .	117
10.4 Analysis of solvent region contribution . . . . .	118
<b>11 Outlook</b>	<b>120</b>
<b>12 Materials and Methods</b>	<b>121</b>
12.1 MAD phasing . . . . .	121
12.1.1 List of python scripts . . . . .	121
12.1.2 Filtering of datasets . . . . .	122
12.1.3 Conversion of data and labeling of columns . . . . .	124
12.1.4 Preparation via SHELXC . . . . .	125
12.1.5 Substructure identification via ANODE . . . . .	126
12.1.6 Phasing via SHELXE . . . . .	127
12.1.7 Quality assessment via cross correlation . . . . .	127
12.1.8 Automated preparation for SHARP . . . . .	128
12.1.9 Phasing via SHARP . . . . .	129

---

12.2 Comparison of electron density maps . . . . .	130
12.2.1 Surface water comparison . . . . .	131
12.2.2 Separation of protein and solvent region . . . . .	131
12.2.3 Cross correlation of partial maps . . . . .	133
12.3 Refinement against electron density map from SHARP . . . . .	134
12.4 Analysis of solvent region contribution . . . . .	135
<b>References</b>	<b>136</b>
<b>A Appendix</b>	<b>151</b>
A.1 Appendix of Part I: Beyond the Surface . . . . .	156
<b>A Acknowledgements</b>	<b>158</b>

### III. List of Abbreviations

<b>MHV</b>	the murine hepatitis
<b>NSP</b>	non-structural protein
<b>DMV</b>	double membrane vesicle
<b>Ubl1</b>	ubiquitin like domain 1
<b>HVR</b>	hyper variable region
<b>Mac1</b>	macrodomain 1
<b>Mac2</b>	macrodomain 2
<b>Mac3</b>	macrodomain 3
<b>DPUP</b>	domain preceding Ubl2 and PL2 <sup>pro</sup>
<b>Ubl2</b>	ubiquitin like domain 2
<b>PL1pro</b>	papain-like protease 1
<b>PL1pro</b>	papain-like protease 2
<b><math>\beta</math>SLD</b>	betacoronavirus-specific linker domain
<b>NAB</b>	nucleic acid binding domain
<b><math>\beta</math>SM</b>	betacoronavirus-specific marker domain
<b>TM1</b>	transmembrane domain 1
<b>TM2</b>	transmembrane domain 2
<b>AH1</b>	amphipathic helix 1
<b>Y1</b>	nidovirus-conserved domain of unknown function
<b>CoV-Y</b>	Coronavirus-specific C-terminal domain
<b>RMSD</b>	root mean square deviation

## Zusammenfassung

Proteine erfüllen in jeder bekannten Lebensform zahlreiche essentielle Funktionen, finden durch die moderne Biotechnologie aber auch in experimentellen, medizinischen und industriellen Anwendungen Verwendung. Wissen über die Struktur und Funktion kann unter anderem genutzt werden, um Proteine nach unseren Anforderungen zu optimieren und auch um Medikamente gegen Proteine zu entwickeln, welche einem Pathogen entstammen.

Im Rahmen dieser Arbeit wurde das Multi-domain Protein NSP3, ein Coronavirusprotein, mithilfe von bioinformatischen Methoden und Strukturvorhersage untersucht. Die Ergebnisse umfassen den Nachweis einer bisher unentdeckte Domäne und diverse Erkenntnisse über mögliche Funktionen weiterer Domänen. Zudem wurde der manuelle Prozess der hier zur Domänenidentifizierung genutzt wurde automatisiert, um auch bei beliebigen Multi-domain Proteinen ein besseres Verständnis der Struktur zu bieten.

Im zweiten Teil dieser Dissertation wurde die Elektronendichte des Solvensbereichs von Proteinkristallen untersucht. Speziell wurden Reflexdaten automatisiert experimentell phasiert und die resultierenden Dichtekarten wurden mit veröffentlichten Karten verglichen. Die gewonnen Erkenntnisse wurden genutzt, um neue Wassermodelle zu konstruieren, welche zu einer verbesserten Struktur mit niedrigeren R-werten geführt hat.

**Keywords:** SARS-CoV-2, NSP3, AlphaFold2, Macromolecular X-ray Crystallography, Experimental Phasing, Solvent Model.



## Abstract

Proteins fulfil various essential functions in all known life forms and modern biotechnology enabled their use in experimental, medical, and industrial applications. Knowledge about the structure and function can be used to optimize proteins for these applications and finds use in drug discovery, if a the protein originates from a pathogen.

In this work, the multidomain protein NSP3 from coronaviruses was examined with bioinformatical methods and with structure prediction. The results include the identification and experimental validation of a previously unknown domain and a new hypothesis about functions of additional domains. Furthermore, the manual process of a new domain determination technique was automated to enable use in any multidomain protein.

In the second part of this dissertation, the electron density of solvent regions from macromolecular crystals was examined. Specifically, reflection data was experimentally phased in a newly developed automated pipeline. The resulting density maps were compared to published maps and the new insights enabled the construction of new water models, which improved one structure and lowered its R-values.

**Keywords:** SARS-CoV-2, NSP3, AlphaFold2, Macromolecular X-ray Crystallography, Experimental Phasing, Solvent Model.



## **Part I**

# **Unlocking the Secrets of SARS-CoV-2: AlphaFold2-assisted Domain Determination of NSP3**

# 1 Introduction and Objectives

## 1.1 Introduction

Coronaviruses are a recurring global threat, with the Severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) causing COVID-19 (Coronavirus disease 2019) being the most recent example. Infected patients suffer from severe symptoms, which can result in lethal damage of lung tissue and the cardiovascular system [1]. Therefore, medical treatments that reduce the severeness of COVID-19 infections are in high demand.

While numerous potential drugs have been identified against the main protease [2, 3] and against PL2<sup>pro</sup> [4, 5], alternative therapeutics are still desired. First of all, viruses mutate rapidly and vaccine-induced antibodies provide a worldwide selection pressure, raising the risk of immune and vaccine escape variants [6–8]. The availability of multiple drug treatments would reduce the success of such escape variants. Secondly, some proteins and domains are conserved to varying degrees among related viruses. Identifying drug targets which are also well conserved could therefore lead to therapeutics that are potentially applicable against multiple viruses and even future pathogens [9, 10]. For example, current treatments against PL2<sup>pro</sup> target SARS-CoV-2 and SARS-CoV-1, but are ineffective against MERS-CoV [5], illustrating the narrow focus of certain drugs. Finally, therapeutics can cause severe side effects in some patients and may be incompatible with other prescribed medications [11, 12]. In such cases, access to alternative treatments is a necessity.

Considering all these aspects, a complete exploration of the structure and function of all proteins of SARS-CoV-2 is of high relevance. This work focuses on non-structural protein 3 (NSP3), the largest *Coronavirus* protein that consists of 17 domains, of which only a fraction have been structurally solved for SARS-CoV-1 and SARS-CoV-2 [13]. In addition, most research has focused on two enzymatic domains of NSP3 with essential functions, macrodomain 1 and PL2<sup>pro</sup>, as reflected in the number of structures in the Coronavirus Structural Task Force database [14]. Of the remaining domains, only a few have been biochemically characterized and for some domains neither the structure nor the function is known [13]. However, understanding the function of these domains might help to identify potential targets.

Determining the protein structures of the remaining domains can be achieved by several orthogonal techniques. The most recent technology is AI-based structure prediction, where artificial neural networks are trained on the structure models obtained from experimental methods over the past decades, where AlphaFold2 was the first reliable structure

prediction tool [15, 16]. Since the introduction of this technology, we have to differentiate between purely computationally generated structure models and models based on experimentally obtained data, where only the latter works with the real sample and all of its properties. Although predicted structures lack data of the actual protein and its environment, the method shows practical use cases in supporting the structure solution by experimentally based methods, where the design of crystallizable protein constructs is one example.

## 1.2 Objectives

The current research on NSP3 is spread across numerous publications and the latest review paper summarizing these efforts to a practical overview was published before the pandemic [17], thus excluding SARS-CoV-2. Furthermore, the available technology and definition of NSP3 domains have changed over the decades, leading to conflicts in naming conventions and residue ranges of domains. Therefore, this work aimed at expanding the available knowledge about NSP3 with published structures of SARS-CoV-2 and the assistance of structure prediction, which became available at the beginning of this project.

Currently available sequence information as well as structure prediction results suggested the presence of an unnoticed folded domain in a linker sequence, next to the well-researched drug target domain PL2<sup>pro</sup>. This domain, later specified as the betacoronavirus-specific linker domain, was investigated experimentally in order to validate its folding nature. The sequence and predicted structure of another domain, known as "nidovirus-conserved domain of unknown function" [18], were analysed for possible functions. The results of this work indicated that it was likely involved in the assembly of a multimeric complex and in the export of viral RNA, which had not been reported previously. Another goal of this work was to fit structure models of each domain into an available cryo electron tomography map of the pore complex, which consists mainly of NSP3 [19]. However, the used integrative modeling approach required precise knowledge about the length of each domain and their connecting linkers. Therefore, the first objective was to map each residue of NSP3 to a single domain or linker, which was done in parallel with summarizing all current information about NSP3 in a new review paper and resolving nomenclature conflicts. Accomplishing this goal is possible with AlphaFold2 [15], which is able to indicate intrinsic disorder in predicted structures [20]. Finally, it is desirable to automate this new method of domain determination and enable this approach for any multidomain protein.

## 2 Theoretical Background

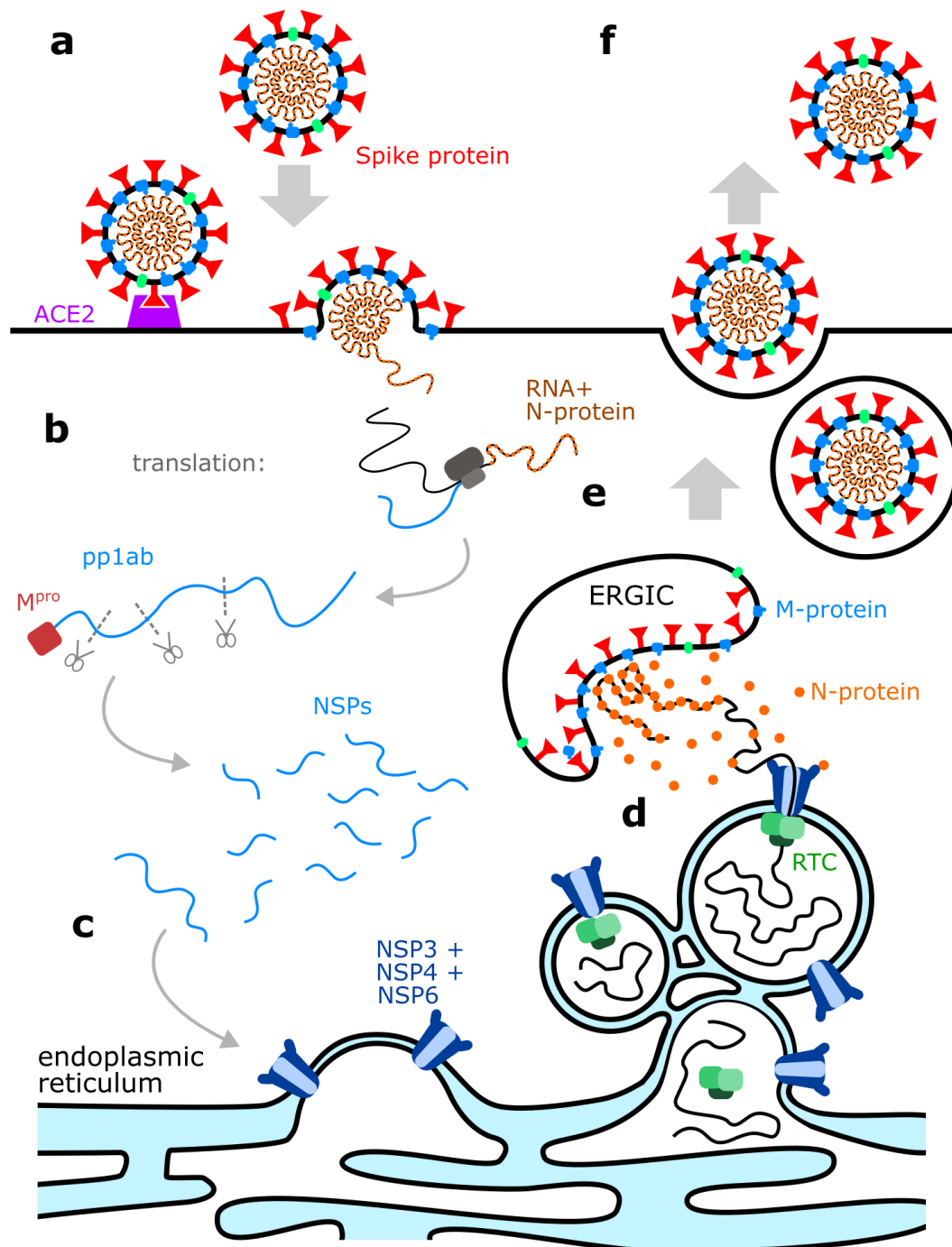
### 2.1 The coronaviral multi domain protein NSP3

#### 2.1.1 Coronavirus biology

SARS-CoV-2 (Severe acute respiratory syndrome coronavirus type 2) is the famous pathogen that caused a global pandemic in early 2020. It belongs to the viral family *Coronavirinae* and to the genus *Betacoronavirus*. These viruses are positive-stranded RNA viruses with genome sizes of 27kb to 32 kb, which are among the largest RNA virus genomes [21]. Related viruses are SARS-CoV-1 and the murine hepatitis virus (MHV), which are of great interest in SARS-CoV-2 research as they share similar proteins and biochemical functions. MHV is also practical to use for experiments, since it is not harmful to humans and requires reduced safety measures [22].

Upon infection, coronaviruses release their RNA genome into the cytosol of the host cell (Figure 1a), where it is translated into the viral proteins [23, 24]. The first proteins to be translated are the large polyproteins pp1a and pp1ab (Figure 1b), which are chains of so-called non-structural proteins (NSPs). The first one, pp1a, is cleaved into the proteins NSP1 to NSP10, whereas pp1ab is an extended version of pp1a and comprises NSP1 to NSP16. Structural proteins on the other hand, are translated as individual proteins and function as structural components of new virions [23], which are the spike-, nucleocapsid-, membrane-, and envelope proteins. The NSPs play various roles in the viral replication cycle, with the viral proteases M<sup>Pro</sup> and papain-like proteases on NSP3 performing one of the most important tasks: cleaving the polyproteins into individual NSPs [25, 26].

Replication of the viral genome takes place inside double membrane vesicles (DMVs). These vesicles are the product of interaction between NSP3 and NSP4 (Figure 1c), which cause the membranes of the endoplasmic reticulum to rearrange [27]. Inside the DMVs, RNA is replicated by the viral replication-transcription complex (Figure 1d) consisting of the proteins NSP7, NSP8, NSP12, and NSP13 [28, 29]. The DMVs protect the replication-transcription complex and the viral genome from the host's viral interference system [30], but to progress the infection cycle, the RNA must be exported into the cytosol where the structural proteins assemble into new virions [31]. Key to the export of genome copies are the hexameric pore complexes formed by NSP3, NSP4, and NSP6 [19, 32]. NSP3 makes up the majority of these large complexes and interacts with the nucleocapsid proteins [33], which attach themselves to the viral RNA to prevent RNA degradation by host proteins [34, 35]. In the cytosol, the viral genome coated with nucleocapsid proteins is moving towards the assembly-site of new virions [36], where the structural proteins induce



**Figure 1:** *Replication cycle of SARS-CoV-2.* Virions enter the host cell via interaction between the spike proteins and host's ACE2 receptor (a). Upon infection, the viral RNA genome coated with N-proteins (nucleocapsid) is released into the cytosol and translated into viral proteins by the host's ribosomes, where the NSPs (non-structural proteins) are translated into a large polyprotein, pp1ab (b). The polyprotein is cleaved into individual NSPs by the viral proteinase M<sup>pro</sup> (c). The proteins NSP3, NSP4, and NSP6 integrate into the membrane of the endoplasmic reticulum, where NSP3 cleaves NSP1 to NSP4 from the remaining polyprotein. The interaction between NSP3 and NSP4 induces membrane curvature, which creates double membrane vesicles (d). The RTC (replication-transcription-complex), ribonucleotides, and viral genome are inside these vesicles, allowing replication of the viral genome (d). New copies of the genome are exported into the cytosol through the pore complex consisting of NSP3, NSP4, and NSP6 (e). Here, the RNA is coated by N-proteins for protection from degradation and the RNA-N-protein complex is bound by M-proteins at the ERGIC (endoplasmic reticulum-Golgi intermediate compartment) membrane during virion assembly (e). These compartments enclose the genome, resulting in new virions encapsulated in vesicles, which allows secretion out of the cell (f).

the formation of virions from the membrane of the host’s endoplasmic reticulum-Golgi intermediate compartment (Figure 1e) and fix the genome inside the virion’s interior [37–39]. Last but not least, the virions induce a self-secretion out of the host cell [40, 41], from where they can infect new cells and restart the infection cycle (Figure 1f).

### 2.1.2 Domains and functions of NSP3

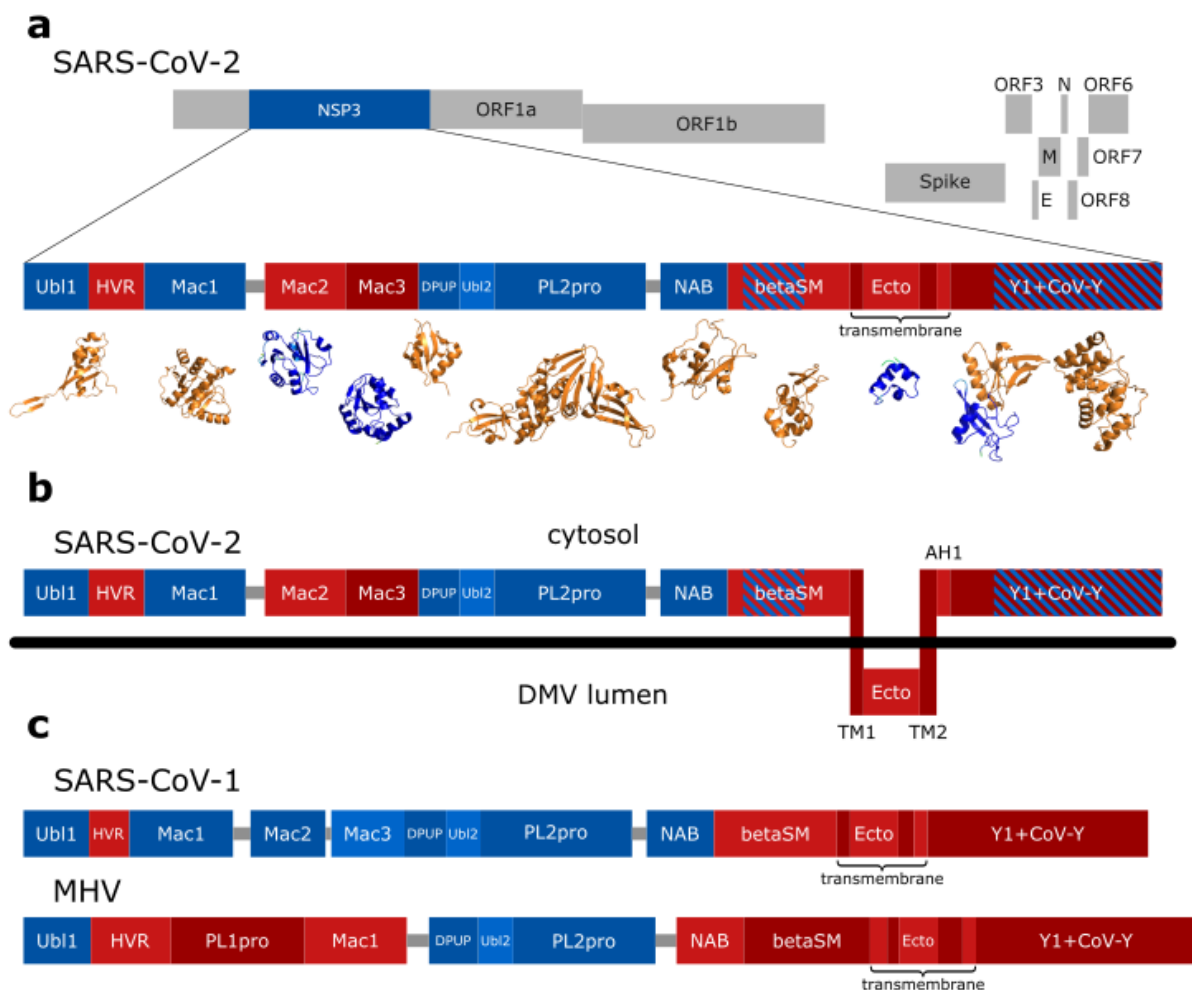
Non-structural protein 3 (NSP3) is the Swiss army knife of SARS-CoV-2, as its 1945 amino acid residues fold into 17 domains, each performing a unique task [13]. It is the largest of all coronavirus proteins and many of its domains are essential for viral replication [17]. Among the domains are two transmembrane helices, which anchor the entire protein to the membrane of the DMVs (double membrane vesicles) [42], where most of NSP3 is on the cytosolic side and only the ectodomain is in the lumen of the DMV [17]. This ectodomain interacts with NSP4 to form the DMVs [27], and together with NSP6 they assemble into the RNA-exporting pore complex with sixfold symmetry [19, 32]. Without NSP3, the infection cycle cannot progress due to its important role in RNA export and its ability to cleave NSPs from the polyprotein [19, 25]. It also facilitates evasion of the immune response by reversing processes used in anti viral defense [43–45]. All in all, the domains of NSP3 provide several drug targets in the fight against COVID-19 [4, 5, 17, 27, 46, 47].

The 17 domains of NSP3 fulfill various roles, but only a few domains are understood in detail. The N-terminal domain is the **ubiquitin like domain 1 (Ubl1)**, which is essential for the virus [47]. It reaches far into the cytosol [19] where it interacts with the nucleocapsid protein [33] and probably brings it into contact with the exported RNA [19, 33]. Ubl1 is followed by a disordered domain of 95 residues, the **hypervariable region (HVR)** [13]. It has a high content of glutamic acid and aspartic acid, allowing it to mimic DNA and RNA [48]. A recent study shows that HVR, like Ubl1, also interacts with the nucleocapsid protein, although the interaction is weaker [49] and HVR is dispensable for the virus [47]. Despite its potential intrinsic disorder and large variations in length and sequence, it is present in all coronaviruses [17].

The next domain in SARS-CoV-2 NSP3 is **macrodomain 1 (Mac1)**, which reverses a mechanism involved in antiviral defense, the human PARP14-derived ADP-ribosylation [50]. It is often referred to as ADP-ribose phosphatase [17] and its role in immune evasion has made it one of the two well-studied drug targets of NSP3 [51]. While it is dispensable for viral replication *in vitro* [52], it is shown to be essential *in vivo* [53]. Macrodomain 1 is



followed by a 33-residue linker and a region previously known as SARS-unique domain [54], consisting of the three domains **macrodomain 2 (Mac2)**, **macrodomain 3 (Mac3)**, and the **domain preceding Ubl2 and PL2<sup>pro</sup> (DPUP)**. However, murine hepatitis virus (MHV) has been shown to possess a DPUP-like domain [55], leaving only Mac2 and Mac3 unique to *Sarbecovirus*. While the function of DPUP is unclear, Mac2 and Mac3 are binding oligo(G)-nucleotides [52, 56]. This property makes Mac3 indispensable for viral replication, in contrast to Mac2 and DPUP [52]. DPUP is suspected to assist in RNA binding as it shows a positive charge at neutral pH [57].



**Figure 2:** Figure 1 Domain overview of NSP3 from both sarbecoviruses and MHV. (a): Position and size of NSP3 on the polyprotein, as well as all domains (blue/red boxes) and larger linkers (gray lines). Blue domains correspond to experimentally solved structures in the PDB; red domains are not experimentally solved; red domains with blue stripes are partially solved. 3D structures shown in orange are from the PDB (see Table 1 for PDB codes), while the other structures are predicted by AlphaFold2 and coloured according to their pLDDT, with blue representing high confidence with pLDDT values greater than 90. (b): Membrane topology of SARS-CoV-2, based on the results of Oostru et al. [42]. (c): Domains of SARS-CoV-1 and MHV. Depicted domain ranges are mentioned as preliminary ranges in the results and are listed in Table 2. Domain ranges of SARS-CoV-1 are according to Lei et al. [17]; ranges of SARS-CoV-2 and MHV are based on sequence alignments with those of SARS-CoV-1 and experimental structures.

DPUP is followed by the **ubiquitin-like domain 2 (Ubl2)** and the **papain-like protease 2 (PL2<sup>pro</sup>)**, whereas Ubl2 is seen as a subdomain of PL2<sup>pro</sup>. This protease cleaves the non-structural proteins NSP1, NSP2, NSP3, and NSP4 from the polyproteins, making PL2<sup>pro</sup> an essential domain for the virus [17]. This property has made it one of the two most researched drug targets on NSP3 [25].

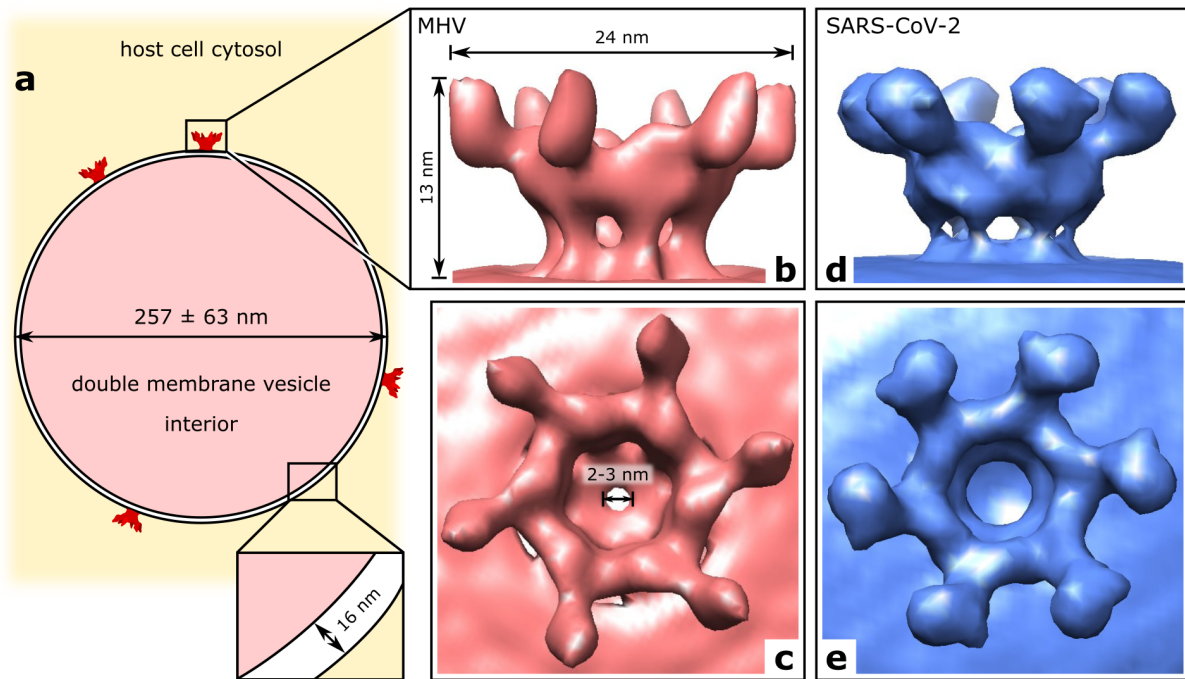
PL2<sup>pro</sup> and its subsequent domain, the **nucleic acid binding domain (NAB)**, are connected by a 34-residue linker. However, as shown later in this work, this "linker" is a folded and conserved domain specific to *Betacoronavirus*. Therefore, it has been labeled as **betacoronavirus-specific linker domain ( $\beta$ SLD)**. NAB is also specific to *Betacoronavirus* and is followed by a third region specific to *Betacoronavirus*, the **betacoronavirus-specific marker domain ( $\beta$ SM)**, which consists of a folded core surrounded by large disordered linkers [17, 58].

After  $\beta$ SM begins the transmembrane region followed by the C-terminal domains, where the former consists of two transmembrane domains, an amphipathic helix, and the only luminal domain of NSP3, known as **ectodomain (3Ecto)** [17, 42]. The C-terminal region consists of the **Y1** domain, known as nidovirus-conserved domain of unknown function, and the **CoV-Y** domain, which is specific to *Coronavirus* [18]. Recently published structures (PDB codes 7RQG [59]; 8F2E [60]) suggest that both domains consist of two subdomains.

While SARS-CoV-1 NSP3 is highly similar in its sequence compared to SARS-Cov-2 NSP3 and identical in domain composition [13], murine hepatitis virus (MHV) NSP3 differs by lacking the domains Mac2, Mac3, and by having the domain PL1<sup>pro</sup> between HVR and Mac1 in addition to PL2<sup>pro</sup> [17, 56] (Figure 2). PL1<sup>pro</sup> cleaves NSP1 from the polyprotein, whereas MHV PL2<sup>pro</sup> frees NSP2, NSP3, and NSP4 [17]. In *Sarbecovirus*, PL2<sup>pro</sup> takes both functions and cleaves NSP1 to NSP4, from the polyprotein [17].

### 2.1.3 Double membrane vesicles and virion assembly

Several taxonomic orders of positive single stranded RNA (+ssRNA) viruses initiate the formation of replication organelles in the infected host cell [61–64]. Upon infection, the membrane of the endoplasmic reticulum is rearranged and forms spherical double membrane vesicle (DMVs), which shield the replication of viral RNA from interference by host enzymes or RNA, concentrate all components of replication in one place, and potentially prevent negative ssRNA copies of the genome from entering the cytosol and assembling into new virions [27, 63]. In betacoronaviruses, DMV formation is induced by interaction



**Figure 3:** Schematic illustration of a double membrane vesicle induced by murine hepatitis virus with RNA-exporter pore complexes embedded into its surface (a). The luminal spacing between the two membranes is 16 nm wide [32]. These transmembrane pore complexes consist of NSP3, NSP4, and NSP6, where the large cytosolic structure consists of a hexameric assembly of NSP3 (b and c). The shown structure volume is a cryo electron tomography map from Wolff et al. [19] at 30 Å resolution, shown with a contour level of 2.89. Both, side view (b) and top view (c) show the same map. The depicted diameter numbers are taken from Wolff et al. [19]. The blue map (d, e) shows the pore complex made of NSP3 and NSP4 from SARS-CoV-2 (EMDB entry EMD-15963) at a contour level of 0.47 [32]. Maps are visualized in UCSF Chimera [66], where contour levels use arbitrary units.

between the ectodomain of NSP3 and the ectodomain of NSP4, probably after both membrane proteins are embedded into the membrane of the endoplasmic reticulum [27, 65]. Mutation experiments have shown that this interaction is essential for DMV formation and that its absence disrupts viral replication in MHV [27], making both ectodomains potential drug targets. However, the presence of such replication organelles raises the problem of exporting the positive RNA strand from the DMV’s interior into the cytosol and transporting it to the virion assembly site, where in SARS-CoV-2 this includes packaging of RNA into a coat of nucleocapsid proteins (N-protein) [36].

While it is not clear how RNA export from DMVs works in most viruses, the responsible structure in MHV was successfully imaged via cryo electron tomography (cryoET) [19], which is depicted in Figure 3c-b. It is a hexameric pore complex consisting of NSP3, NSP4, and NSP6, the three membrane NSPs of coronaviruses [19]. It contains six copies of NSP3, NSP4 and is assumed to contain also six copies of NSP6, but molecular weight

estimates leave room for additional proteins, which could in theory originate from the host cell [19, 32] as the complex is located on the cytosolic surface of the DMVs (Figure 3a). The exact position and orientation of each NSP3 domain within the complex is currently unknown. The location of Ubl1 in the outermost extensions is roughly known, as an NSP3 mutant which is N-terminally fused to GFP (green fluorescent protein) was also imaged using cryo electron tomography and fluorescence imaging [19]. A recent study [32] identified NSP3 and NSP4 as minimum constituents of the pore complex and provided a cryoET image of the complex from SARS-CoV-2 (Figure 3d-e). Furthermore, they identified via deletion mutants the N-terminal domains Ubl1 to Ubl2 to be crucial for the formation of the complex's crown and identified the domain Mac1 to be also located in the prolongations together with Ubl1 [32].

After exporting the RNA into the cytosol, it must be packaged with N-proteins, which protect the viral genome from interaction with host enzymes and interference RNA [34]. This step must occur before the RNA is transported to the virion assembly site, because M-protein, which fixes the genome in the virion's interior, only binds RNA packed in N-protein [39]. The binding of N-protein by Ubl1 is likely driving the RNA packaging right after export across the DMV membranes [33].

In SARS-CoV-2 and potentially in all coronaviruses, new virions are formed by the interaction between the structural proteins, where the membrane protein (M-protein) and envelope protein (E-protein) induce membrane rearrangement of the endoplasmic reticulum-Golgi intermediate compartment (ERGIC) membrane [38, 67]. The M-proteins arrange in a scaffold formation on the virion membrane, which enhances the assembly of spike proteins [31, 68] and the M-protein's endodomain, located in the virions interior, binds the viral RNA packaged in N-proteins [39]. Conclusively, NSP3 plays an important role in the viral replication cycle by exporting the RNA and providing contact to N-protein, before the genome is packaged and moved towards the virion assembly site.

## 2.2 Protein analysis techniques

Protein structures and sequences can be analyzed experimentally and with bioinformatical tools [69]. Because each methods has its own advantages and disadvantages, one must adapt the choice of tools and techniques to the current problem and the desired information outcome. Structures of highly complex proteins such as NSP3 cannot be solved by one or few techniques. Therefore, it is necessary to apply multiple methods, to work in collaborations, and to combine information from various sources.

While sequence analysis reveals information about the evolutionary history of a protein, it does not directly provide details about a protein's fold. Structure prediction makes the best use of this sequence information by identifying co-evolving residues, but it is limited and biased by the available data from structure databases. Furthermore, it lacks a connection to experimental evidence of the target structure. Atomic resolution structures are obtainable from macromolecular crystallography, but these require the structure to be crystallizable in the first place. While this is sufficient for single proteins and domains, it becomes an obstacle for multidomain proteins, membrane proteins, molecules with large amounts of intrinsic disorder, and assemblies of multiple proteins or domains - all properties which are given for NSP3 [17]. Nuclear magnetic resonance (NMR) spectroscopy can handle proteins with high flexibility and intrinsic disorder, but is limited to smaller structures. Solution small-angle X-ray scattering also works with non-crystallizable proteins, but it only provides limited structural information. For the analysis of large complexes, cryo electron microscopy and cryo electron tomography are the methods of choice, but in such cases they are limited to relatively low resolution. Finally, integrative modelling combines the best of all worlds. Here, cryo electron microscopy provides the low resolution scaffold, which is then augmented by atomic-resolution structures from macromolecular crystallography, NMR, and structure prediction. Additional information from experiments about the interaction and assembly of proteins is used for modelling restraints. All data combined leads to a model of a protein complex which serves as the most plausible hypothesis.

### 2.2.1 Sequence analysis

Proteins are chains of amino acids, with each of the twenty amino acids having its own chemical properties. Ultimately, the sequence of these amino acids defines a protein's three dimensional shape and thus its function. Therefore, it makes sense to compare amino acid sequences when comparing two similar proteins. A rule of thumb is that highly similar amino acid sequences result in highly similar protein folds [70]. However,

similar folds can also have very different sequences [71]. In any case, evolutionary related proteins with similar sequences, known as homologs, can be expected to not only share a similar fold, but also similar functions and some bioinformatical tools aim at exploiting this property [72, 73].

The most common measurement is to perform a sequence alignment, which is an algorithm used to identify and overlap identical or similar regions of two or more sequences. From such sequence alignments, one can calculate the sequence identity, which is the percentage of identical amino acids at the same common position, or the sequence similarity, where amino acids must share common chemical properties. Sequence similarity is thus always equal to the sequence identity or greater. Sequence alignments are either performed locally or globally, where the former is used for pairwise alignments of two sequences and the latter is optimized for the alignment of multiple sequences [74, 75]. In addition to analyzing the similarity between known sequences, one can also perform a BLAST search [76] on the NCBI servers [77] to find potential homologs from other viruses or organisms by searching against a large database.

Often, only a fraction of a protein's sequence shows high similarity to other proteins. This is especially true for multidomain proteins, which consist of multiple folded structures connected by linkers and are part of the same amino acid chain. While folding regions of such sequences are considered ordered, other segments may not fold because they are intrinsically disordered [78]. Such disordered regions are often found in linkers between domains or as disordered tails at the N- or C-terminus, but large disordered regions can also be found as non-folding domains if they serve a biological function [79]. Because ordered regions have a high selection pressure against mutations to maintain a functional fold, especially against insertions or deletions, such regions share usually a higher sequence similarity with homologs in contrast to disordered regions, where all kinds of mutations are allowed because there is no fold to disrupt [78]. However, it became evident that disordered domains which serve a function are also susceptible to mutations and are thus in contrast to linkers, which are primarily maintaining flexibility and a certain length. [80].

Multidomain proteins are difficult cases, as they have large sequences which often comprise ordered folds and disordered regions in the form of linkers. The *Coronavirus* protein NSP3 is such a protein and contains with HVR and  $\beta$ SM also two large regions of intrinsic disorder [17]. Sequence alignments and BLAST searches [76] are a common way of identifying domains in a sequence-based approach, but they are limited to the availability

of sequence homologs in the searched databases [81] and as mentioned before, similar folds may result from different sequences. Structure-based domain identification relies on similar structures, but these are available in fewer numbers than similar sequences [81]. On the opposite, domains can be validated, but also identified by experimental structure solution methods, as explored in the following subsections. In contrast to the conventional methods, this work describes the development of a sequence- and structure-based approach, where the structure prediction tool AlphaFold2 [15] is utilized.

### 2.2.2 Experimental structure determination

Most of the protein structures deposited in the Protein Data Bank (PDB) [82] are modelled from experimental data collected on protein crystals at a beamline, where this technique is known as macromolecular crystallography. Data from this method enables the reconstruction of electron density maps, which then serve as guideline for construction of a protein model. The resulting protein models reach atomic resolution, but to get there, several conditions must be fulfilled. One of the most relevant ones is the need for crystallizable protein constructs [83]. A protein crystal is an assembly of proteins that can periodically expand and repeat its pattern in all three dimensions. However, this is only possible, if the protein of interest consistently folds into the same shape and into a stable state at all [84, 85]. Regions of intrinsic disorder may hinder the formation of large crystals due to their flexibility, thus breaking the repeated periodic expansion. Therefore, X-ray crystallography is not suited for solving the structure of disordered proteins and entire multidomain proteins which contain flexible elements.

Since growing crystals which diffract well to high resolution can be a time and cost consuming process, it is helpful to know beforehand some properties of the target protein such as the presence of intrinsic disorder. Solution small-angle X-ray scattering (SAXS) solves this problem by measuring protein samples in solution and thus much closer to their native state, but comes with the trade-off of lower resolution ranging from 10 Å to 20 Å [86]. One major difference to X-ray crystallography is that the molecules are not fixed in position and orientation in a 3D lattice, but that they float freely in the solution. Therefore, measuring data by directing an X-ray beam of a controlled wavelength at the sample does not result in discrete reflections in the three-dimensional reciprocal space. Instead, the freedom in movement of the molecules leads to a scattering signal that can be averaged radially, resulting in a one-dimensional curve [86]. Because not only the target molecule, i.e. the protein, but also solvent molecules such as water and buffer are

measured, a second measurement of just the solution allows to subtract the background information and obtain a signal related to the protein only [86].

### 2.2.3 Prediction of protein structures

The amino acid sequence of a protein is the most important factor in the folding of the chain of residues, which results in a functional, three-dimensional protein structure [87, 88]. The other important factor is the environment in which an amino acid chain folds [87, 88]. Knowing how these factors drive the folding could therefore lead to computational simulations and predictions of folded polypeptides and since solving a protein structure by experimental methods is a costly and time-consuming procedure, such technology is very desirable.

Algorithmic approaches for protein fold prediction come in various ways, such as simulation of thermodynamics and molecular dynamics [89–91] or knowledge-based approaches [92–95] utilizing protein structure data available from decades of experiments. Today, convolutional neural networks are a common architecture of artificial neural networks [96, 97], which are capable of learning a function that can reliably map input data to the expected output data. This works also for data that was not present in the dataset used to train the neural network, thus allowing the network to not only reproduce known results, but also to make predictions from unseen information. In the case of protein structure prediction, the protein data bank (PDB) [82] provides structures combined with sequences and additional information for thousands of proteins, which could be used as such training data.

All attempts of developing a reliable protein structure predictor remained unsuccessful until 2018, when Deepmind’s AlphaFold [98] achieved the highest scores at CASP13, the Critical Assessment of Structure Prediction [99]. Two years later, AlphaFold2 [15] reached even higher scores at CASP14 and for the first time in history, a software was able to predict protein structures with near experimental accuracy from sequence information alone in most of the test cases [15, 16]. Free public access launched in 2021 and alternative structure prediction tools, namely RoseTTAFold [100] and ESMFold [101], followed shortly after. Due to some additional features, the focus of this thesis is put on AlphaFold2.

Most, if not all prediction tools for protein structures have in common that they predict a 3D protein model from a given amino acid sequence. AlphaFold2 provides in addition a confidence score for each residue, the predicted Local Distance Difference



Test (pLDDT), which is used for quality control [15]. This score ranges from 0 to 100, where values above 90 are considered highly confident and values below 50 are considered non-confident, meaning that the respective region of the fold is likely predicted incorrectly. Values above 80 are still considered confident, although such residues may deviate from experimentally determined structures in comparison. Another metric, the predicted Aligned Error (pAE) assesses how confident the fold is globally for single domain proteins [102]. For multidomain proteins, it is used to evaluate how two domains are arranged relative to each other. Furthermore, AlphaFold2 is able to predict multimeric assemblies of proteins [103]. Here we get the same metrics, but the pAE can also be used to evaluate whether the arrangement of monomers to each other is random or confident.

In a nutshell, AlphaFold2 performs a multiple sequence alignment (MSA) with the query sequence against a sequence database and feeds this information into the first of two neural networks, the Evoformer block [15]. This network with attention and non-attention architecture [104] recognizes conserved and co-evolving residues from the MSA and feeds the processed information into the second network, the structure module, which translates and rotates each residue in 3D space [15]. The structure prediction runs both networks in several iterations and refines the resulting model with each additional loop. The pLDDT is calculated by an additional small network at the end for each residue [15].

The predicted structures are highly accurate and comparable to experimentally determined structures when sufficient sequence coverage is available in the MSA [15, 105]. Additional limitations are alternative conformational states, highly flexible regions, and transmembrane domains. Protein folds depend not only on the amino acid sequence but also on the environment, which can lead to alternative stable conformations [85, 106–108]. These, however, can still be obtained by tweaking AlphaFold2’s settings, although the mechanics behind this procedure are not straightforward and are limited to proteins with specific properties [109–111]. Also, important environmental factors such as membranes are absent, which leads to incorrect relative arrangements of domains when predicting a multidomain protein with transmembrane domains [112]. Furthermore, the PDB contains primarily globular and crystallizable proteins solved by X-ray crystallography [113], which do not represent all possible forms of proteins. The training data on such proteins is therefore limited, which can have a negative impact on the prediction of such structures. However, low pLDDT values were shown to correlate with disordered regions [20], which can be used in the design of crystallizable constructs [114]. A recent assessment of the accuracy of AlphaFold2 [115] shows that more than 80 % of predicted residues match near perfectly with experimentally determined folds and more than 94 % match those at

least roughly within an error margin. This high accuracy enables application in molecular replacement, where predicted structures are used for phase estimates in macromolecular crystallography. The downsides of this method are explored in the second part of this dissertation (section 7).

All in all, AlphaFold2 predicts protein structures in most of the cases correctly, especially if it is a globular protein under common physiological conditions. Together with its non-intended feature of predicting intrinsic disorder via low pLDDT values makes it a useful tool in the design of crystallizable constructs. However, its nature as a black-box has made it a new subject of research and its behaviour and limitations are still being explored, making it only an assisting tool instead of a replacement for experimental methods.

#### 2.2.4 Integrative modelling

Integrative modeling combines data from various experimental techniques and theoretical models to determine the structure of proteins or protein complexes, which are otherwise unobtainable [116, 117]. Low-resolution models such as volumetric data from cryo electron microscopy or tomography serve as restraints for the position and orientation of individual proteins or domains within the volume [116, 118].

One integrative modelling tool is Assembline from the Kosinski lab [119]. In a nutshell it uses a 3D map reconstructed from single particle cryo electron microscopy (cryoEM) of a multidomain protein or a protein complex and fits experimentally determined or predicted structures into this map, where additional information such as domain ranges or cross-linking results from experiments are considered as restraints [120].

In cross-linking mass spectrometry, covalent bonds are formed at the interface between two interacting proteins or domains [121, 122]. The amino acid chains of all proteins are then digested, where the previously established covalent cross links remain. Analysis of all peptides via mass spectrometry provides a library of short sequences and where cross links have occurred, which can be traced back to exact location of cross links and bound protein partners in the structures.

For the fitting into the cryoEM map, a low-resolution volume is generated for each structure [119]. These are then randomly placed throughout the map and evaluated by cross correlation, which generates a table of best fits, the fit library. This step is performed for each structure separately. In the global optimization step, best fits and constraints from cross-linking are combined to find optimal fits of all structures together,

which results in rough positions and orientations of each structure within the complex. Finally, all parameters are refined to obtain atomic-resolution models of the complex. The success of such a modeling approach depends heavily on the complexity of the assembly, the available data for restraints, and the resolution of the cryoEM map [116, 119].

### 2.2.5 Structure similarity search

While the previous methods are suitable for determining a protein's structure, they do not provide much information about the function of a protein. Sequence analysis can potentially identify active sites of an enzyme by finding conserved amino acid residues [123]. However, it is more practical to work with a structure, since the fold defines the function.

A data base search for similar structures is a promising option, as similar structures may consist of vastly different sequences and are hence invisible in sequence-based searches [71]. Because it is not feasible to perform such a search on structures composed of hundreds to thousands of atoms, a density distribution volume is generated and compared instead [124]. The Protein Data Bank [82] provides such a search tool [124], where the volumes are described as BioZernike descriptors [125], which allow an efficient search across all deposited structures of the PDB [124]. In the case of enzymes, one can alternatively search directly for similar active sites, which does not limit the search to similar protein shapes. One method enabling such a search is the genetic active site search, which uses genetic algorithms to find non-exact amino acid matches [126].

## 3 Results

Non-structural protein 3 (NSP3), the Swiss army knife of *Coronavirus* [13], plays an essential role in the viral replication cycle with its 17 domains [4, 5, 25, 27, 33, 44–47, 52, 53, 56, 65, 127]. Not all domain structures are experimentally solved, but structure prediction can provide insights into the unexplored folds. Since previous domain ranges were partially contradictory with experimentally solved structures, incomplete for SARS-CoV-2, or included large regions of intrinsic disorder, the domain ranges were redefined. Furthermore, the used technique was developed into a new method for *ab initio* domain determination in multidomain proteins. Moreover, the presence of a previously unnoticed domain was identified and experimentally validated. With current experimental data from multiple sources, prediction of a hexameric complex, bioinformatical analysis, and logical reasoning, it was also possible to deduce one potential function of the C-terminal Y1 domain, which was previously known as "nidovirus-conserved domain of unknown function".

### 3.1 Utilizing AlphaFold2 for domain boundary determination and construct design

#### 3.1.1 Preliminary domain ranges

The definitions of number, location, and residue ranges of domains from *Coronavirus* NSP3 changed significantly over time [18]. Experimental results led to many revisions, which also affect nomenclature. Today, unifying all present information is a difficult task as most domains are described under multiple names [13, 17] in addition to contradictory nomenclature in the current literature and incomplete gene annotations. Especially for SARS-CoV-2, no summarizing review paper was available and since the last review for SARS-CoV-1 NSP3, new domain structures were solved. Therefore, the first goal of this project was to update the domain ranges of SARS-CoV-1 NSP3 with recent experimental results and utilize structure prediction to identify domain borders in regions where structural information is absent. The murine hepatitis virus (MHV) is a close relative of the two sarbecoviruses and is considered harmless to humans, which makes it to a practical model organism [22]. Hence, domain ranges for MHV were updated as well.

At the start of this work, NSP3 from SARS-CoV-2 and MHV had the least number of experimentally solved domains in comparison with SARS-CoV-1 NSP3. Therefore, the domain ranges for SARS-CoV-1 were updated first, where ranges from Lei et al. [17] were taken as a starting point and if new experimental structures were available,

the respective domain ranges were updated. Transmembrane domain predictions and sequence alignments between NSP3 of SARS-CoV-2 or MHV and the updated NSP3 domain ranges of SARS-CoV-1 led to domain ranges of all three viruses. The ranges for SARS-CoV-2 are listed in Table 1, while the ranges for SARS-CoV-1 and MHV are listed in Table 17 of the appendix. Since no domain separation between Y1 and CoV-Y was provided for any of the three viruses, this domain border remained undefined. These preliminary ranges were then further analyzed and used as input sequences for structure prediction.

**Table 1:** *Preliminary residue ranges of NSP3 domains from SARS-CoV-2 in comparison to ranges defined by NCBI gene annotations with the reference id YP\_009742610.1. a: These ranges were not stated explicitly and are derived from the surrounding ranges. b: These ranges had no defined start/end. c: These ranges were predicted by TMHMM 2.0. d: This domain was not determined in the gene annotations. Domains colored in red are essential for viral replication in vivo.*

Complete Name	Preliminary residue ranges for SARS-CoV-2	Previously defined residue ranges for SARS-CoV-2
Ubiquitin-like domain 1 (Ubl1)	1-111	??? <sup>d</sup>
Hypervariable region (HVR)	112-206	??? <sup>d</sup>
Papain-like protease 1 (PL1 <sup>pro</sup> )	-	-
Macrodomain 1 (Mac1)	207-379	234-359
Linker Mac1-Mac2/Linker Mac1-DPUPlike	380-412	360-414 <sup>a</sup>
Macrodomain 2 (Mac2)	413-550	415-540
Linker Mac2-Mac3	-	-
Macrodomain 3 (Mac3)	551-675	533-675
Domain preceding Ubl2 and PL2 <sup>pro</sup> (DPUP)	676-745	678-743
Ubiquitin-like domain 2 (Ubl2)	746-804	748-??? <sup>b</sup>
Papain-like protease 2 (PL2 <sup>pro</sup> )	805-1063	???-1050 <sup>b</sup>
Linker PL2 <sup>pro</sup> -NAB	1064-1088	1051-1994 <sup>a</sup>
Nucleic-acidic-binding domain (NAB)	1089-1203	1095-1201
Betacoronavirus-specific marker domain ( $\beta$ SM)	1204-1412	1226-1341
Transmembrane domain 1 (TM1)	1413-1435 <sup>c</sup>	1414-1435
3Ecto	1436-1531 <sup>c</sup>	1436-1518 <sup>a</sup>
Transmembrane domain 2 (TM2)	1532-1554 <sup>c</sup>	1519-1541
Linker TM2-AH1	1555-1560 <sup>c</sup>	??? <sup>d</sup>
Amphipathic helix 1 (AH1)	1561-1583 <sup>c</sup>	??? <sup>d</sup>
Nidovirus-conserved domain of unknown function (Y1)	1584-??? <sup>b</sup>	??? <sup>d</sup>
Coronavirus-specific C-terminal domain (CoV-Ya)	??? <sup>b</sup> -1843	??? <sup>d</sup>
Coronavirus-specific C-terminal domain (CoV-Yb)	1844-1945	??? <sup>d</sup> -1944

### 3.1.2 Sequence alignments of preliminary domain ranges

The sequence similarity between domains from the three examined viruses based on the updated domain ranges was calculated with local pairwise sequence alignments [74]. The resulting sequence identities are listed in Table 2. In all cases, sequence identities between sarbecoviruses are highest, with values above 70 % for folding domains. From the 17 domains, the HVR, Linker-Mac1-Mac2, Linker-PL2<sup>pro</sup>-NAB, and  $\beta$ SM were considered to be intrinsically disordered [17]. During this work, however, the PDB structure 7T9W of the  $\beta$ SM domain [58] was published, which covers the central  $\sim$ 80 residues of this otherwise disordered domain. Therefore, this domain shows with 69.2 % a higher sequence identity than the disordered hyper variable region (43.8 % sequence identity) and the linker between Mac1 and Mac2 (41.4 % sequence identity). The other larger linker located between PL2<sup>pro</sup> and NAB, however, is with a sequence identity of 80 % more conserved than several folded domains. Later, it is shown that this linker is in fact folding to a stable protein structure. Two of the highest sequence identities are found at the C-terminus beyond the transmembrane region, with 88.1 % for Y1+CoV-Ya and 90.2 % for CoV-Yb, which together comprise 362 residues in SARS-CoV-2. Their role is also explored in more detail.

### 3.1.3 Comparison between AlphaFold2 predictions and experimentally determined structures

All sequences from the preliminary domain ranges were submitted to AlphaFold2 [15] via ColabFold [128] to predict structure models. Where possible, the predicted structures were aligned to their experimentally determined counterpart, where the root mean square deviations (RMSDs) between experimental structure and predicted model are listed in Table 3. The RMSD values are in all cases below 1 Å except for the Ubl1 domain from SARS-CoV-1 (RMSD of 1.3 Å) and MHV (RMSD of 2.7 Å). The major difference is in the disordered N-terminus [129, 130], which is predicted as loop with a low confidence (Figure 4a), reflected in low pLDDT (predicted local distance difference test) values. While secondary structure elements deviate slightly in orientation and position, the number of such elements was predicted always correctly for all cases in Table 3. Four of the structure superimpositions are shown in Figure 4.

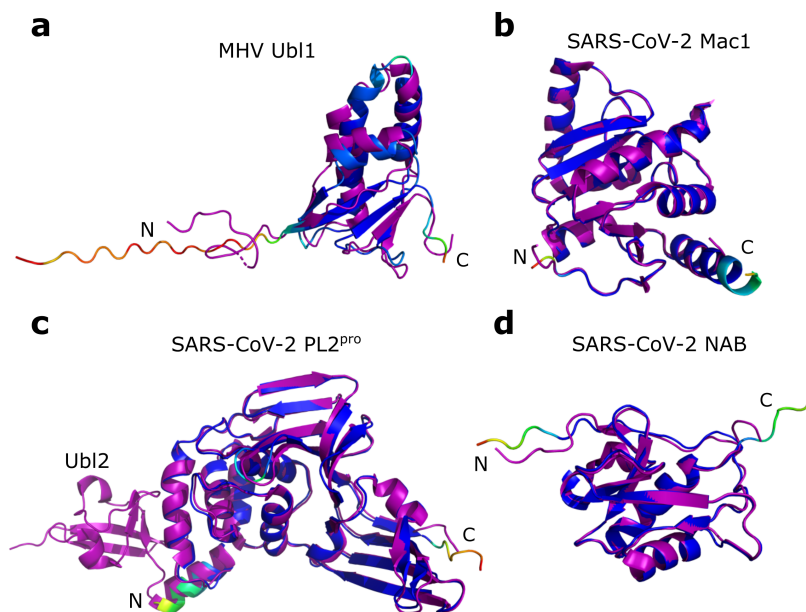
**Table 2:** Results from local pairwise sequence alignments with domain sequences from preliminary ranges. The alignments were calculated via EMBOSS Water version 6.6.0 [74]. Domains not present in a virus are marked with “-”. “x” indicates a failed alignment with less than 10 residues aligned. Domains shaded in blue are assumed to have a stable fold; red indicates linkers and domains assumed to be mostly disordered; yellow indicates domains in the transmembrane region; information about order and disorder according to Lei et al. [17].

Domain name	Sequence identity SARS-CoV-2 / SARS-CoV-1	Sequence identity SARS-CoV-2 / MHV	Sequence identity SARS-CoV-1 / MHV
Ubiquitin-like domain 1 (Ubl1)	78.60%	38.30%	35.40%
Hypervariable region (HVR)	43.80%	21.70%	25.60%
Macrodomain 1 (Mac1)	73.80%	32.50%	34.90%
Linker Mac1-Mac2	41.40%	38.50%	x
Macrodomain 2 (Mac2)	71.10%	-	-
Macrodomain 3 (Mac3)	81.60%	-	-
Domain preceding Ubl2 and PL2 <sup>pro</sup> (DPUP)	73.00%	25.00%	40.00%
Ubiquitin-like domain 2 (Ubl2)	90.00%	31.00%	33.00%
Papain-like protease 2 (PL2 <sup>pro</sup> )	81.10%	32.60%	35.10%
Linker PL2 <sup>pro</sup> -NAB	80.00%	40.00%	36.00%
Nucleic-acidic-binding domain (NAB)	81.70%	32.00%	29.00%
Betacoronavirus-specific marker domain ( $\beta$ SM)	69.20%	20.30%	22.80%
Transmembrane domain 1 (TM1)	77.00%	x	x
Nsp3 ectodomain	71.00%	32.00%	38.00%
Transmembrane domain 2 (TM2)	80.00%	x	x
Amphipathic helix 1 (AH1)	x	29.00%	29.00%
Nidovirus-conserved domain of unknown function (Y1) + Coronavirus-specific C-terminal domain a (CoV-Ya)	88.10%	37.70%	37.90%
Coronavirus-specific C-terminal domain b (CoV-Yb)	90.20%	40.60%	39.60%
full NSP3 sequence	76.00%	23.60%	23.40%

**Table 3:** Structural similarity based on alignments between AlphaFold2 prediction and corresponding experimental structure. Listed are domain names, abbreviation, and RMSD values for domains from SARS-CoV-2, SARS-CoV-1, and MHV, with PDB code of used experimental structure in parentheses after the RMSD value in Å. Alignments and RMSD were calculated in PyMOL [131]. Predicted structures are always rank1 models from AlphaFold2 [15], accessed via ColabFold [128]. For Mac1 and PL2<sup>pro</sup> from SARS-CoV-2 the ten experimentally determined structures with highest resolution were used, with the lowest RMSD listed in the table. For all other domains, all PDB entries were used, with the lowest RMSD listed here. Since 8F2E consists of multiple domains, unused atoms were removed in PyMOL prior alignment.

Name	Abbreviation	SARS-CoV-2	SARS-CoV-1	MHV
Ubiquitin-like domain 1	Ubl1	0.5 Å (7KAG)	1.3 Å (2GRI)	2.7 Å (2M0A)
Macrodomain 1	Mac1	0.3 Å (7KQP)	0.3 Å (2ACF)	-
Macrodomain 2	Mac2	-	0.4 Å (6YXJ)	-
Macrodomain 3	Mac3	-	0.7 Å (2JZD)	-
Domain preceding Ubl2 and PL2 <sup>pro</sup>	DPUP	0.3 Å (7THH)	0.6 Å (2KAF)	0.4 Å (4YPT)
Ubiquitin-like domain 2	Ubl2	0.2 Å (7D6H)	0.2 Å (2FE8)	0.3 Å (5WFI)
Papain-like protease 2	PL2 <sup>pro</sup>	0.7 Å (7D6H)	0.5 Å (5TL7)	0.5 Å (5WFI)
Nucleic-acidic-binding domain	NAB	0.4 Å (7LGO)	0.8 Å (2K87)	-
Betacoronavirus-specific marker domain	$\beta$ SM	0.9 Å (7T9W)	-	-
Nidovirus-conserved domain of unknown function subdomain b	Y1b	0.4 Å (8F2E)	-	-
Coronavirus-specific C-terminal domain	CoV-Ya	0.4 Å (8F2E)	-	-
Coronavirus-specific C-terminal domain	CoV-Yb	0.4 Å (7RQG)	-	-





**Figure 4:** Superimposition of experimentally determined PDB structures (purple) and AlphaFold2 prediction (colored according to pLDDT, with deep blue being highly confident and red non-confident). AlphaFold2 prediction are based on the preliminary domain ranges. a: Ubl1 domain from MHV, superimposed with PDB structure 2M0A. b: Mac1 domain from SARS-CoV-2, with PDB structure 7KQP. c: PL2<sup>pro</sup> domain from SARS-CoV-2, with PDB structure 7D6H. The Ubl2 domain is part of the PDB structure, but was excluded from prediction. d: NAB domain from SARS-CoV-2, with PDB structure 7LGO. 2M0A was solved via solution NMR, while the other three structures were solved via X-ray diffraction.

### 3.1.4 Classification into regions of order and disorder

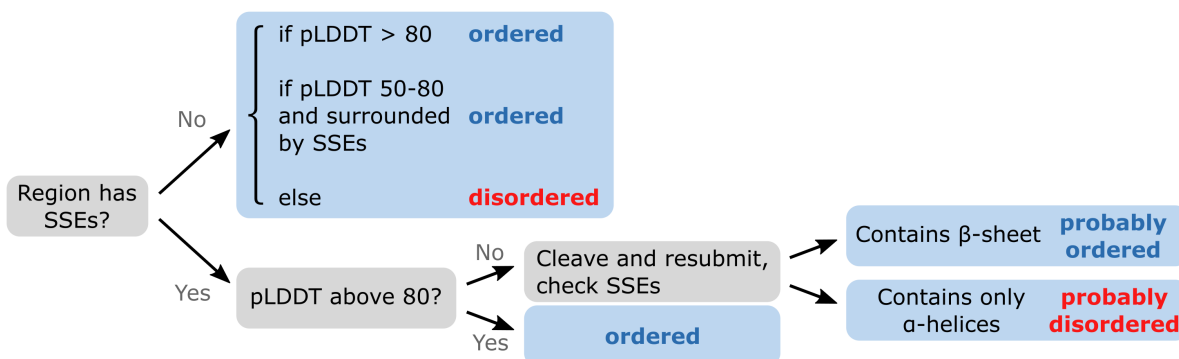
Since no experimentally determined structures are available for the remaining domains, the confidence metrics from AlphaFold2 are used to evaluate the plausibility of predicted models. These metrics are the per residue pLDDT (predicted local distance difference test) and the predicted aligned error (pAE), which assess the correctness of the distance between two residues from the sequence and is used to assess the relative spatial arrangement between each residue pair.

Since pLDDT values below 50 often indicate a region to be intrinsically disordered [20], the pLDDT could be utilized to distinguish between regions of order and disorder. However, low pLDDT values were also present in ordered regions. Therefore, information about secondary structure elements was utilized as well. Figure 5 illustrates which decisions were met when classifying a region into ordered or disordered. Loops and turns were fairly easy to categorize, since they are only ordered if they have high pLDDT values above 80 or are surrounded by secondary structure elements while having pLDDT values above 50. Otherwise they are disordered and tend to be predicted as pointing away from the main fold and into random directions, as described by Williams et al. as "barbed

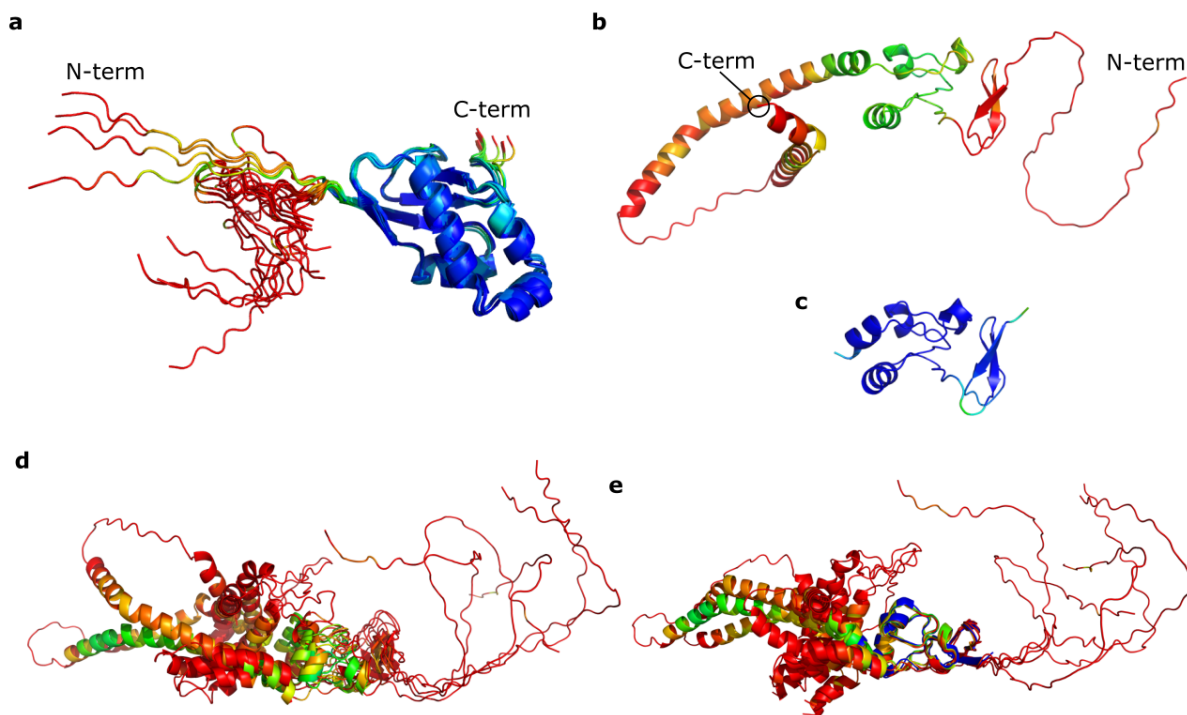
wire" conformations [132]. These conformations disagree with regular torsion angles and are depicted in Figure 6. If a region had secondary structure elements such as  $\alpha$ -helices or  $\beta$ -sheets, and pLDDT values above 80, it was considered ordered.

Figure 6a shows the Ubl1 domain of SARS-CoV-1, which consists of an ordered region with high pLDDT values (blue) and a low pLDDT region (red) at the 16 N-terminal residues. The superimposition of twenty predicted models emphasizes the differences between tiny deviations in the ordered region and the large deviations at both termini, where the N-terminal tail expands away from the fold into random directions.

The last case that is covered in the classification diagram (Figure 5) handles the prediction of secondary structure elements with low pLDDT values, which is best explained on the example of the Betacoronavirus-specific marker domain ( $\beta$ SM). This domain consists of a central folded domain (Figure 6c) surrounded by two large, disordered termini (Figure 6b). The prediction of the whole region comes with overall low pLDDT values and especially the termini, but also nearby secondary structure elements, have pLDDT values below 50 (colored red). Aligning the experimentally determined structure (PDB code 7T9W [58]) of the central section to the whole prediction in Figure 6b (alignment in Figure 6e), comes with an RMSD of 0.9 Å. Cleaving the low pLDDT termini beyond the well defined secondary structure elements, however, results in an increase of pLDDT values to 90 or higher for the central fold (Figure 6c) and RMSD values drop to 0.5 Å. Thereby, cropping certainly disordered regions can improve the prediction quality of nearby regions. Aligning predictions with ordered and disordered parts can help to identify a fold with the same orientation among all predictions (Figure 6d), but it is not clear where to draw the borders around the ordered fold from this.



**Figure 5:** Decision tree for classifying regions from fold-predictions into ordered or disordered. SSEs stands for “secondary structure elements”; pLDDT for “predicted local distance difference test”. Multiple iterations are only required if SSEs with low pLDDT are present, where certainly disordered regions must be cropped, which can increase prediction confidence.



**Figure 6:** Ordered and disordered regions in relation to confidence scores in predicted models. Residues are colored according to their pLDDT value, with blue representing high confidence with values above 90, red representing low confidence with values below 50, and other colors representing values in between. (a): Ensemble of twenty predicted models of SARS-CoV-1 Ubl1 emphasizing the difference between disorder (red, left) and order (blue, right). Predictions are made with AlphaFold2 by prediction five models per run for a total of four seeds. (b): Prediction of complete  $\beta$ SM domain from SARS-CoV-2 with low confidence secondary structure elements and large disordered termini. (c): Prediction of a cropped  $\beta$ SM domain sequence lacking the termini, showing an increased overall pLDDT. (d): Ensemble of five predicted models aligned to the rank1 model, with the central fold recognizable. (e): Same ensemble as in (d) aligned to the experimentally determined structure 7T9W colored in blue. Disordered regions, despite containing predicted  $\alpha$ -helices, point away in various directions.

### 3.1.5 Final determination of domain boundaries

The described method of classifying segments of NSP3 into ordered or disordered regions based on models and pLDDT values from AlphaFold2 predictions was applied on the sequences from the preliminary domain ranges listed in Table 1. Regions of order were then always containing one or more domains, where Table 4 and Table 5 list these regions and the final domain ranges. In total, 357 residues of SARS-CoV-2 NSP3 (18.4 % of NSP3) were classified as disordered. For SARS-CoV-1 and MHV, 344 (17.9 %) and 506 residues (25.2 %) were classified as disordered, respectively.

The last step was to utilize the predicted aligned error (pAE) matrix to identify sub-domains or multiple domains within one ordered region. In NSP3, a large ordered region ranged from Mac3 to NAB (Table 4). However, since the ranges based on experimentally determined structures were already available, the region was simply split into domains

according to the preliminary domain ranges.

Noteworthy insights gained from the classification are the division of the Betacoronavirus-specific marker domain ( $\beta$ SM) into two disordered regions surrounding one ordered fold (Table 4) and the identification of the domain ranges for the C-terminal domains. Towards the end of this work, both findings were validated by independent research through published experimental structures. The PDB structure 7T9W [58] resembles the central fold of  $\beta$ SM, which was named  $\beta$ SM-M (Table 4). The structures with the PDB codes 7RQG [59] and 8F2E [60] validate the ranges of the C-terminal subdomains Y1b, CoV-Ya, and CoV-Yb, which are further explored in section 3.3.

Additional discoveries, which have not been experimentally observed before are a predicted  $\alpha$ -helix in the hyper variable region for MHV (NSP3 residues 230-241); a folding linker between PL2<sup>PRO</sup> and NAB, which is explored in section 3.2; and the predicted structure of the ectodomain.

### 3.1.6 Structural similarity of final domains

Since the new ranges listed in Table 4 and Table 5 comprise complete regions of order or disorder, new local pairwise sequence alignments were conducted for each domain and with sequences from all three examined viruses. Furthermore, the structure predictions were aligned and their RMSD was calculated. Both results, for sequence and structure similarity, are listed in Table 6. The results for alignments between SARS-CoV-1 and MHV are found in Table 18 of the appendix.

The ordered NSP3 domains of the sarbecoviruses are highly similar in sequence and structure with sequence similarities above 77 % and with RMSD values of 0.7 Å or lower (excluding transmembrane domains). The C-terminal domains Y1a, Y1b, CoV-Ya, and CoV-Yb stand out with sequence similarity from 92 % to 100 %, despite consisting together of more than 340 residues. Potential functions of these domains are explored in section 3.3.

Between SARS-CoV-2 and MHV, sequence similarities for ordered domains are above 42 % and RMSD values range from 0.2 Å up to 7.4 Å, where only four domains have RMSD values above 1.1 Å. The exceptional cases are CoV-Y (RMSD of 2.1 Å), DPUP (2.9 Å RMSD),  $\beta$ SM-M (3.4 Å RMSD) and the folded core of the ectodomain (7.2 Å RMSD).

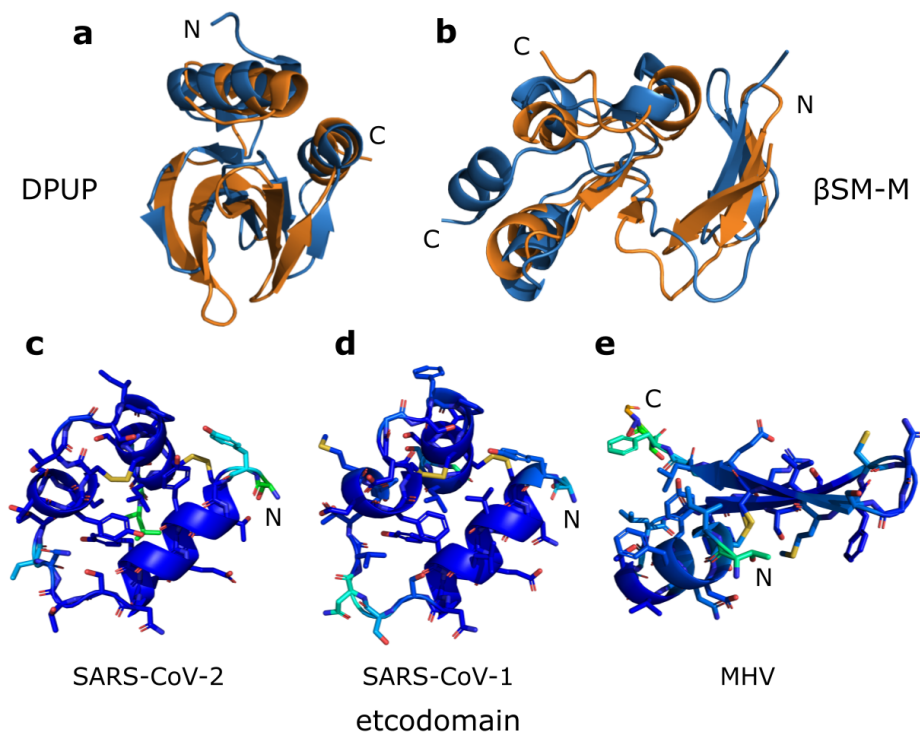
The high RMSD of the CoV-Y domain is due to an alternative arrangement of the subdomains as both subdomains individually have RMSD values of 0.7 Å or lower. A noticeable difference for DPUP are longer  $\beta$ -strands in MHV, which leads to two non-

**Table 4:** Ranges of all domains and linkers prior the transmembrane region. Ranges are determined by classifying NSP3 segments into ordered or disordered via AlphaFold2. Entries shaded in blue are classified as ordered, while non-shaded entries describe regions predicted to be disordered. Domains discussed in more detail are shaded in orange.

Complete Name	Abbr.	SARS-CoV-2	SARS-CoV-1	MHV
N-terminal loop of Ubl1	Ubl1-N	1-16	1-17	1-16
Ubiquitin-like domain 1	Ubl1	17-111	18-107	17-113
Hypervariable region	HVR	112-208	108-186	114-272
Papain-like protease 1	PL1 <sup>pro</sup>	-	-	273-476
Linker PL1 <sup>pro</sup> -Mac1		-	-	477-487
Macrodomain 1	Mac1	209-377	187-355	488-644
Linker Mac1-Mac2 (Sarbecovirus) / Linker Mac1-DPUP-like (MHV)		378-412	356-390	645-666
Linker Mac1-DPUP-like helix		-	-	667-680
Linker pre DPUP		-	-	681-703
Macrodomain 2	Mac2	413-540	391-517	-
Linker Mac2-Mac3		541-550	518-526	-
Macrodomain 3	Mac3	551-675	527-651	-
Domain preceding Ubl2 and PL2 <sup>pro</sup> / DPUP-like domain (MHV)	DPUP	676-745	652-722	704-777
Ubiquitin-like domain 2	Ubl2	746-804	723-781	778-837
Papain-like protease 2	PL2 <sup>pro</sup>	805-1056	782-1036	838-1084
betacoronavirus-specific linker domain	$\beta$ SLD	1057-1090	1037-1067	1085-1116
Linker $\beta$ SLD-NAB		-	-	1117-1135
Nucleic-acidic-binding domain	NAB	1091-1196	1068-1174	1136-1211
Betacoronavirus-specific marker domain N-terminal subdomain	$\beta$ SM-N	1197-1239	1175-1217	1212-1292
Betacoronavirus-specific marker domain Folded Core	$\beta$ SM-M	1240-1325	1218-1304	1293-1369
Betacoronavirus-specific marker domain C-terminal subdomain	$\beta$ SM-C	1326-1412	1305-1390	1370-1448

**Table 5:** Ranges of all domains and linkers starting from the transmembrane region. Ranges are determined by classifying NSP3 segments into ordered or disordered via AlphaFold2. Entries shaded in blue are classified as ordered, while non-shaded entries describe regions predicted to be disordered. Domains discussed in more detail are shaded in orange.

Complete Name	Abbr.	SARS-CoV-2	SARS-CoV-1	MHV
Transmembrane domain 1	TM1	1413-1435	1391-1413	1449-1471
Linker TM1-Ecto		1436-1442	1414-1419	1472-1504
Ectodomain core fold	EctoCore	1443-1477	1420-1453	1505-1534
Ectodomain linker / Linker Ecto-TM2 (MHV)	EctoL	1478-1499	1454-1475	1535-1564
Ectodomain TM-like helix	EctoTM	1500-1522	1476-1492	-
Linker EctoTM-TM2		1523-1531	1493-1495	-
Transmembrane domain 2	TM2	1532-1554	1496-1518	1565-1587
Linker TM2-AH1		1555-1560	1519-1522	1588-1607
Amphipathic helix 1	AH1	1561-1583	1523-1545	1608-1630
Linker AH1-Y1		1584-1598	1546-1575	1631-1660
Nidovirus-conserved domain of unknown function	Y1	1599-1759	1576-1736	1661-1819
Y1-subdomain a	Y1a	1599-1664	1576-1646	1661-1731
Y1-subdomain b	Y1b	1665-1759	1647-1736	1732-1820
Linker Y1-CoV-Y		1760-1765	1737-1742	1821-1824
Coronavirus-specific C-terminal domain	CoV-Y	1766-1945	1743-1922	1825-2006
CoV-Y-subdomain a	CoV-Ya	1766-1847	1746-1824	1825-1908
CoV-Y-subdomain b	CoV-Yb	1848-1945	1825-1922	1909-2006



**Figure 7:** Comparison of domains with high RMSD values. *a:* DPUP from SARS-CoV-2 (PDB structure 7THH) in blue and MHV (PDB structure 4YPT) in orange. *b:*  $\beta$ -SM-M from AlphaFold2 predictions for SARS-CoV-2 (blue) and MHV (orange). *c-e:* AlphaFold2 predictions of ectodomain for SARS-CoV-2 (c), SARS-CoV-1 (d), and MHV (e), all with backbone colored according to pLDDT with deep blue for pLDDT values above 90. Non-carbon atoms are colored by element, with oxygen in red, nitrogen in blue, and sulfur in yellow.

overlapping regions, while the remaining arrangement of secondary structure elements looks similar. Comparing the experimentally solved equivalents also leads to a high RMSD of 2.4 Å between the PDB structures 7THH and 4YPT, for SARS-CoV-2 and MHV respectively (Figure 7a). The  $\beta$ SM-M domain is a similar case, where the fold and composition of secondary structure elements is similar, but orientation and length of individual elements differ (Figure 7b). Since an experimentally determined structure of this domains exists only for SARS-CoV-2, no validation of this prediction was possible.

With an RMSD of 7.2 Å the MHV ectodomain is the strongest outlier. Despite a sequence similarity of 50 % and pLDDT values above 90, it is predicted with a completely different fold in MHV compared to the sarbecoviruses (see Figure 7c-e). While the sarbecovirus fold is predicted with three short  $\alpha$ -helices and two disulfide bonds, the MHV fold prediction comes with only one  $\alpha$ -helix and with two  $\beta$ -strands. Furthermore, the four cysteines are present in all three viruses, but in MHV only one pair is forming a disulfide bond, while the other pair is out of reach to form such bond.



**Table 6:** Sequence similarities and RMSD values between predicted folds of domains from SARS-CoV-2 and SARS-CoV-1 or MHV. Only domains predicted to fold into a defined structure and large regions of disorder are listed. RMSD values are calculated with PyMOL [131] for folded domains. Results are sorted by decreasing sequence similarity between both sarbecoviruses. Domains of the transmembrane region are listed below and are not sorted, since only short alignments were found. For these cases, the alignment-length is given in parentheses. Sequence similarity was used over sequence identity due to also listing disordered domains, which show high similarity and low identities. Domains marked with an asterisk consist of subdomains, which are also listed individually.

Domain	SARS-CoV-2 to SARS-CoV-1		SARS-CoV-2 to MHV	
	Sequence similarity	RMSD	Sequence similarity	RMSD
CoV-Yb	100%	0.1 Å	58%	0.4 Å
Y1b	98%	0.1 Å	61%	0.6 Å
Y1*	97.5%	0.1 Å	65.8%	0.7 Å
Y1a	97%	0.1 Å	73%	0.3 Å
Ubl2	97%	0.1 Å	54%	0.6 Å
CoV-Y*	96.5%	0.3 Å	56.5%	2.1 Å
CoV-Ya	92%	0.2 Å	58%	0.7 Å
$\beta$ SLD	90%	0.3 Å	54%	0.2 Å
Ubl1	90%	0.5 Å	53%	0.4 Å
Mac3	89.6%	0.2 Å	-	-
PL2 <sup>pro</sup>	89.3%	0.3 Å	47.1%	1.1 Å
NAB	88.7%	0.3 Å	45.6%	1.0 Å
$\beta$ SM-M	88%	0.7 Å	42%	3.4 Å
DPUP	87%	0.3 Å	50%	2.9 Å
$\beta$ SM-N	86%	-	53%	-
Mac1	85.5%	0.2 Å	52.6%	0.7 Å
Mac2	84.9%	0.4 Å	-	-
3Ecto core	77%	0.5 Å	50%	7.2 Å
$\beta$ SM-C	75%	-	63%	-
HVR	62%	-	35%	-
TM1	82% (22)	0.3 Å	67% (6)	0.2 Å
EctoL + EctoTM	87% (38)	-	62% (13)	-
EctoTM	75% (16)	0.1 Å	-	-
TM2	80% (10)	1.7 Å	83% (6)	0.7 Å
AH1	100% (8)	0.5 Å	48% (21)	0.5 Å



**Table 7:** Results from structure similarity search on the PDB [124]. For NAB and  $\beta$ SM-M the PDB structures 7LGO and 7T9W were used as templates, respectively. For the remaining domains AlphaFold2 predictions were used. Best hits are the next highest scoring match which are not other assemblies deposited under the same PDB code. RMSD values were calculated in PyMOL [131] after structure alignment via cealign [133].

Domain	Best hit	Structure match score strict	Structure match score relaxed	RMSD
$\beta$ SLD	2K6R	no hits	37.60	4.7 Å
NAB	1QW1	8.99	43.90	3.2 Å
$\beta$ SM-M	1L2N	no hits	37.45	7.6 Å
EctoCore	2E2F	0.07	38.40	5.1 Å
Y1	6A6I	7.51	42.99	7.9 Å

### 3.1.7 Structure similarity search

For the domains  $\beta$ SLD, 3Ecto, and Y1a, no structures were available in the PDB. Furthermore, the domains NAB and  $\beta$ SM are of great interest as research regarding their functions is scarce. Therefore, a structure similarity search was performed on the whole PDB [124] in order to identify similar folds which result from vastly different sequences and are thus not identifiable via a BLAST search [76].

The results are listed in Table 7. While structures with a roughly similar shape were identified as indicated by high structure match scores in the relaxed search, the best hits had very different folds and composition of secondary structure elements, which is reflected in high RMSD values and low scores or no hits in the strict structure similarity search.

## 3.2 Experimental validation of the Betacoronavirus-specific linker domain ( $\beta$ SLD)

The Betacoronavirus-specific linker domain describes the segment of NSP3 which is preceded by PL2<sup>pro</sup> and followed by NAB. It was previously not described to exhibit a fold, nor was it ever stated as a unique domain or part of the nearby domains. It was only indirectly shown as a linker by defining the ranges of its surrounding domains [17]. The following results validate its folding nature and show its specificity to *Betacoronavirus*.

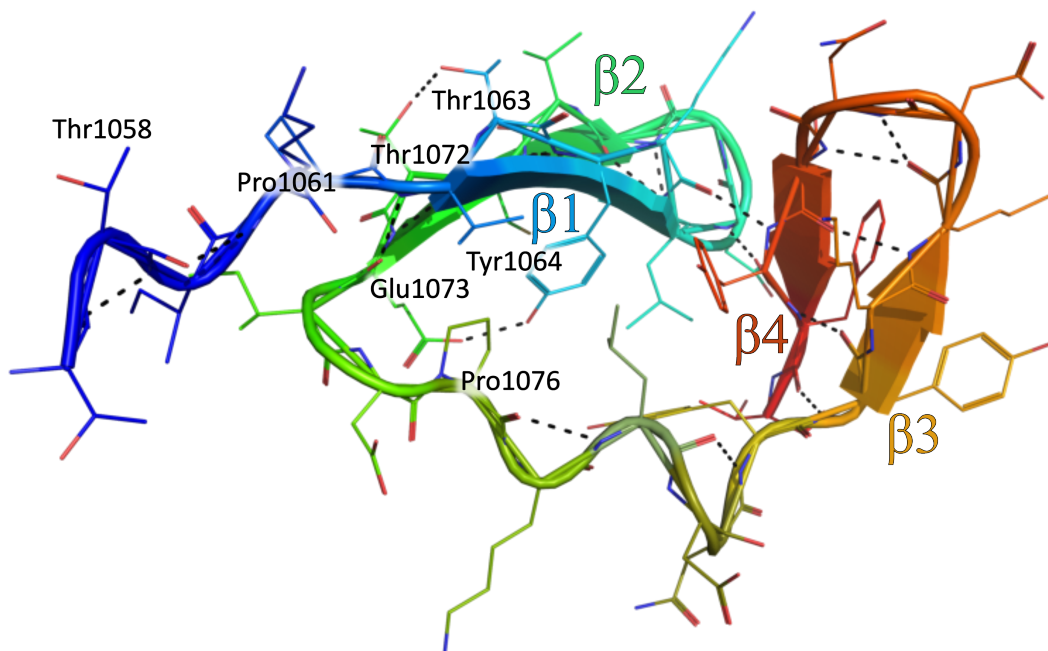
### 3.2.1 Structure prediction of linker domain between PL2<sup>pro</sup> and NAB domain

The AlphaFold2 structure prediction of the SARS-CoV-2 linker domain is illustrated in Figure 8. The fold consists of 34 residues, forming two loops and four  $\beta$ -strands in the sequence Nterm- $\beta$ 1- $\beta$ 2-loop- $\beta$ 3- $\beta$ 4-Cterm. The strands  $\beta$ 1 and  $\beta$ 2 are interconnected and

form an anti-parallel  $\beta$ -sheet. The same situation is true for  $\beta 3$  and  $\beta 4$ . Furthermore,  $\beta 1$  and  $\beta 4$  are connected via two hydrogen bonds in parallel direction, leading to an overall compact fold. Additional connections between  $\beta 1$  and  $\beta 2$  are provided by the hydrogen bonds of Thr1063-Thr1072 and Tyr1064-Glu1073.

Between  $\beta 2$  and  $\beta 3$  is the 8-residue long central loop, which exhibits a helix-like structure element containing two hydrogen bonds with itself. The residue Pro1076 could assist in the formation of this structure element by limiting the freedom in conformational space. A similar structure element containing the residue Pro1061 is found at the N-terminal loop, which connects the main fold to PL2<sup>PRO</sup>.

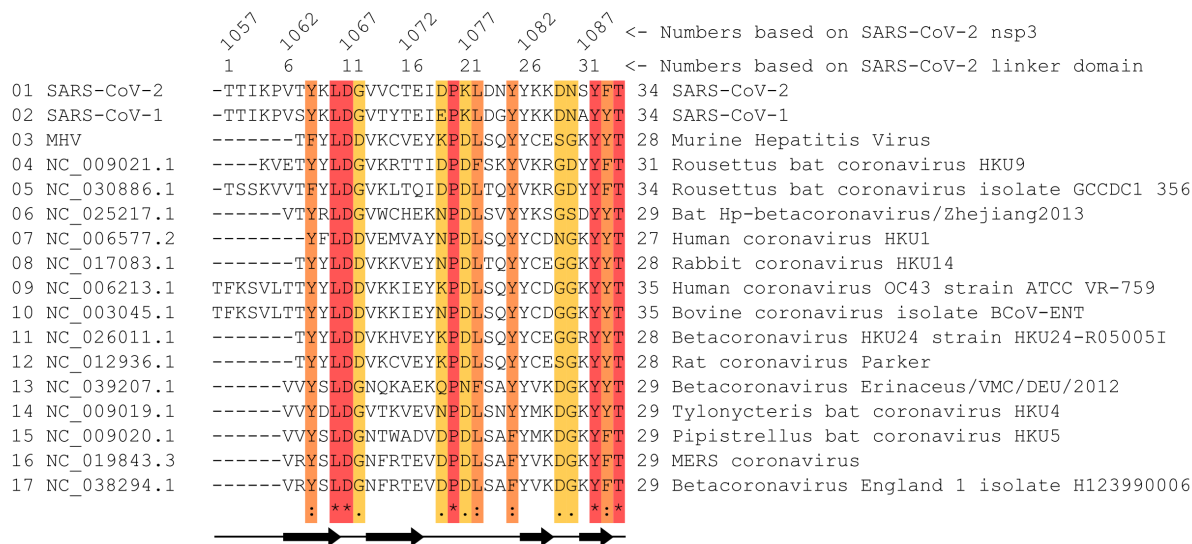
From the 34 residues, 26 are in favoured regions of the Ramachandran plot (see Figure 43 in the appendix). Only the residue Thr1058 (Figure 8) is recognized as a Ramachandran outlier.



**Figure 8:** *AlphaFold2* structure prediction of the SARS-CoV-2 NSP3 linker domain. N-terminus is coloured in blue, C-terminus in red. Hydrogen bonds are depicted as black dotted lines. The labelled proline residues potentially restrict the conformational freedom of the loops in which they are located, with Pro1076 being located in the central loop. Thr1058 is the only Ramachandran outlier. The other labelled residues interact via hydrogen bonds with each other. Residue numbers are based on SARS-CoV-2 NSP3 sequence.

### 3.2.2 Conservation of the linker domain

The sequence of the SARS-CoV-2 linker domain shares 80 % sequence identity with SARS-CoV-1 and 40 % with MHV (Table 2). Regarding sequence similarity, these values raise to 90 % and 54 % respectively. While similar sequences were found in other betacoronaviruses

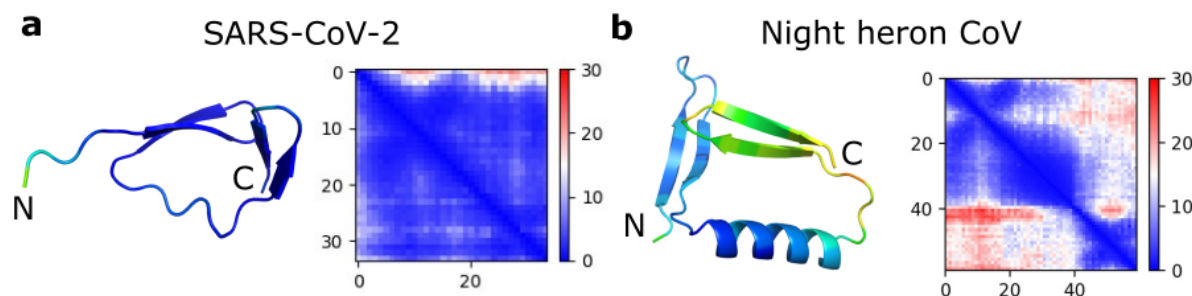


**Figure 9:** Multiple sequence alignment between the linker domain of 17 betacoronaviruses. The sequences for the linker domains were first identified with a local pairwise sequence alignment between the linker domain sequence of SARS-CoV-2 and the *orf1ab* sequence of the respective virus. Afterwards, the identified sequences were used in a global multiple sequence alignment, resulting in this figure. Residues conserved in all examined viruses are highlighted with red (marked with '\*'), residues with strongly similar chemical properties at the same position are highlighted with orange (marked with ':'), and weakly similar properties are highlighted with yellow (marked with '.'). Secondary structure elements are indicated at the bottom.

in a BLAST [76] search, no hits were found outside of this genus.

To identify conserved residues, multiple sequence alignments between the linker domain sequence from SARS-CoV-2 and the respective region from 16 additional betacoronaviruses were performed (Figure 9), which identified the residues Leu1066, Asp1067, Pro1076, Tyr1088, and Thr1090 as conserved among this group (residue numbers based on SARS-CoV-2 NSP3). Residues conserved among the majority of this group are the residue Leu1078, which is not leucine in only two cases, and the residues 1064, 1081, and 1089, which are in all cases either phenylalanine or tyrosine. In total, 14 out of 34 residues (including the fully conserved ones) retain similar chemical properties, of which 12 are located in a loop or at the edge of a  $\beta$ -strand adjacent to a loop residue. Comparison of structure predictions from all 16 viruses to the prediction from SARS-CoV-2 show RMSD values below 0.2 Å in all cases. The sequence identity and similarity range from 37 % to 55 % and from 46 % to 69 %, respectively (excluding SARS-CoV-1), and are listed in Table 19 of the appendix.

Additional sequence alignments with the unclassified shrew coronavirus, alpha- (24 viruses), gamma- (5 viruses), and deltacoronaviruses (10 viruses), revealed no hits located on NSP3. Therefore, the domain was named "Betacoronavirus-specific linker domain".



**Figure 10:** Comparison of the predicted Betacoronavirus-specific linker domain from SARS-CoV-2 (a) to the prediction of a similar sequence from Night heron coronavirus HKU19 (b), which is located in the endoRNase of this deltacoronavirus. Both structures are colored according to their pLDDT values, with deep blue representing values above 90 and red for values below 50, where other colors are interpolating for values in between. Next to the structure, the respective pAE matrix is shown, where lower values mean higher confidence of the spatial arrangement of two residues. The structure from Night heron coronavirus contains an additional  $\alpha$ -helix.

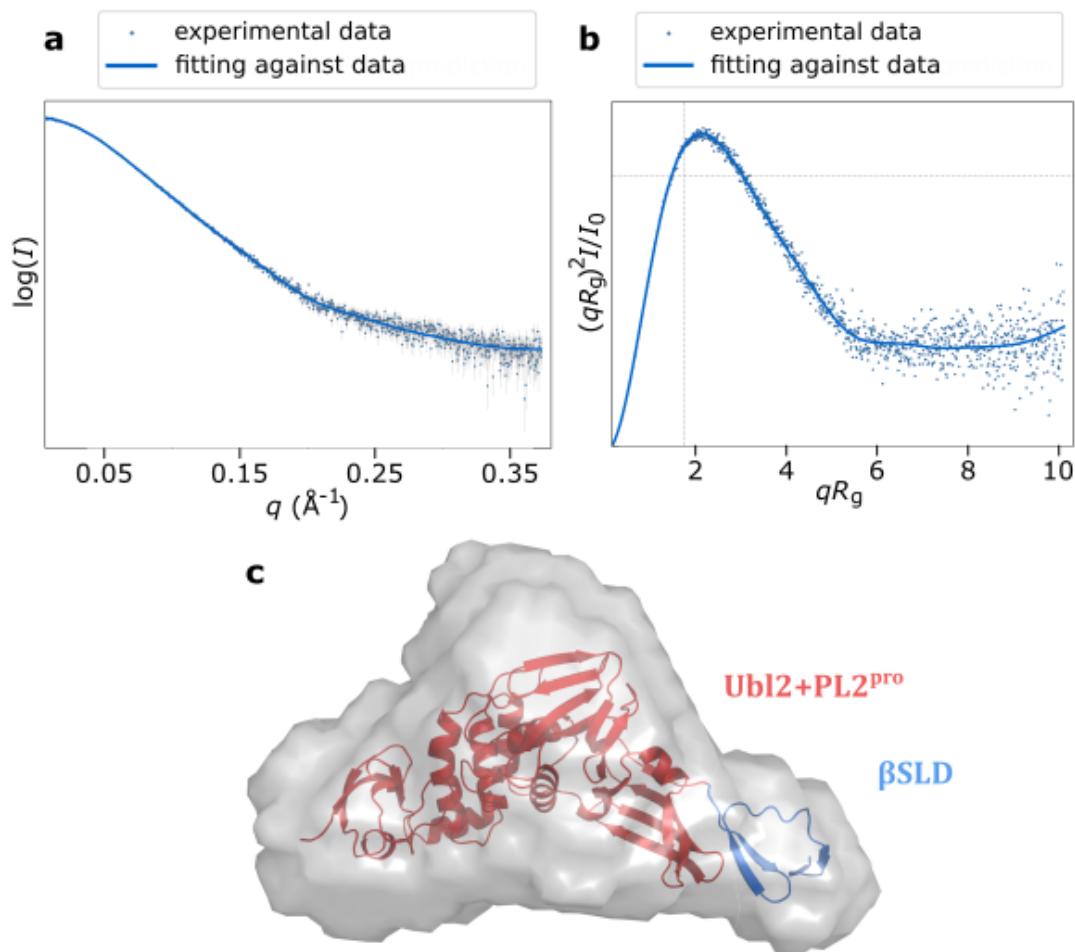
However, numerous hits were located in the endoRNase (NSP15 in SARS-CoV-2), which was the case for fifteen alpha-, three gamma-, and all deltacoronaviruses. In *Deltacoronavirus*, the linker domain sequence aligns in all cases with the region that covers Thr275 to Lys345 in the SARS-CoV-2 endoRNase structure 6VWW. The alignment with the Night heron coronavirus HKU19 (NCBI accession number NC\_016994) is shown below:

SARS-CoV-2_CL	14	VCTEIDPKLDNY-----YKKDNSYFT	34
		:..: ...:	:.....:
NC_016994	5943	VCSVVDLTIDNYIDIIRQAHSVYETKSKVFT	5973

An AlphaFold2 prediction of this region does not look like the respective region of the endoRNase, but consists of the same secondary structure elements as the linker domain with an additional  $\alpha$ -helix inserted between Tyr1081 and Tyr1082 (the gap region in the alignment). This helix is located between the central loop and the third  $\beta$ -strand, which prevents hydrogen bonds between  $\beta_4$  and  $\beta_1$ . This composition of secondary structure elements prevents the formation of a compact fold. This is also reflected in lowered pLDDT values and in the pAE matrix, which suggests two subdomains that are arranged with low confidence to each other (Figure 10).

### 3.2.3 Experimental validation

To assess if the Betacoronavirus-specific linker domain ( $\beta$ SLD) is in fact folded, single-crystal X-ray diffraction and solution small-angle X-ray scattering (SAXS) experiments were conducted. Due to the small size of the linker domain and the high confidence prediction of a close arrangement between this domain and PL2<sup>PRO</sup>, a multidomain construct was designed, which comprised the SARS-CoV-2 residues 746 to 1090 and the domains



**Figure 11:** Results from SAXS experiment. (a): The fitting result of relaxed AlphaFold2 model against the experimental SAXS data. Dots are the SAXS data with the relative errors. Solid line is the estimated scattering curve of the relaxed AlphaFold2 prediction. (b): The dimensionless Kratky plot. The peak maximum largely shifts away from the theoretical value for a compacted globular (marked by the two dashed grey lines), suggesting the solution structure of Ubl2-PL2pro- $\beta$ SLD is a multidomain protein. convergence of the Kratky plot at higher  $q$  suggests the absence of long flexible linkers or termini. (c): A projection of the envelope of the *ab initio* model (grey volume) and the relaxed AlphaFold2 prediction (cartoon), with the linker domain at the right side in blue and the remaining part being Ubl2+PL2pro in red. (Figure caption was primarily written by Yunyun Gao for a publication related to this work. Plots and structure model image are by Yunyun Gao.)

Ubl2, PL2<sup>pro</sup>, and  $\beta$ SLD. For better crystallization chances the common C111S mutant [25] was used.

Crystallization trials led to thin crystals, from which data could be collected and processed into an electron density map. Unfortunately, the model building revealed that cleavage took place and only half of the construct assembled into the crystal. It was the undesired half covering Ubl2 and approximately half of PL2<sup>pro</sup>. Despite the surprising result that half of PL2<sup>pro</sup> crystallizes, no useful information about the structure of  $\beta$ SLD could be gathered. It is unclear whether PL2<sup>pro</sup> was involved in the proteolysis or other factors introduced a systematic cleavage. However, it is unlikely that PL2<sup>pro</sup> is able to self-cleave *in vivo* due to its sequence-specificity [17] and the assembly of NSP3 into large hexameric complexes [19, 32], as well as no reports about such behaviour. Analysis of the construct prior crystallization showed a band at the expected size. Also, remaining protein sample was used in a SAXS experiment, where cleavage of the construct could be excluded.

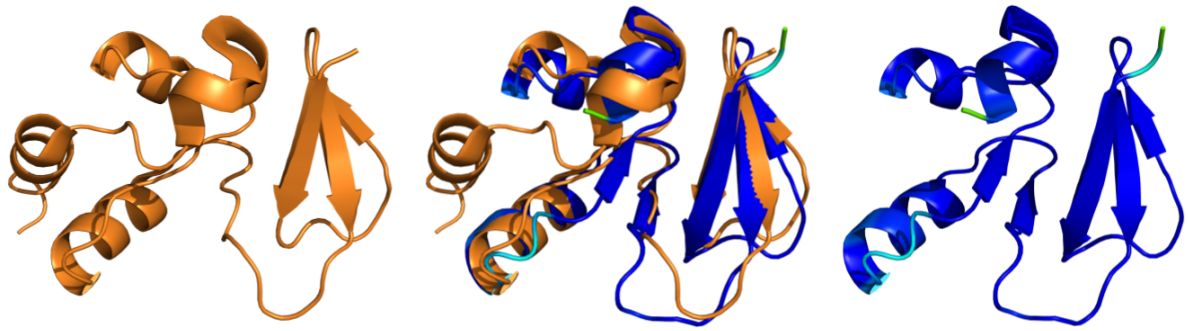
*(The following paragraph was primarily written by Yunyun Gao for a publication related to this work.)*

The SAXS result shows a good agreement between the relaxed AlphaFold2 prediction and the experimental data, with a  $\chi^2$  of 0.98 (Figure 11a) (The relaxed AlphaFold2 prediction is the most fitted structure found by the end of a single SREFLEX run [134]). The dimensionless Kratky plot suggests that the solution structure of Ubl2-PL2<sup>pro</sup>- $\beta$ SLD is a rather compacted multidomain entity (Figure 11b). It is highly unlikely that a long disordered tail exists, according to the flat plateau at the high scattering vector in the Kratky plot (Figure 11b); instead, short flexible regions between domains are expected. The relaxed AlphaFold2 prediction also fits well into the envelop of the *ab initio* SAXS model (Figure 11c).

*(From here on, texts are written again by Maximilian Edich.)*

### 3.2.4 Nucleic acid binding domain and Betacoronavirus-specific marker domain

The two domains following the Betacoronavirus-specific linker domain, namely the nucleic acid binding domain (NAB) and the Betacoronavirus-specific marker domain ( $\beta$ SM), are also specific to *Betacoronavirus* [17, 18]. Sequence alignments between NAB and  $\beta$ SM from SARS-CoV-2 and NSP3 from 15 other betacoronaviruses (listed in Table 19 in the



**Figure 12:** Comparison of Betacoronavirus-specific marker domain from SARS-CoV-2 (left, PDB structure 7T9W) with AlphaFold2 prediction of Gammacoronavirus-specific marker domain from Canada goose coronavirus (right, colored according to pLDDT with deep blue for values above 90). Both, experimental structure and high-confidence prediction, are superimposed in the central image with an RMSD of 1.5 Å.

appendix), excluding SARS-CoV-1, show sequence identities of 27 % to 38 % for NAB and identities of 18 % to 33 % for  $\beta$ SM-M. Structure predictions lead to similar folds, overall indicating that both domains are present in all betacoronaviruses. In *Gammacoronavirus*, a region similar to  $\beta$ SM is present, the Gammacoronavirus-specific marker domain ( $\gamma$ SM) [17]. However, past research compared the whole  $\beta$ SM and  $\gamma$ SM domains [17, 18, 57], which consist both of large proportions of intrinsic disorder. For the new domain ranges listed in Table 4,  $\beta$ SM was divided into the disordered subdomains  $\beta$ SM-N and  $\beta$ SM-C and the folding subdomain  $\beta$ SM-M. Both, NAB and  $\beta$ SM-M, were compared to gammacoronaviruses, where the highest sequence identities were 22 % for NAB and 31 % for  $\beta$ SM-M, both with sequences from the Canada goose coronavirus. Structure prediction of the  $\gamma$ SM-M from this virus resembles closely that of  $\beta$ SM-M from SARS-CoV-2 with an RMSD of 1.5 Å and comes with overall high pLDDT values (Figure 12). However, no similar results were found for other gammacoronaviruses. For NAB, no similar folds were predicted from *Gammacoronavirus*.

An interesting observation is the high sequence similarity of 56 %, when aligning the sequences of NAB and  $\beta$ SM-M from SARS-Cov-2 with an alignment length of 25 residues. While no structural similarity between NAB and  $\beta$ SM-M is present, this high sequence similarity could hint at a common origin via gene duplication, as both domains emerged during the recent evolution of *Betacoronavirus*.

### 3.3 Hexameric pore complex

The interaction between NSP3 and NSP4 induces the formation of double membrane vesicles (DMVs), which serve as a safe encapsulated environment for replication of viral RNA [27]. Furthermore, a hexameric pore complex is assembled at the surface of the



DMVs, which exports the replicated RNA into the cytosol and consists of NSP3, NSP4, and NSP6 [19, 32].

During the time of this work, a density volume at 30.5 Å resolution obtained from cryo electron tomography [19] was the only structural data available of this complex and to this date, no arrangement of the NSP3 domains within this complex was experimentally determined except for the Ubl1 domain being located at most distant part from the DMV surface [19]. Hence, the next steps of this thesis describe how models from structure prediction, information about ordered and disordered domains, and logical reasoning can be combined with available experimental data to state plausible hypotheses about the arrangement of few of NSP3's domains. After the practical part of this work, a second paper regarding the pore complex was published together with a cryo tomography map of the SARS-CoV-2 pore complex at 20.3 Å resolution and a manual placement of certain NSP3 domains into the map [32]. Because the study was published recently and the integrative modelling required a lot of time and computing power, most of the following steps were performed only on the map of the MHV pore complex.

Key to the presented assembly hypothesis is the prediction of the hexameric assembly of Y1, which is known as “nidovirus-conserved domain of unknown function” [18]. The structure and conservation of the C-terminal domains Y1 and CoV-Y are analysed here in regard to potential functions and the complex assembly. Additionally, a nomenclature conflict between two recently published structures of Y1b and CoV-Y is resolved.

### 3.3.1 Predicted arrangement of domains

Sequences covering multiple adjacent domains were submitted to ColabFold [128] to analyse the confidence of their predicted arrangement to each other, which can be used to identify quaternary domain assemblies [103]. A full length prediction of NSP3 was performed as well. Figure 13 shows the predicted aligned error (pAE) matrices for the predictions of the multidomain segments. High confidence values were only observed for the arrangement of PL2<sup>pro</sup> with Ubl2 and  $\beta$ SLD, between the transmembrane domains, and for the subdomains of Y1 and CoV-Y (covered in the next section). The arrangement of PL2<sup>pro</sup> with Ubl2 and  $\beta$ SLD was examined in the SAXS experiment in section 3.2.3, where the data indicates some flexibility between the domains. The arrangement of the transmembrane domains cannot be considered realistic, as the prediction was made in absence of a membrane. The high confidence may come from the domains being hydrophobic and arranging those close together improves the score.

Moderate confidence was observed for the arrangement of the three domains Mac2,



Mac3, and DPUP, and for the arrangement of Y1 and CoV-Y (Figure 13). However, only the latter arrangement was consistently predicted with moderate to high confidence, while Mac2 to DPUP is often predicted with high pAE values and thus low confidence.

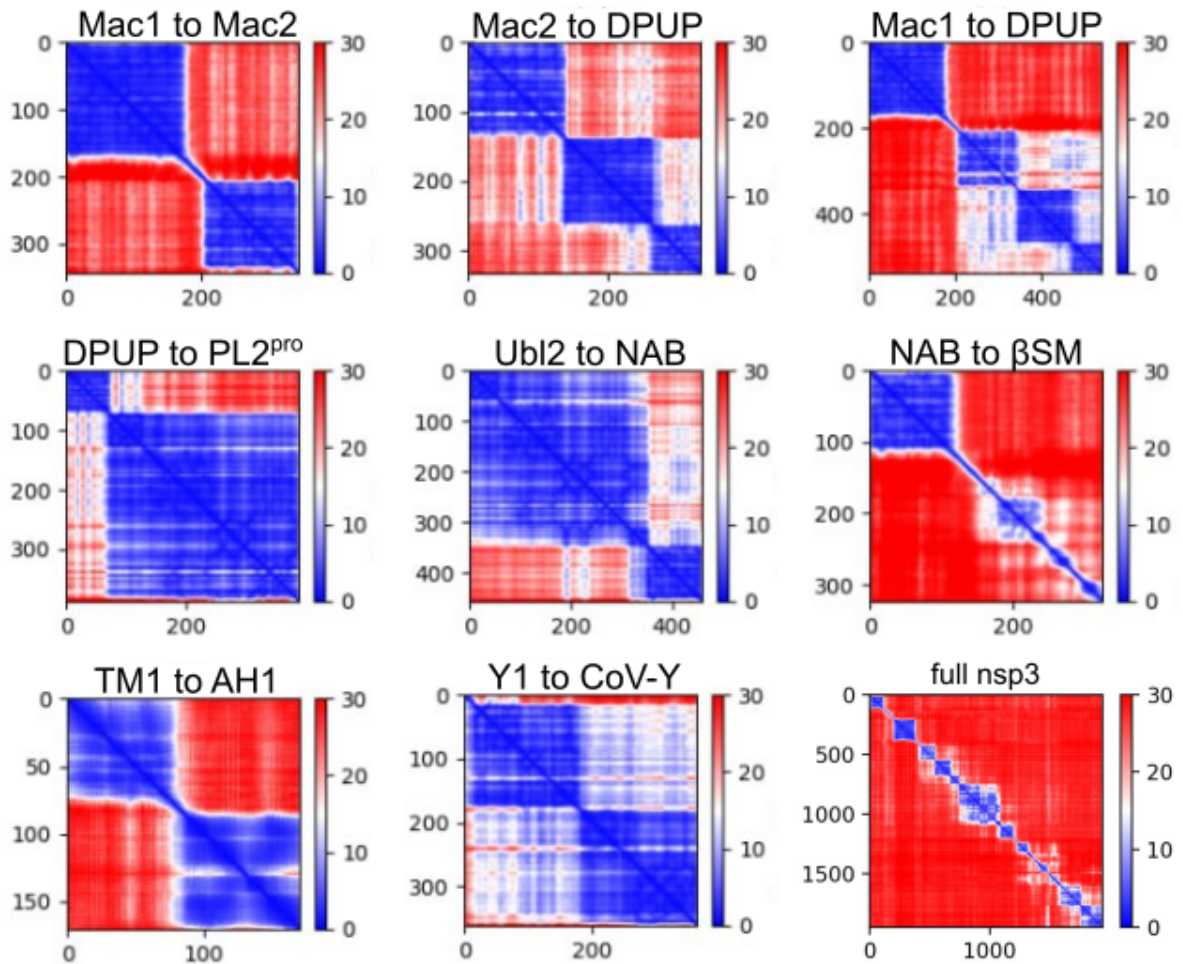
The lower right matrix from Figure 13 shows the pAE for the full length prediction of NSP3 (Figure 14a). Most of the domains are clearly recognizable from this matrix, but the pAE values are often worse compared to predictions from smaller sequences. However, this results in certain subdomains becoming more visible such as Ubl2 as subdomain of PL2<sup>pro</sup> or the subdomains CoV-Ya and CoV-Yb. Nevertheless, the full length prediction does not provide any other domains to assemble into a quaternary structure with high confidence. Some of the domains from the full length prediction show also lowered pLDDT values compared to the prediction of the single domain (Figure 14b-c), potentially due to large nearby loops and regions of disorder as outlined in section 3.1.4. This affects also the fold, as seen in the Mac2 domain lacking one  $\beta$ -strand and a short  $\alpha$ -helix at the N-terminus (Figure 14b-c).

### 3.3.2 Structure prediction of Y1 and CoV-Y

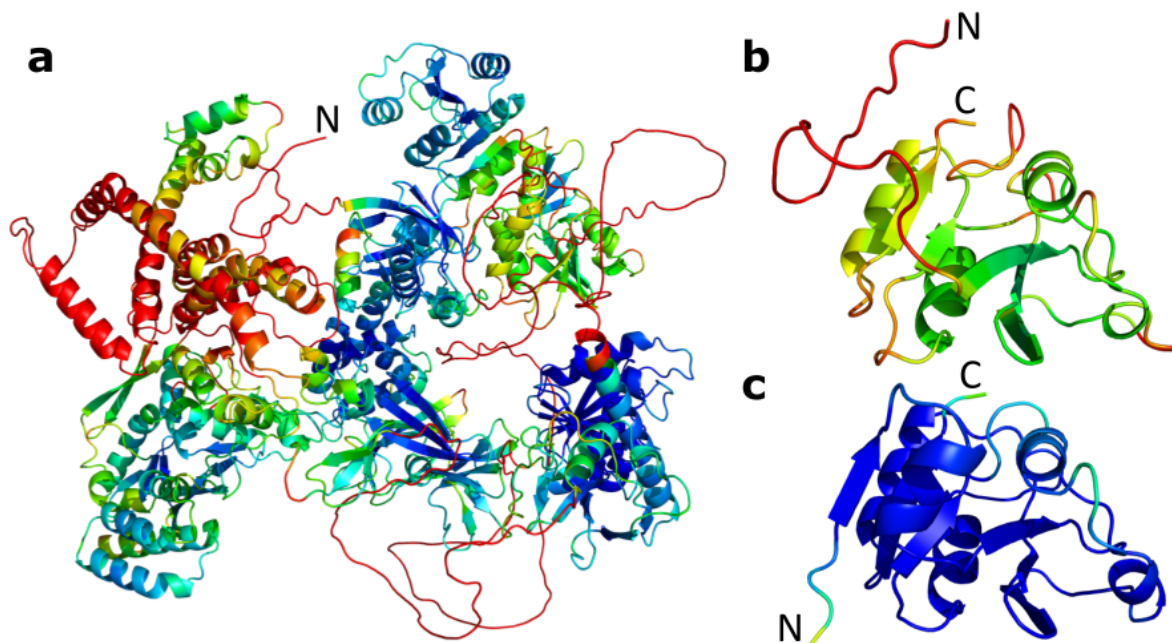
The structure prediction for the segment beyond the transmembrane region, shows for both sarbecoviruses and MHV an N-terminal helical loop followed by a high confidence fold of the Y1 and CoV-Y domains (Figure 15). In total, the segment consists of 362 to 377 residues for all three viruses.

The predicted structure has overall high pLDDT values (Figure 15a) except for the N-terminal helical tail, the C-terminal tail, and a few residues at a linker around residue Gly1763 (SARS-CoV-2 NSP3 numbering). The predicted aligned error (pAE) matrix (Figure 15b) suggests a separation into two domains exactly at residue Gly1763, which marks the border between the domains Y1 and CoV-Y. Furthermore, the values in the pAE matrix between both domains have moderate confidence values between 10 and 20, which speaks for some uncertainty in the relative arrangement of both domains, but not for pure randomness. Running several predictions results always in a similar arrangement.

Two globular folds are recognizable for each domain, which would divide the C-terminus into the subdomains Y1a + Y1b (Figure 15c) and CoV-Ya + CoV-Yb (Figure 15a). A separation of those, however, is not clear from the pAE matrix as the arrangement of the subdomains is highly confident. The experimental structure 8F2E supports the division into these four subdomains [60], where all subdomain structures except for Y1a are experimentally determined. The predicted structures are highly similar to the experimentally solved domains with RMSD values of 0.4 Å to the PDB structures 7RQG



**Figure 13:** Predicted Aligned Error (pAE) matrices from predicted structures which cover multiple domains from SARS-CoV-2 NSP3. X and Y axis show residue number, colored is according to pAE confidence score. Lower pAE values describe a higher confidence in the prediction of the spatial arrangement of a residue pair. Subsequences which show low pAE values to each residue within itself result in blue squares in the matrix and describe folded domains. Red regions describe non-confident spatial arrangements between residues. The predicted structure is thus only a possible, but uncertain outcome and may not reflect reality. White regions describe moderate confidence. Note that the lower right matrix of the full length NSP3 prediction makes the domains and their location on NSP3 visible but shows AlphaFold2 failing to predict a confident fold of the complete protein.

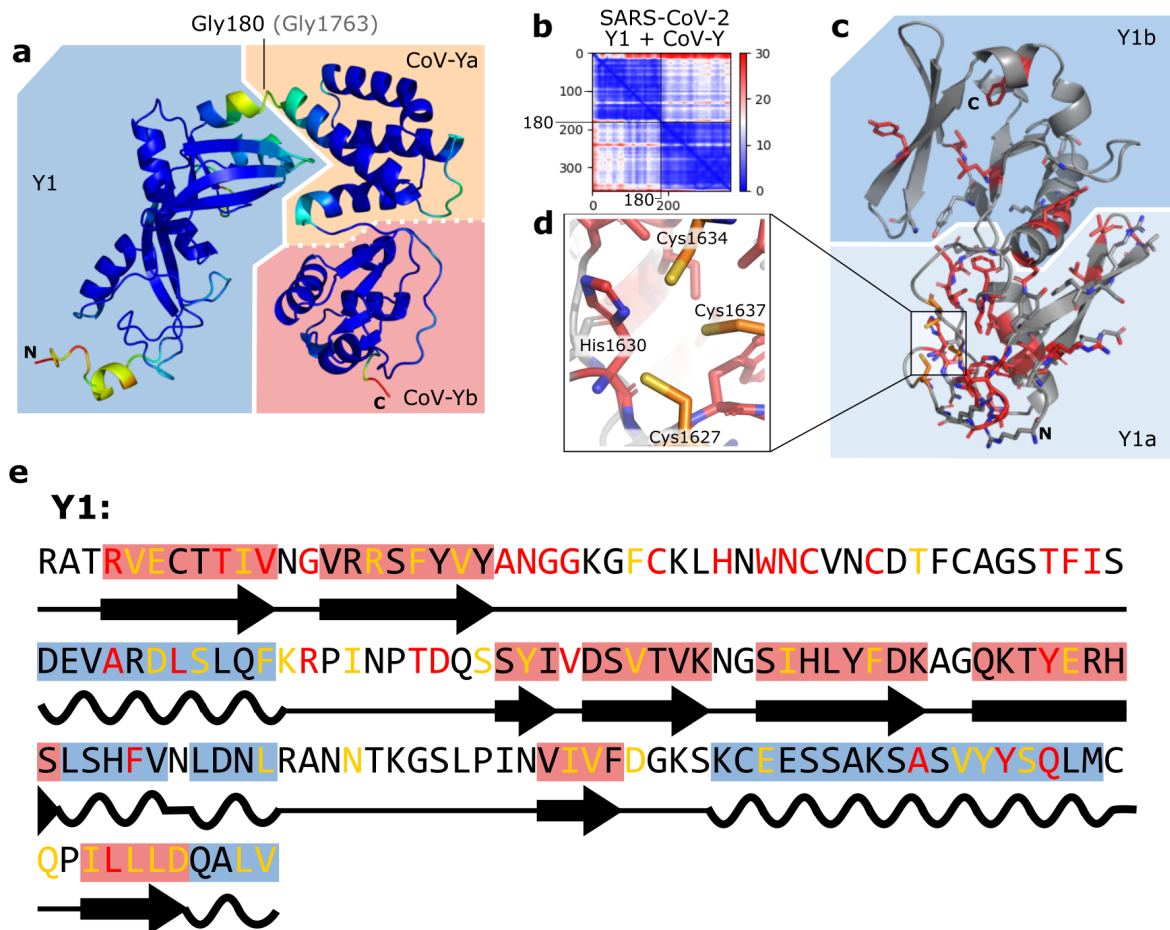


**Figure 14:** *AlphaFold2 prediction of full length SARS-CoV-2 NSP3 (a) and the differences in fold and pLDDT between Mac2 from the full length prediction (b) and from the Mac2 sequence alone (c). Residues of all structures are colored according to their pLDDT, where low confidence values below 50 are red, high confidence values above 90 are blue, and values in between are colored accordingly on a rainbow interpolation. The C-terminus in (a) is obscured by the structure.*

[59] for CoV-Yb and the structure 8F2E [60] covering the subdomains Y1b, CoV-Ya, and CoV-Yb. However, the arrangement of subdomains in 8F2E is much more compact due to Y1a missing. Hence, low RMSD values are only achievable when aligning individual subdomains. Aligning the whole prediction of Y1+CoV-Y with 8F2E results in an RMSD value of 2.3 Å.

The globular fold of Y1a begins with two large  $\beta$ -strands (21 residues in total) followed by 50 residues arranged in interconnected loops and an 11-residue long  $\alpha$ -helix (15c). The loop region contains a conserved cysteine-histidine cluster (15d), which was proposed as a potential zinc-binding cluster [135]. The fold of Y1b is intertwined by numerous structure elements and contains a triplet of parallel  $\beta$ -strands, another triplet of anti-parallel  $\beta$ -strands, one  $\alpha$ -helix, sharp turns, and a loop stabilized by hydrogen bonds (15c). CoV-Ya consists of four  $\alpha$ -helices, while CoV-Yb comprises four short  $\alpha$ -helices, four  $\beta$ -strands, and loops (15a).

The classification into ordered and disordered regions based on pLDDT and secondary structure elements excludes the low pLDDT linker prior Y1 (Figure 15) from the Y1 domain. However, a second potential zinc binding site, also consisting of one histidine and three cysteines, was identified [135] in the linker between the last domain of the transmembrane region, AH1 (amphipathic helix 1), and Y1. In MHV and SARS-CoV-1,



**Figure 15:** *a*: Prediction of linker+Y1+CoV-Y from SARS-CoV-2 colored according to pLDDT with blue as high confident. *b*: the respective pAE matrix showing two confident local folds, one with residues before Gly180 and one with residues after. *c*: Y1 monomer from hexamer prediction, with amino acids conserved among all betacoronaviruses highlighted in red, and the conserved cysteine cluster from (d) in orange. Other displayed residues are involved in hexamer formation but are not conserved. *d*: prediction of conserved cysteine-histidine cluster. *e*: sequence of SARS-CoV-2 Y1. Residues in red are conserved among all betacoronaviruses, while residues in yellow show conserved chemical properties. Red shades indicate  $\beta$ -strands, while blue shades indicate  $\alpha$ -helices.

this linker contains the entire cysteine-histidin cluster, which could fulfil a biochemical function and could thus be considered as part of Y1 or AH1. In SARS-CoV-2, the linker is shorter and the histidine is located inside AH1, making the situation more difficult to resolve without any experimental evidence. Therefore, the linker is neither associated with AH1, nor with Y1 in the domain boundary definitions in Table 5.

### 3.3.3 Conservation of Y1 among nidoviruses

The sequence similarity between both sarbeviruses is outstandingly high for the C-terminal subdomains, ranging from 92 % to 100 % (Table 6). Between SARS-CoV-2 and MHV, the sequence similarity ranges from 58 % to 73 % and RMSD values between predicted structures are always below 0.8 Å (Table 6).

The function of Y1 is unknown, but identifying highly conserved residues can give important clues. Therefore, sequence alignments were first performed with other betacoronaviruses, then with viruses from *Orthocoronavirinae*, and finally with various viruses from *Nidovirales*, as Y1 was previously known as "nidovirus-conserved domain of unknown function" [18].

For the first global multiple sequence alignment, Y1 from SARS-Cov-2 was aligned with the 16 betacoronaviruses listed in Table 19 of the appendix. From the 161 residues, 29 amino acids were conserved with 22 of those being located in the subdomain Y1a (Figure 15c, e). These 22 residues include also the conserved potential zinc binding site [135] consisting of three cysteines and a histidine shown in Figure Figure 15d. Since the cysteines are apart by 3.4 Å to 4 Å, they are unlikely to form disulfide bonds. Furthermore, many of the conserved residues are like the cysteine-histidine cluster located in the loop region of Y1a and form hydrogen bonds, which could potentially stabilize such a large region lacking secondary structure elements. Additional 45 residues show conserved chemical properties.

Sequence alignments for Y1 with viruses from *Orthocoronavirinae* outside of *Beta-coronavirus* (examined viruses are listed in the methods in section 6.3.2) show sequence similarities in the range of 25.9 % to 39.1 %. Structure predictions based on these alignments, however, show high predicted structural similarity with RMSD values from 0.6 Å to 1.3 Å. A BLAST search [76] excluding *Coronaviridae* resulted in 27 hits of which all are artificial sequences such as a recombinant SARS-CoV. Although this search showed no similar sequences from nidoviruses, sequence alignments with subsequent structure prediction were performed with nidoviruses from seven related families (listed in the methods in section 6.3.2). While sequence identities ranged from 18 % to 24.4 %, comparison of

predicted structures showed RMSD values from 5.96 Å to 15.8 Å. If a confident fold resulted from the prediction, it did not resemble the fold of SARS-CoV-2 Y1. In the end, no folds similar to Y1 were identified outside of *Coronaviridae*.

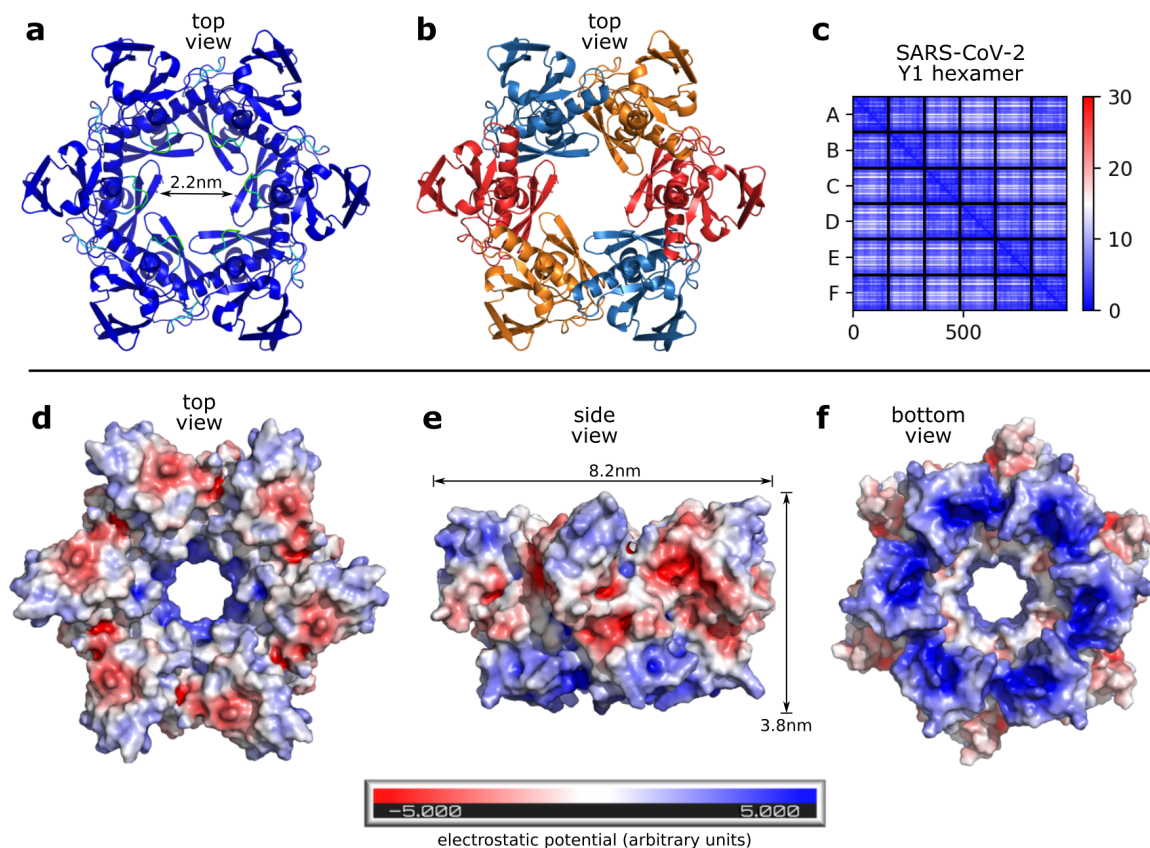
### 3.3.4 Multimer prediction of Y1

NSP3 assembles with NSP4 and NSP6 to a hexameric pore complex, which allows the replicated viral RNA to leave the double membrane vesicles (DMVs) into the cytosol [19]. For MHV, a low resolution map of this complex and the DMV membranes from cryo electron tomography from Wolff et al. [19] is available. Figure 19g shows a side view of this complex sitting on the surface of the DMV, where it forms a crown-like shape that becomes narrower towards the DMV membrane. Since Y1 immediately follows the transmembrane region, it must be located close to a membrane surface and due to the hexameric nature of the pore complex, six copies of Y1 are expected in the complex.

An AlphaFold2 multimer prediction [103] delivers a hexameric assembly of Y1 of high confidence in pLDDT and pAE values (Figure 16a-c). The average pLDDT of 93.6 for the hexamer is highly confident and only lowered by 1.2 in comparison to the average pLDDT of the monomer prediction. Median pAE values per cell in the multimer pAE matrix (16c) range from 2.6 to 7.3, and are with a median of 5.4 for the whole matrix in the confident range of below 10. For reference, the median pAE value for the Y1 monomer is 2.7. Hexamer predictions of Y1 for SARS-CoV-1 and MHV are also highly confident with median pAE values of 7.3 and 4.7, respectively.

The minimal channel diameter of the Y1 hexamers is 2.2 nm wide (Figure 16a) and fits into the proposed diameter of 2-3 nm [19, 32]. Monomer contact is made via hydrogen bonds, which includes also residues conserved in *Betacoronavirus*, which were identified in section 3.3.3. In SARS-CoV-2 Y1, the conserved residues involved in monomer contact are Arg1602, Thr1607, and Gly1611; the non-conserved ones are Ser1615, Arg1642, Asp1654, Leu1718, Ser1729, Ser1734, Lys1737, and Leu1746, adding up to eleven contact residues. Residues involved in hydrogen bonds were located on both subdomains and submitting only the sequence of Y1a or Y1b to ColabFold for hexamer predictions, did not lead to any hexameric ring structure similar to the high confidence Y1-hexamer prediction. Replacing the contact residues with alanine still led to hexamer predictions, since some of the hydrogen bonds were between backbone atoms and where it was not the case new hydrogen bonds emerged at other residues emerged. Replacing the initial monomer-contact residues with proline, however, did no longer result in the prediction of hexameric assemblies.





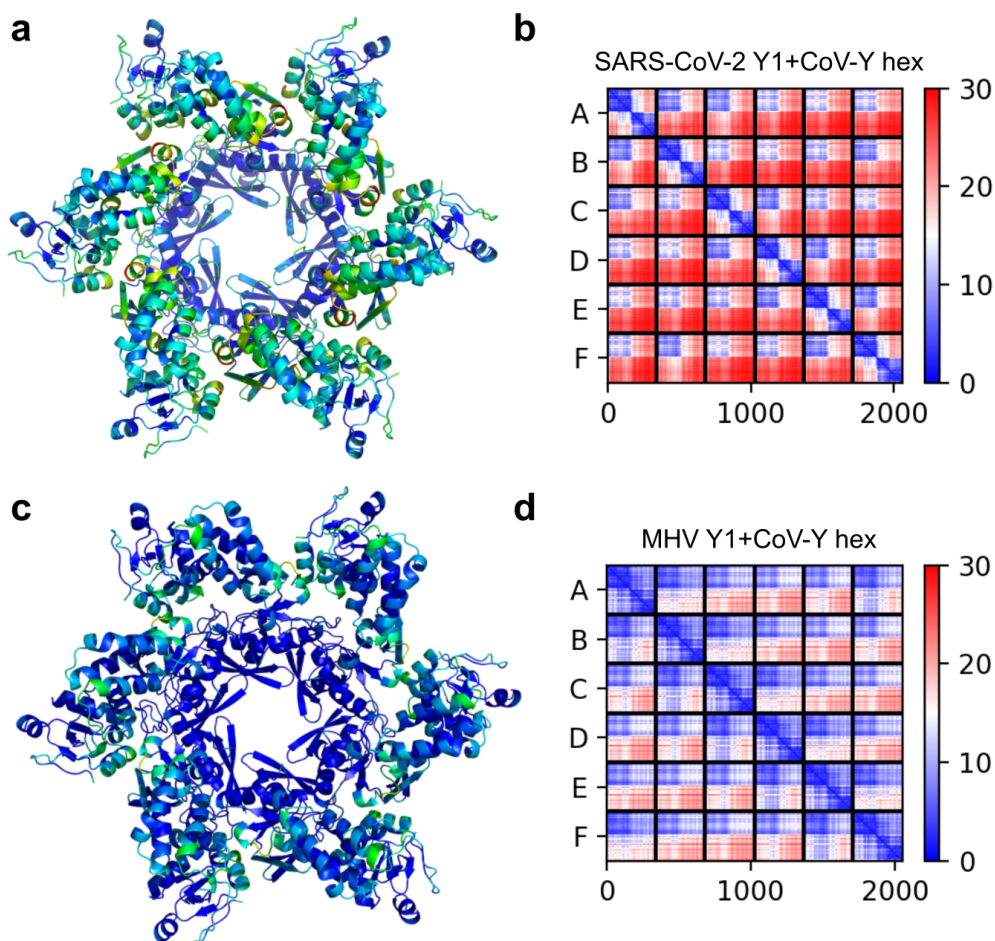
**Figure 16:** *a-c*: high confidence prediction of Y1 hexamer, colored according to pLDDT (*a*) and by chain (*b*). The narrowest diameter of the inner channel is 2.2 nm. The pAE matrix shows a highly confident arrangement of all monomers (*c*). *d-e*: electrostatic surface calculation of Y1 hexamer via PyMOL [131] viewed from different angles. The bottom surface (*f*) and the inner surface along the channel (*d, f*) are positively charged (blue), while the top surface (*d*) and side surface (*e*) are mostly neutral (white) or negatively charged (red).

Vacuum electrostatics of the hexamer (Figure 16d-f) show primarily positive charges at the channel’s inner surface and at the bottom surface towards the membrane, while the top and side surface show mostly negative or neutral charges.

Validation of the Y1 structure and its hexameric assembly were attempted. However, the expression organism BL21 Gold *E. coli* synthesized only insufficient amounts of the Y1-construct.

### 3.3.5 Multimer prediction of Y1+CoV-Y

Predictions of Y1+CoV-Y hexamers result in a similar structure as for the hexamer consisting only of Y1 (Figure 17a), but for SARS-CoV-2 and SARS-CoV-1, median pAE values rise from 5.4 and 7.3 to 20.4 and 21.6, respectively, which is beyond the range of confident predictions. This results primarily from the bad alignment between each CoV-Y and each other domain (Figure 17b). For MHV, however, the overall structure is predicted with much higher pLDDT values (Figure 17c) and the pAE matrix shows a

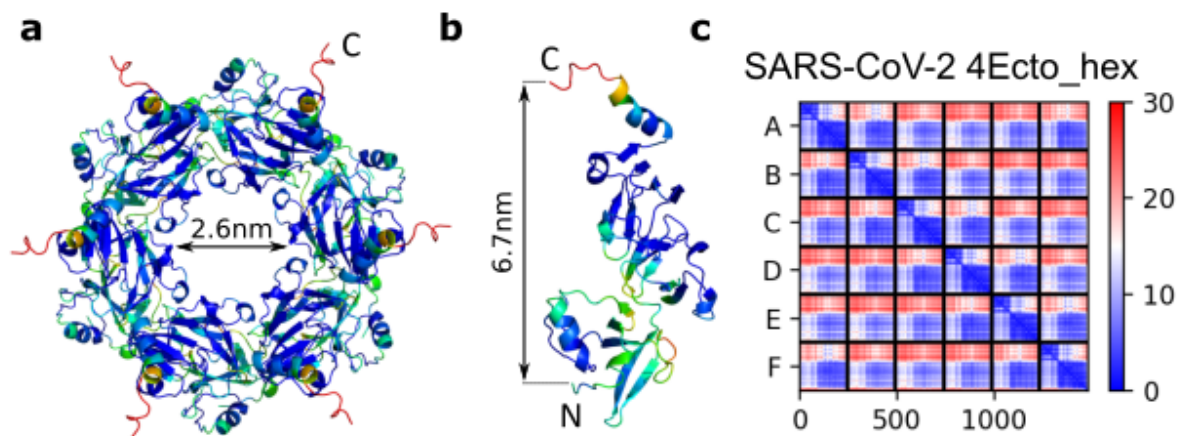


**Figure 17:** Comparison of SARS-CoV-2 and MHV hexamer prediction of Y1+CoV-Y. *a:* hexamer prediction of SARS-CoV-2 Y1+CoV-Y. *b:* pAE matrix for structure shown in (a). *c:* hexamer prediction of MHV Y1+CoV-Y. *d:* pAE matrix for structure shown in (c). (a) and (c) are colored according to pLDDT, with values above 90 shown in deep blue and values below 50 shown in red, with other colors interpolated in between. Y1 domains in (c) are arranged well according to pAE values in (d). CoV-Y is only arranged well to the Y1 of the same monomer and in some cases also well to Y1 of nearby monomers.

better arrangement among all Y1 monomers, between each Y1 and CoV-Y or a monomer, and between several Y1 and CoV-Y domains from different monomers (Figure 17d). The median pAE for MHV is with 11.7 also close to the high confidence range and the average pLDDT remains almost the same, with an average pLDDT of 89.8 for the Y1+CoV-Y monomer and 88.0 for the Y1+CoV-Y hexamer. The major differences between the hexamer from MHV and SARS-CoV-2 are in the orientation of CoV-Y domains, resulting in RMSD values of 2.5 Å.

For all remaining NSP3 domains, no hexameric structure with high confidence was predictable. For ectodomain of NSP4, however, a confident hexamer was predicted (Figure 18). NSP4 follows the domain-sequence Nterm-TM1-4Ecto-TM2-TM3-TM4-Cterm\_domain (where TM stands for transmembrane domain) and its ectodomain is like 3Ecto located in the ER lumen [27]. 4Ecto is predicted with two subdomains (Figure 18b), which are





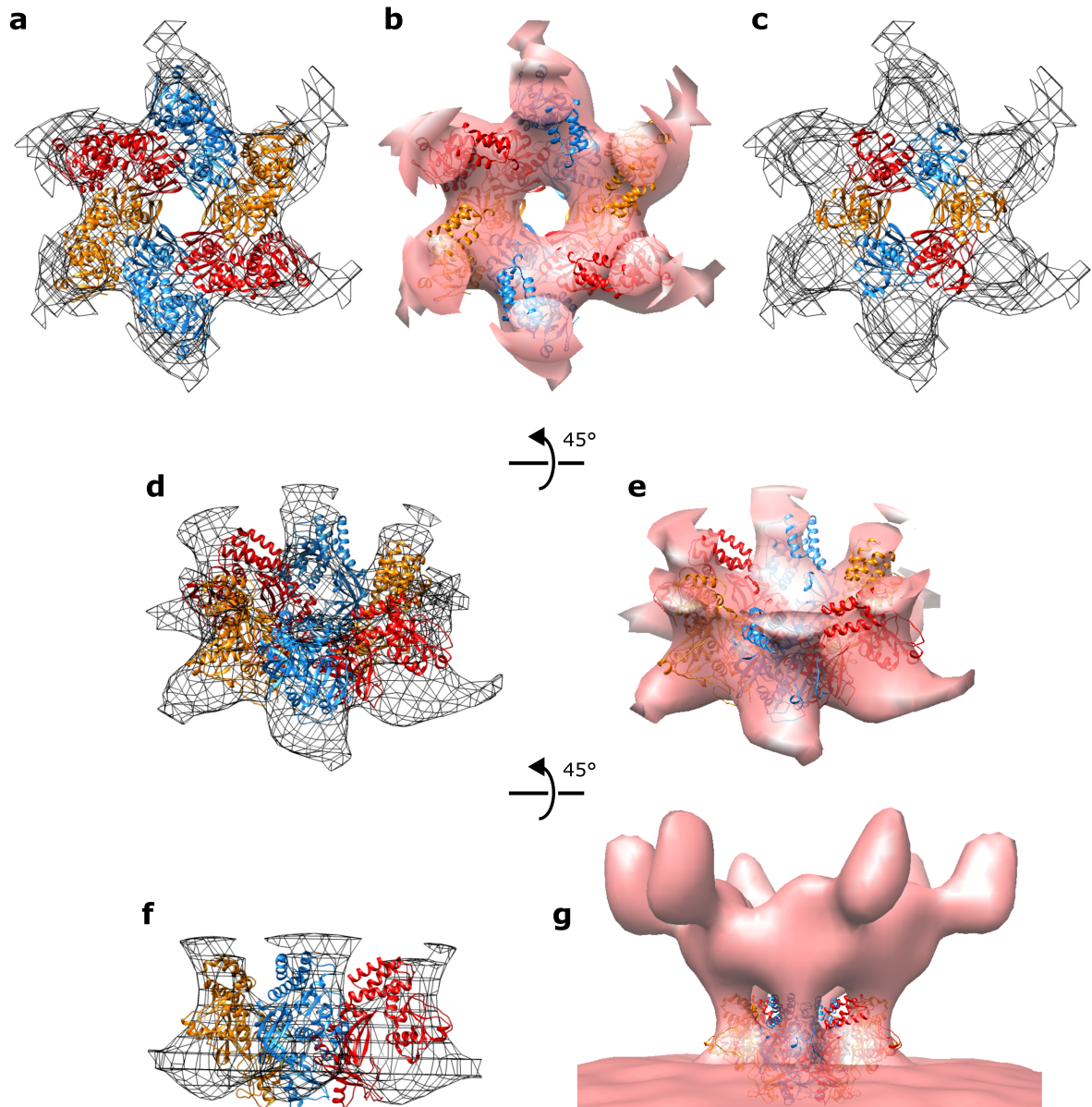
**Figure 18:** Prediction of the SARS-CoV-2 4Ecto hexamer. *a:* Hexameric structure of the large ectodomain from NSP4 colored according to pLDDT. *b:* side view of a single monomer from the hexameric structure. *c:* respective pAE matrix of the hexameric 4Ecto prediction. Note that the termini of 4Ecto are connected to transmembrane helices in the full structure.

flexibly linked as suggested by the pAE matrix (Figure 18c). The C-terminal subdomain is predicted with higher pLDDT values and while the pAE matrix suggests no confident arrangement between the N-terminal subdomains within the hexamer, a high confidence assembly is predicted for the C-terminal subdomain. The median pAE value of the entire matrix is 12.8 and the diameter of the emerging inner pore is 2.6 nm. In the predicted orientation, both termini, which are attached to transmembrane helices in the full NSP4 structure, are at opposed sides and are 6.7 nm apart.

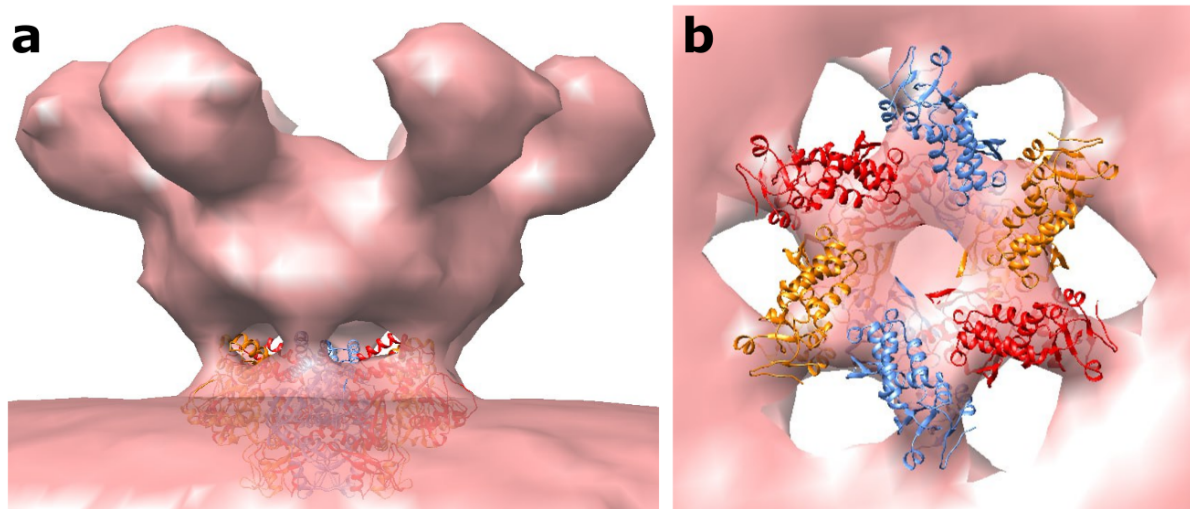
### 3.3.6 Integrative modeling

Since the prediction of MHV Y1+CoV-Y had the highest confidence (see Figure 17), it was fitted into the cryo electron tomography map from the NSP3+NSP4+NSP6 pore complex from Wolff et al. [19]. Due to the low resolution of 30.5 Å of the map, the atomic structure and map volume did not overlap perfectly. Nevertheless, the fitting algorithm of UCSF Chimera [66] positioned the structure at the complex's base, where Y1 is embedded partially in the membrane region and CoV-Y extends into the pillars of the base (Figure 19). A very similar result is observable for the fit into the subtomogram averaged cryoET map from Zimmermann et al. [32], which shows the hexameric assembly of SARS-CoV-2 NSP3 and NSP4 at 20.3 Å resolution (Figure 20). Moreover, Figure 20b shows well, how at high contour level density is lacking around the disc in the membrane at the base of the pore, but that stronger density is still visible in a shape that agrees with the hexameric protein structure.

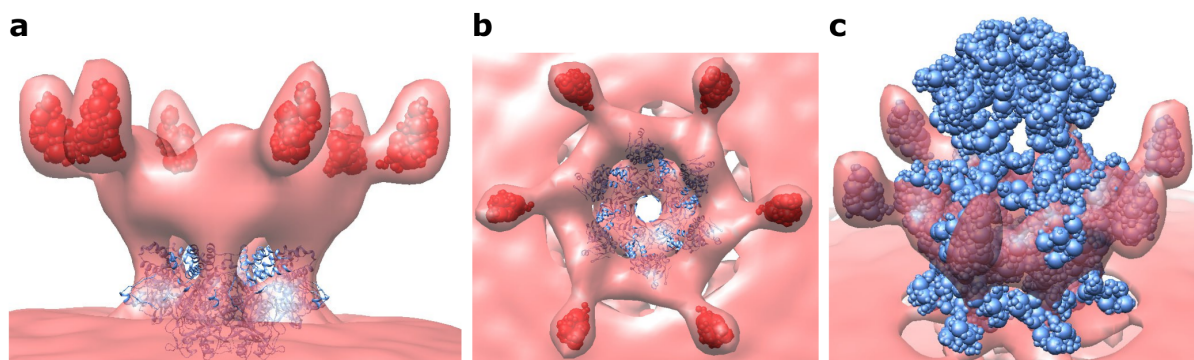
In order to fit all cytosolic NSP3 domains into the map volume, experimentally deter-



**Figure 19:** *a-b: Top-view of predicted model of Y1+CoV-Y hexamer from MHV fitted into the cryo electron tomography map EMD-11514, with volume shown as grid (a) and surface (b). c: same fit, but only Y1 hexamer is shown. d-e: same model-fit as in a-b with rotated camera angle by 45°. f: same model-fit as in d with rotated camera angle by 45°. For more clarity, only half of the hexamer is shown. g: side view of the model fitted into the map with the whole pore-complex visible. All images were generated in UCSF Chimera [66] with a contour level of 2.89. For images a-f, volume was truncated to highlight the complex's base.*



**Figure 20:** *Fit of SARS-CoV-2 Y1+CoV-Y hexamer prediction into the 20.3 Å subtomogram averaged cryoET map EMD-15963. Fit and visualization were made via UCSF Chimera [66]. Side view (a) shows the structure fitted into the lower disk and columns of the complex's density at a contour level of 0.354 (Chimera uses arbitrary units). Top view (b) shows the same fit at less visible density at a contour level of 1.85.*



**Figure 21:** *Results from integrative modelling of MHV NSP3 with Assemblin [119]. a: side view of fitted Y1+CoV-Y hexamer and Ubl1 domain into cryo tomography map of MHV pore complex (EMD-11514). The predicted model of the MHV Y1+CoV-Y hexamer is shown as blue atomic model, while the Ubl1 domain is shown in red as coarse grained model. b: same fits viewed from above. c: attempted fit of all cytosolic NSP3 domains. The coarse grained models occupied the entire volume, which forced some of the domains to be located outside the volume at low density regions above the center of the complex. All images were generated in UCSF Chimera [66] with a contour level of 2.89.*

mined or predicted structures were combined with additional restrictions in an integrative modelling approach. In this work, the software Assemblin [119] was used to fit the NSP3 domains of MHV into the cryoET map of the MHV pore complex [19]. Data from cross-linking experiments was not available. Instead, the length of linkers were utilized for distance restraints between domains as these were defined in the section 3.1.5.

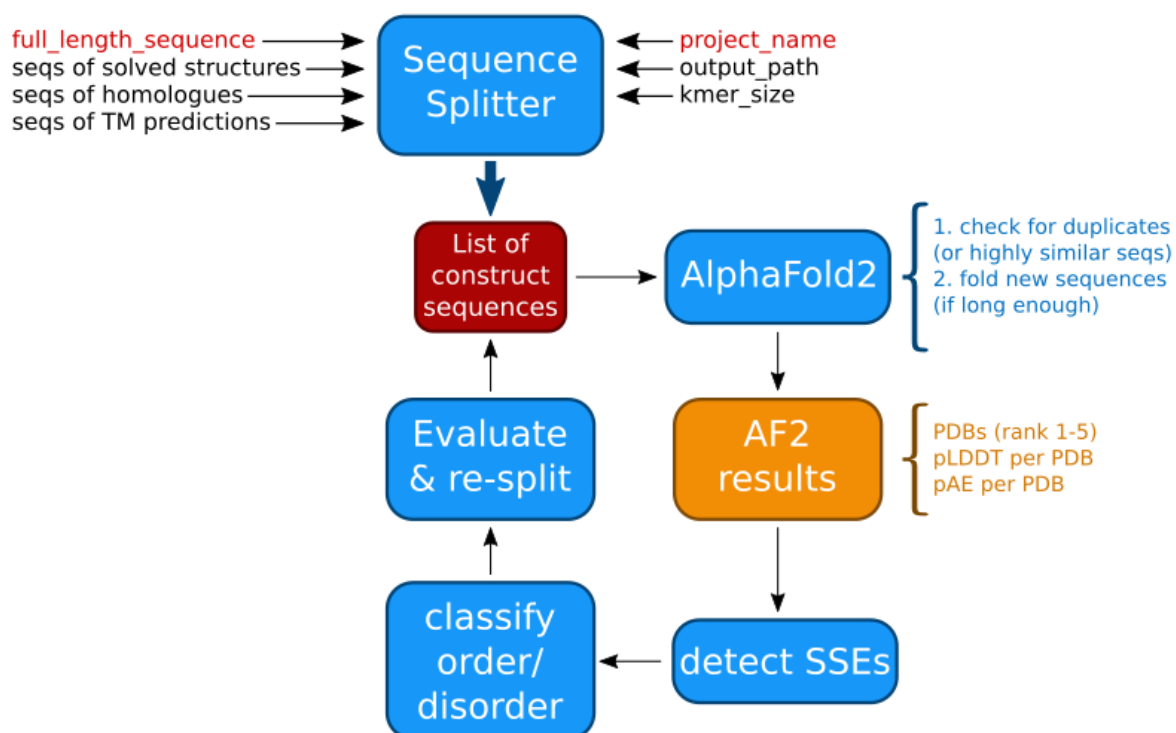
First, Assemblin [119] was tested to fit the Ubl1 domain and the Y1+CoV-Y into the expected locations, which also generated fit-libraries for the subsequent steps. From experiments with GFP fused to Ubl1 it was shown that Ubl1 is located in the outermost prolongations of the pore complex’s crown [19] and indeed, the six Ubl1 domains were fitted into the prolongations without any further restrictions (Figure 21a, b). However, the Ubl1 domains were positioned in the center of the volume of the prolongations, leaving unoccupied space. The hexameric Y1+CoV-Y structure was positioned by Assemblin at almost the same location as the structure placed by UCSF Chimera’s fitting algorithm [66] (Figure 21a, b), where the RMSD between both positions and orientations is 1.2 Å. While the fitting of these domains followed the expectations, the fitting of all cytosolic domains remained unsuccessful after a large number of runs. The fitted models occupied the entire volume, which led to the positioning of few domains outside the strongest density of the map (Figure 21c). Above the center of the map was a disconnected volume of weak density, which was filled with domains. The density of the membrane regions was excluded, hence no domains were fitted there. Since the global optimization fitting failed and distance restraints were violated, no refinement runs were performed.

### 3.4 Automated domain determination

The here presented method of determining domain ranges via AlphaFold2 based on the classification of segments into ordered or intrinsically disordered, was applied manually on NSP3 from SARS-CoV-2, SARS-CoV-1, and the murine hepatitis virus (MHV). However, the individual decisions leading to the final domain ranges are based on thresholds and numeric data and the whole process can therefore be automated. The next sections describe the implementation and validation by comparing the automatically generated results to the manually determined domain ranges of NSP3.

#### 3.4.1 Implementation

The developed python script used AlphaFold2 to predict structures from initial sequences, where the structure prediction results were used to classify sequences into ordered or disordered, as described in previous parts of this work. Figure 22 illustrates the workflow.



**Figure 22:** Workflow of the automated domain detection. The sequence splitter provided a list of construct sequences, which were given to AlphaFold2. Predicted structures, pLDDT and pAE values were used to detect secondary structure elements via DSSP. The sequences were classified into ordered or disordered, which was evaluated and leading to the next iteration of construct sequences. Required inputs are shown in red.



The minimum inputs are the full multidomain protein sequence and a project name, but sequences from solved structures, homologs, and transmembrane domain predictions can also be added. The sequence splitter processed this information to provide initial sequences, which should in the end be either completely disordered or crystallizable constructs without any disorder. To accomplish this, AlphaFold2 was used to predict structures from these sequences. A DSSP algorithm identified secondary structure elements, which were used together with pLDDT to determine whether a region is ordered or disordered. The pAE values were used to identify subdomains. All of this information was evaluated and sequences were redefined to either contain only ordered regions or disordered region. If no changes were made, a previous set of sequences was generated, or after a fixed number of iterations, the algorithm terminated.

### 3.4.2 Benchmarking on database of multidomain proteins

The method was developed and tested on SARS-CoV-2 NSP3, where fine tuning of the parameters was made until the domain ranges matched those listed in Table 4 and Table 5 with an error margin of 5 residues in each direction.

For a full validation via a benchmark study, the MDB carroll database from the Columbus State University (<http://csc.columbusstate.edu/carroll/MDB>) was used, since it contained protein sequences and information about the number and ranges of domains. A custom script was written to automatically test the domain detection algorithm on the database sequences and evaluate automatically how much the results deviated from the database entries. However, it turned out that many database entries had domain ranges not matching with the structures, which were often only solved after the database was set up. Since it was unclear how many cases were incorrect, it was not possible to determine the accuracy of the domain detection algorithm. Due to time limitations and focus on other projects, this project was not further developed.

## 4 Discussion

Non-structural protein 3 (NSP3) is of great interest in the fight against the pandemic caused by SARS-CoV-2, since multiple of its domains serve essential roles in the viral replication cycle. Nonetheless, available information of these domains was incomplete or in some cases contradictory between different sources of literature. And due to its hexameric assembly into a large complex, large amounts of disorder, and transmembrane domains, its complete structure is not easily solvable, nor predictable.

In this work, the domain ranges of NSP3 were updated and completed for the SARS-CoV-1, SARS-CoV-2, and the murine hepatitis virus (MHV). This was accomplished with a new method developed in this work, which utilizes AlphaFold2 to classify protein regions as ordered or intrinsically disordered. Since this method is suitable for domain identification and construct design of any multidomain protein, the entire technique was automated. Furthermore, a hypothesis for the function of the NSP3 domain Y1 was developed and the presence of a folding, previously unnoticed, domain between PL2<sup>Pro</sup> and NAB was bioinformatically analysed and experimentally validated.

### 4.1 Utilizing AlphaFold2 for domain boundary determination and construct design

#### 4.1.1 Classification of regions as ordered or disordered

Although predicted models are not replacing experimental measurements, they can assist in the design of crystallizable constructs [114], development of hypotheses, and provide models for molecular replacement or integrative modeling. In case of AlphaFold2, one of the key features is the predicted local distance difference test (pLDDT), a confidence measure that does not only indicates the quality of a predicted model, but also serves as indicator for intrinsic disorder [20]. Together with the predicted aligned error (pAE), one can utilize these metrics for manual and automated domain determination. For that matter, atomic precision of the entire prediction is not necessary as long as predicted structures are locally correct. Structure prediction can then be reduced to the prediction of secondary structure elements on the sequence level, which builds the foundation for classifying regions as ordered or disordered and is sufficient for construct design followed by experimental structure solution. For other use cases such as integrative modelling, structure prediction must of course be comparable to experimentally determined high resolution structures.

The results from superimposing the AlphaFold2 predictions with experimentally determined structures from the PDB show a correct prediction of secondary structure elements in all examined cases (see Table 3 in the results section). Furthermore, in all but two cases, predicted models showed high structural similarity with root mean square deviation (RMSD) values below 1 Å, making most of these predictions as good as experimentally determined structures. The exceptions, predicted Ubl1 from SARS-CoV-1 and MHV compared to the NMR structures 2GRI and 2M0A, include a flexible N-terminus, which is the main reason for the high RMSD values. From the remaining domains, the structures for Y1b, CoV-Y, and  $\beta$ SM were published after the release of AlphaFold2 and were thus not included in the training data set, which shows that the prediction of unknown structures is generally possible. This is in agreement with the large scale analysis from Terwilliger et al. [115], which identified about 80 % of predicted side chains to closely match those of experimentally determined structures, while 94 % are at least matching roughly. Thus, AlphaFold2 is not only capable of reliably predicting secondary structure elements correctly, but also predicts protein folds with high accuracy.

Williams et al. observed the prediction of so-called "barbed-wire" conformations [132], which are intrinsically disordered regions predicted as long loops extending far from the main fold, often as straight line or barbed-wire-like structures. These come with low pLDDT values [132] and in this thesis such regions were also observed to be predicted falsely as  $\alpha$ -helices as in the case of  $\beta$ SM (see Figure 6 in the results section). Identifying such helices as false  $\alpha$ -helices is therefore necessary when using secondary structure elements for domain determination. In the examined cases, false prediction of  $\beta$ -sheets was not observed, but a large scale study would be required to quantify the error rate for this secondary structure element. Because  $\beta$ -strands only occur in proximity to at least one other  $\beta$ -strand, they are only found in locally compact folds and therefore such an error rate must be much smaller as it is for  $\alpha$ -helices, which are also predicted in disordered loops pointing away from the globular fold.

Predicting and superimposing multiple structures of  $\beta$ SM led to the observation of a decent superimposition of ordered regions, while disordered termini were not superimposed (Figure 6d). An easier approach for identifying false helices was to consider the pLDDT, which is a good indicator for intrinsic disorder [20]. In the case of  $\beta$ SM and full length NSP3 prediction it was observed, that large disordered loops have a negative impact on the pLDDT of nearby folds. For the full length NSP3 prediction, it lowered pLDDT and altered the predicted fold of the Mac2 domain (Figure 14), compared to the



prediction of just the Mac2 domain. Since adding disordered regions worsened the quality of a prediction, removing it from the sequence seemed to be a valid option. For  $\beta$ SM, cropping the regions with barbed-wire conformation at the termini did not only increase the pLDDT of the main fold, but also lowered the pLDDT of falsely predicted  $\alpha$ -helices. This behaviour was not described before and can be applied in additional iterations of domain determination, if a region could not be classified unambiguously as ordered or disordered. It also highlights, that predicting secondary structure elements alone is not sufficient for a reliable classification and that the pLDDT must be considered. Likewise, the pLDDT alone is not sufficient as nearby disorder negatively affects a fold. Secondary structure elements and pLDDT together, allowed to define the classification tree depicted in Figure 5 based on observations. This also allows to differentiate between predicted transmembrane helices and false  $\alpha$ -helices, since transmembrane helices come in random arrangements but with high pLDDT values in contrast to the falsely predicted  $\alpha$ -helices, which also come in random arrangements but with low pLDDT.

It is important to note that pLDDT calculation depends on the coverage of the sequence during the multiple sequence alignment step of AlphaFold2 [15], as low sequence coverage by similar sequences comes with lower pLDDT values compared to sequences with high sequence coverage. The lower boundary for the number of similar sequences must be evaluated by large scale studies, but until now this fact can explain cases where plausible looking folds are not predicted with high pLDDT values. In the case of NSP3, a coverage by 14 sequences with sequence identities from 50 % to 100 % (for DPUP-like domain from MHV) was sufficient to reach pLDDT values above 80 for most of the residues. For lower sequence identities, starting from 20 %, a coverage by 30 sequences was sufficient for the same manner (for SARS-CoV-2 Mac2).

In conclusion, predicted secondary structure elements, consideration of "barbed-wire" conformations [132], and the pLDDT are together providing sufficient information to classify a protein region into ordered or disordered, which in turn enable design of crystallizable constructs for experimental structure determination. While large scale studies are still required to identify error rates for false predictions of secondary structure elements, iterative improvement of sequences by cropping disordered termini can lead to improved prediction confidence. While the predictions remain hypotheses, they are cheap to generate and have large potential of assisting in subsequent validation experiments.

### 4.1.2 Domain boundary determination of NSP3

NSP3 is highly relevant in the fight against COVID-19 and due to the many open research questions surrounding this protein's structure and functions, prediction software holds the potential of providing new insights. Furthermore, NSP3 is for several reasons an interesting protein for exploring the capabilities of structure prediction software. First, NSP3 consists of several domains, which provide a great variety of folds. While some domains consist mostly of  $\alpha$ -helices such as CoV-Ya, others contain mostly loops such as Y1a. The domains HVR and  $\beta$ SM introduce large regions of intrinsic disorder and complicate the arrangement of domains to each other. Other domains, such as Ubl2 and PL2<sup>pro</sup> test the capability of differentiating closely arranged subdomains. The transmembrane region adds another factor of difficulty and the multimerization of NSP3 expands the complexity of predicting a full length NSP3 structure correctly.

Only a fraction of all domains was structurally solved for SARS-CoV-2, SARS-CoV-1, and the murine hepatitis virus (MHV), and not all residue ranges of domains were defined clearly. With the capability of AlphaFold2 to classify regions as ordered or disordered, a new method of determining domains opened up. Prior structure prediction, sequence based techniques were a viable option. Since disordered regions tend to show low sequence identities compared to sequence identities of folded domains, sequence alignments even allowed to differentiate roughly between disordered regions and regions potentially folding into domain [78]. Because folded domains have a selection pressure for maintaining the fold and function, they accumulate mutations slower than disordered regions, which usually lack this kind of selection pressure. The results from sequence alignments between both sarbecoviruses and MHV agree with this assumption, where the *Betacoronavirus*-specific linker domain stood out as its sequence identities were comparable with those from folded domains, thus making it look more like a domain instead of a flexible linker (see Table 2). The presence of the fold was later validated in a SAXS experiment and is discussed in section 4.2. This result is contradictory with current literature, where the assumption of this region being a linker may have emerged from the incremental annotation of domains over the past. PL2<sup>pro</sup> and NAB were annotated first [54] and the region in between is with 34 residues relatively small. Since all linkers were not explicitly stated as linkers, such a small non-annotated region may easily be confused as a linker instead of a potential domain. Furthermore, the sequence analysis revealed this domain to be specific to *Betacoronavirus*, which decreases chances of finding this domain accidentally. This case emphasizes the importance of clear annotations, which indicate whether a segment of a multidomain protein was predicted to be ordered or disordered, experimentally validated

to be such, or if it is still unexplored.

At the start of this work, ColabFold [128] had a limitation of input sequence length, which did not allow a full length prediction of NSP3 and set therefore the requirement of dividing it into suitable input sequences, listed as preliminary sequences in Table 1. Later, that limit was raised to 4000 residues, but the preliminary domain ranges turned out to be still useful as large regions of disorder negatively impact the prediction quality of nearby folds. Predicting full length NSP3, for example, lowered the pLDDT of domains compared to the predictions from single domain sequences (Figure 14). Nevertheless, predicting the complete structure of a multidomain protein is a good starting point for defining domain boundaries without any prior knowledge, as the pAE matrix and secondary structure elements give good indications on how to define preliminary domain ranges, which can then be refined. The final domain ranges for NSP3 of SARS-CoV-2, SARS-CoV-1, and MHV listed in Table 4 and Table 5 are the first complete lists of ranges for each domain and linker, with each single residue being assigned to one of both. Nonetheless, only experiments can validate these partially theoretical ranges. Furthermore, the ranges define crystallizable constructs, but domains may also be defined by functionality, which could involve interaction between ordered and disordered segments.

### 4.1.3 Automated domain determination

The automated domain determination is a promising tool, as the here presented method of classifying regions of a protein into ordered or disordered based on AlphaFold2 predictions had a straight forward decision tree (Figure 5), which could be easily automated. However, the large benchmark set contained entries with domain ranges not agreeing with experimentally determined structures. Since it was unclear how many of such cases exist, it was not possible to determine the performance of the newly developed algorithm. Unfortunately, time was scarce and the project was put on hold. A robust predictor of domains, however, is still desirable, which gives this project some relevance.

## 4.2 Experimental validation of the Betacoronavirus-specific linker domain ( $\beta$ SLD)

### 4.2.1 Structure analysis

The segment between the domains PL2<sup>pro</sup> and NAB was previously not described as a domain [17, 54] and can therefore be confused as a linker, which is reinforced by its small

size of 34 residues. The results of this thesis show that this sequence is conserved among *Betacoronavirus*, and AlphaFold2 predicts in all examined cases a plausible fold with high confidence. The folding nature of this segment was later validated via solution small-angle X-ray scattering (SAXS), but an atomic-resolution structure of the linker domain awaits to be solved. Nonetheless, such a structure is expected to be highly similar to the predicted model as the arrangement of secondary structure elements, hydrogen bonds, and location of conserved residues within the fold are all plausible.

The two  $\beta$ -strand pairs  $\beta 1+\beta 2$  and  $\beta 3+\beta 4$  should be stable in isolation, especially  $\beta 1+\beta 2$  with the additional hydrogen bonds between the side chains. The connection between both pairs relies on two hydrogen bonds between  $\beta 4$  and  $\beta 1$ , which force the domain to fold into a compact shape. Due to the small size, there is no hydrophobic core that could drive the initial folding process. The SAXS experiment, however, suggests that this domain is rather folded than a loose linker. Furthermore, 12 out of the 14 residues preserving chemical properties (according to PAM matrix [136]) among betacoronaviruses are found in loops or directly adjacent to those, which indicates a strong selection pressure on the loop regions. This pattern fits well with the predicted structure, since  $\beta$ -sheets are formed by hydrogen bonds between backbone atoms and experience thus lower selection pressure. The short loops of two residues between the  $\beta$ -strands of a  $\beta$ -sheet are also very restrictive for possible amino acids. While it is not obvious which roles the chemically conserved residues play in the maintenance of the fold, the conserved proline in the central loop comes with restricted torsion angles and is thus limiting the flexibility of the loop. This may play a role in stabilizing the fold, hence the conservation of this residue. Furthermore, two hydrogen bonds within the loop may add additional restriction of movement. The solution small-angle X-ray scattering measurements agree with the predicted structure on a low resolution level. While this does not validate the structural details, it validates the presence of a globular fold between PL2<sup>pro</sup> and NAB. All in all, the structure seems plausible and the pattern of conserved residues agrees with the secondary structure elements. Only one residue was recognized as Ramachandran outlier, but due to its location at the loose N-terminus, its predicted torsion angles are not relevant. Due to its previous history as a linker and its specificity to *Betacoronavirus*, the name "Betacoronavirus-specific linker domain", in short " $\beta$ SLD", felt appropriate.

Although this domain was shown to be present only in NSP3 of betacoronaviruses (section 3.2.2), similar sequences were found beyond NSP3 in the endoRNase in numerous viruses, when aligning the linker domain sequence with the whole polyprotein. This

approach was necessary since most viruses had no annotations on NSPs, but it led to this surprising observation, which may contain hints about the evolution of this domain. Structure prediction of this sequence from endoRNAses shows a similar composition of secondary structure elements, but with an additional  $\alpha$ -helix between the central loop and the third  $\beta$ -strand (Figure 10). Furthermore, the high pAE values, non-globular shape, and the missing hydrogen bonds between  $\beta 1$  and  $\beta 4$ , indicate this arrangement to be unstable. It is possible, that a gene duplication event occurred in the evolution of *Betacoronavirus*, where part of endoRNase was duplicated and translocated to NSP3. Selection pressure towards a stable domain could then have led to the deletion of the  $\alpha$ -helix. Since this helix consists of only 10 residues in most viruses, the deletion process could have happened relatively quickly. More interestingly is that within *Betacoronavirus*, the sequence lengths are consistent from residue number 8 onwards (in SARS-CoV-2  $\beta$ SLD; Figure 9), speaking for a stable state of that domain which reached an optimal sequence length for its function. What kind of function this domain fulfills, however, cannot be determined from this structure prediction alone.

#### 4.2.2 Potential functions of the linker domain, NAB, and $\beta$ SM

Among 17 betacoronaviruses which were aligned in a multiple sequence alignment, five out of the 34 residues are conserved and nine have similar chemical properties (according to PAM matrix [136]), suggesting a functional role of the linker domain. Besides no similar sequence being present outside of *Betacoronavirus*, no similar folds were identified in a similarity search on the PDB and no region of the fold seems like it could be enzymatically active based on the low number of reactive residues and secondary structure elements. The only remaining clues one could use to construct a hypothesis regarding this domain's function is by considering its following domains, NAB and  $\beta$ SM, which are also conserved in *Betacoronavirus* [17, 18]. The three domains of this large region of roughly 350 residues may have co-evolved due to this common taxonomic specificity, which would mean that the domains interact in a structural or functional way.

The function of the Betacoronavirus-specific marker domain ( $\beta$ SM) remains unknown, as only recently the first structure of the central fold ( $\beta$ SM-M) was deposited to the PDB under the code 7T9W, where the paper remains unpublished. The nucleic acid binding domain (NAB) on the other hand, is able to bind ssRNA and to unwind DNA, with a higher binding affinity to ssRNA [57]. Since NSP3 is part of the hexameric RNA-exporting pore complex [19, 32], NAB is potentially involved in exporting new copies of the viral RNA genome from the interior of the double membrane vesicles into the cytosol. In order

to make use of this functionality, the binding site of NAB must be located towards the inner channel of the pore. The details of possible assembly modes are discussed in section 4.3.6, as the hypotheses require first a discussion of the hexameric assembly of the Y1 domain, which is discussed in section 4.3.4. Here, it is sufficient to note, that  $\beta$ SLD and NAB are likely interacting structurally as they are adjacent to each other. Furthermore, if NAB is involved in RNA export, it must be oriented towards the pore and as discussed later, it cannot be located at the narrowest part of the complex. Since this domain cannot form a hexameric ring of the observed diameter [19, 32] on its own, multiple domains must assemble together. Here,  $\beta$ SLD is a promising candidate due to its proximity to NAB, while the position of  $\beta$ SM-M is not restricted due to the large disordered linker  $\beta$ SM-N. Also, a notable observation is the high sequence similarity of  $\beta$ SM-N and  $\beta$ SM-C (see Table 6), as both domains are assumed to be disordered [17], especially for  $\beta$ SM-N which is predicted without any secondary structure element. This raises the question for potential functions in interaction with other proteins or domains for both subdomains.

It is worth noting, that *Gammacoronavirus* NSP3 possesses the Gammacoronavirus-specific marker domain ( $\gamma$ SM), which was described to be similar to  $\beta$ SM [17]. Because no structure of  $\gamma$ SM was solved, structure prediction was used to compare this region to  $\beta$ SM, revealing high similarity between the experimentally solved structure of SARS-CoV-2  $\beta$ SM-M and the structure prediction of  $\gamma$ SM-M of the Canada goose coronavirus (Figure 12). For all other gammacoronaviruses structures were predicted as well, but these differed from  $\beta$ SM-M. Nonetheless, the result suggest  $\beta$ SM to be not completely specific to *Betacoronavirus*. Because no presence of NAB nor  $\beta$ SLD was identified in Canada goose coronavirus,  $\gamma$ SM has likely evolved first, if the structural similarity can be validated experimentally. In such a case,  $\beta$ SM-M could serve a function independent of the other two domains specific to *Betacoronavirus*.

Regarding future experiments for identification of function, mutation experiments would be the best starting point. Disrupting the folds of the Betacoronavirus-specific domains provides insight into the impact of these domains on the viral fitness, by measuring viral titers or by quantifying the number of replicated RNA via qRT-PCR [137]. If all three domains are only involved in NAB's RNA binding functionality, the impact is expected to be non-lethal, as coronaviruses outside of *Betacoronavirus* are able to export RNA without a NAB domain. However, it is also possible that  $\beta$ SLD and  $\beta$ SM-M are regulating the activity of PL2<sup>PRO</sup> or any other enzymatically active domain which appears close during complex assembly. Another option to shed light on the interaction partners

of these domains are *in vitro* binding assays [138], but these may not reflect the true interaction within the pore complex *in vivo*.

### 4.3 C-terminal domains and the hexameric pore complex

#### 4.3.1 Nomenclature of Y1 and CoV-Y

The cytosolic C-terminal region of NSP3 follows immediately after the transmembrane region and comprises the domains Y1 and CoV-Y [17] (Figure 2). Based on the predicted aligned error (pAE) and low pLDDT values around Gly1763, the exact domain ranges are defined as residues 1599-1759 for Y1 and residues 1766-1945 for CoV-Y by the method described in this work, where the domain names are according to Lei et al. [17]. Furthermore, both domains consist of two closely arranged globular folds. Thus, the whole C-terminus consists of four subdomains and the last three of these globular folds were solved experimentally by independent research during the time of this work. In this dissertation, the four subdomains were labeled according to the two main domains as Y1a (residues 1599-1664), Y1b (residues 1665-1759), CoV-Ya (residues 1766-1847), and CoV-Yb (residues 1848-1945) respectively (see Table 5).

Unfortunately, older literature did not define the border between Y1 and CoV-Y, which led to conflict in recent nomenclature: The PDB structure 7RQG [59] covers the fourth subdomain (residues 1844-1945) of Y1+CoV-Y and is deposited under the name Y3 (paper not published yet). The second experimentally solved structure from a more recent publication (PDB structure 8F2E) [60] labels the same subdomain as Y4 and another domain (residues 1764-1847) as Y3, making the name Y3 ambiguous.

In the publication of 8F2E [60], the C-terminal region is also divided into four domains, which closely resemble the ranges from this work. The main difference lies in their names: Y1 (residue range not given), Y2 (residues 1665-1763), Y3 (residues 1764-1847), and Y4 (residues 1848-1945). This study, however, suggests the three latter subdomains to be part of CoV-Y, while the results from this work suggest the first two subdomains to be part of Y1 and only the last two as part of CoV-Y. The reasons for sticking to Y1a, Y1b, CoV-Ya, and CoV-Yb, are due to the close arrangement between the subdomains and sequence conservation, and are further explored in the following section.

#### 4.3.2 Structure prediction of Y1

Domains should ideally comprise one functional unit and, in the case of globular proteins, form also a tightly packed structure. For the C-terminal region of NSP3, Y1a and Y1b

can be considered as two subdomains of the larger domain Y1, since both are consistently predicted with a high confidence arrangement to each other. The arrangement makes sense from a structural biology perspective and, as seen in the section 3.3.4, the multimer prediction leads only to a hexameric structure if both subdomains are included in the input sequence, where both subdomains participate in hydrogen bonds connecting the monomers. CoV-Y is predicted similarly with two subdomains which are closely arranged to each other, also with high confidence in the pAE values. However, the prediction of both, Y1 and CoV-Y, indicates a clear separation between these two at the residue Gly1763 (SARS-CoV-2 NSP3), as seen in the pAE matrix in Figure 15b. This potentially flexible hinge is also observable in the PDB structure 8F2E, which validates the predicted folds of Y1b, CoV-Ya, and CoV-Yb with RMSD values of 0.4 Å.

Currently, the structure of Y1a has not been experimentally determined and in a collaboration with David Briggs it was not possible to purify a construct comprising Y1a and Y1b, which seemed to be toxic to *E. coli*. However, it is likely that the experimental structure is similar to the prediction. The predicted structure has an unusually high content of loops, but most of these residues are connected via hydrogen bonds to each other. Moreover, this domain is relative to other domains highly conserved and most of the residues of Y1a conserved in *Betacoronavirus* are in the loop region, including a cysteine-histidine cluster made of three cysteines and a histidine, which was identified as a potential zinc binding site [135]. Binding zinc could stabilize the loop region and the conservation could be the result of the very specific fold lacking  $\alpha$ -helices and  $\beta$ -sheets, which would be vulnerable to mutations disrupting this fold, hence creating a strong selection pressure against many possible mutations. This is also supported by the highly similar structure predictions from other betacoronaviruses, where only few structure predictions from showed RMSD values above 1 Å and where the highest difference was 1.3 Å.

Conclusively, the predicted structure seems plausible and the conserved residues make sense in the provided fold. The hexameric assembly, as discussed later in section 4.3.4, seems plausible as well and is supported by experimental and theoretical models, which strengthens the plausibility of the Y1 fold even further.

### 4.3.3 Conservation of Y1 among *Nidovirales*

The domain Y1 shows high conservation in *Betacoronavirus*, but from sequence analysis it was previously known as "nidovirus-conserved domain of unknown function" [18]. Furthermore, double membrane vesicles as replication organelles appear in all positive stranded RNA viruses [63], making the need for an RNA-exporter channel present in a large variety



of pathogens. Multiple results of this work indicate Y1 to form a hexameric ring and to function as the base of this RNA-exporting pore complex (discussed in the next section), which would fit well with a highly conserved domain. However, no BLAST [76] search results and no similar predicted structures were identified outside of *Orthocoronavirinae* and sequence similarities stayed below 30 %. While all related RNA viruses with double membrane vesicles require a similar mechanism of exporting only the positive-sense single strand RNA, the structures driving this process may differ drastically, as it becomes evident from very different proposed domain compositions by B. W. Neuman [18].

One example for alternative pore complexes is the Flock house virus (FHV), which induces replication organelles from mitochondrial membrane instead of endoplasmic reticulum membrane, and forms a dodecameric RNA exporter channel [64]. This virus is from the genus *Alphanodavirus* belonging to the order *Nodamuvirales* and its pore complex was also imaged via cryo electron tomography [64]. Both, SARS-CoV-2 and FHV, are distantly related as the different origins for their replication organelles suggest, but they are still sharing the need for an RNA-exporting pore complex. It is therefore thinkable, that distantly related nidoviruses have alternative structures which fulfill similar roles.

Conclusively, the results suggest that it is unlikely for Y1 to be conserved in *Nidovirales*, but its conservation in *Orthocoronavirinae* and especially *Betacoronavirus* indicate an important function. The most likely function is driving the assembly of the hexameric pore complex, which is explored in the next section.

#### 4.3.4 Multimer prediction of Y1 and Y1+CoV-Y

NSP3, NSP4, and NSP6 are part of the RNA-exporting hexameric pore complex [19, 32], where the interaction between the ectodomains of NSP3 and NSP4 induces the formation of double membrane vesicles [27, 32]. But which domains mediate the assembly of the large cytosolic side of the complex and where are all of NSP3's domains located within the assembly? Since the membrane topology of NSP3 is known [42], Y1 is on the cytosolic side of the DMV. Because Y1 immediately follows the transmembrane region with only a 10-residue linker in between, Y1 cannot be located far from the DMV's outer membrane. Furthermore, one or few domains must form the foundation of the pore complex, and there are no options for Y1 to be located elsewhere. A recent study on this complex from SARS-CoV-2 led also to subtomogram averaged cryoET map, where the domains Ubl1 to Ubl2 are deleted from the NSP3 construct [32]. Since Y1 and CoV-Y were included in the construct, it must be located in the remaining density, which still contains the base plate of the complex. Moreover, it is certain that the domains must assemble either alone or

together with additional domains into a hexameric ring structure. Since Y1 and CoV-Y are predicted as large structures, which is partially validated by the experimental PDB structure 8F2E, it is very likely, that Y1 and potentially Y1+CoV-Y are part of this ring. Conveniently, AlphaFold2 introduced a feature to predict multimeric assembly [103] and delivered a high confidence prediction of a hexameric Y1 structure (Figure 16).

The predicted hexamer of Y1+CoV-Y fits the diameter and shape of the pore's foundation of both complexes, of MHV (Figure 19) and of SARS-CoV-2 (Figure 20), and the structure's channel diameter of 2.2 nm agrees with the observed channel diameter of 2-3 nm [19, 32]. The pAE values of the prediction for the Y1+CoV-Y hexamer from MHV are not perfect, which is potentially due to the hinge between CoV-Y and Y1 and the lack of hydrogen bonds between the CoV-Y domains within the multimer. Much better pAE values between monomers come with the prediction of Y1 hexamers for all three viruses, which surpass any other pAE matrix for multimers observed before (Figure 16c). This high confidence in pLDDT and pAE values comes probably from the hydrogen bonds between 11 residues in each monomer, from which three are conserved among all examined betacoronaviruses and where a fourth one is considered chemically similar according to the PAM matrix [136]. Since hydrogen bonds are only predicted between Y1 domains and not between CoV-Y domains in the Y1+CoV-Y hexamer, the hexameric assembly is held together by Y1 alone, which is supported by experimental evidence from Li et al., where a construct comprising Y1b+CoV-Y (PDB 8F2E) crystallized in monomeric form [60]. Although the ectodomains of NSP3 and NSP4 and their surrounding transmembrane domains are sufficient to let NSP3 and NSP4 bind to each other, it was observed in deletion experiments in MHV that both proteins bind more effectively when Y1 and CoV-Y are present [27]. Furthermore, both cryoET maps do not show density between the pillars in which CoV-Y is located unless the contour level is decreased drastically, which further suggests no interaction between CoV-Y domains and Y1 as a strong contributor to the assembly complex assembly. Due to the low resolution of both cryoET maps, however, the exact position of the structure within the foundation of complex does likely differ, especially along the normal axis of the membrane.

Calculated electrostatics of the predicted Y1 hexamer show primarily positive charges at the bottom surface towards the membrane which includes the N-terminus and is thus attached to the transmembrane domain. This would allow the hexamer to interact with the negatively charged lipids of the DMV and the fits into the cryoET maps shows that membrane contact is unavoidable. Additionally, the hexamer has positively charged surfaces along the inner channel, which would make it suitable for the export of RNA due

to the negatively charged phosphate backbone.

Conclusively, a high confidence hexameric ring structure for Y1 is predicted, which agrees with experimental evidence from membrane topology and fits well into the cryo electron tomography maps of the pore complexes of MHV and SARS-CoV-2. This fits can be extended with the hexamer prediction of Y1+CoV-Y and calculated electrostatics agree with Y1's adjacency to the transmembrane region and the RNA-exporting ability of the pore complex. Finally, Y1 is suggested to be the major contributor to the assembly of the pore complex's foundation, which is supported by monomeric crystal structures of CoV-Y and conserved residues in Y1 which participate in hydrogen bonds between individual units of the multimer.

#### 4.3.5 Experimental validation

Purification of a Y1 construct was attempted by David Briggs, but unfortunately, yielded no sufficient quantity of the protein for further experiments. Of the whole C-terminus, which comprises more than 300 residues, Y1a is the only domain which is not experimentally solved and no other experimental data regarding this domain is published. It is therefore plausible to assume that other researchers were also unsuccessful with this domain. One potential problem could be the large positively charged surface emerging upon hexamer formation, which in theory could bind lipid membranes, DNA, and RNA and interfere thus with the metabolism of the cell. Monomers show also areas of strongly positive charges with computed vacuum electrostatics (Figure 44).

Attempting expression and purification in alternative expression systems would be the next step. Since expression in *E. coli* failed, new attempts in eukaryotic cells may be more promising. If this step remains difficult in other systems and if hexamer formation is a valid reason, redesigning the construct to mutational variants incapable of hexamer formation may be the way to go. In AlphaFold2, mutation of residues which form hydrogen bonds between monomers still led to the prediction of hexameric complexes as alternative hydrogen bonds can be formed with backbone atoms. Mutating to proline, however, leads to monomers which are not assembling together and to high pAE values. Validating this prediction experimentally would be interesting in regard to the capabilities of AlphaFold2, as whether effects from mutations are predicted correctly [105, 139] or incorrectly [140, 141] is still to debate.

With a successful purification, one can attempt structure solution starting with crystallization trials. However, an advisable first step would be to test the sample on a size

exclusion chromatography column to determine whether or not the construct forms a hexamer on its own. Mutational variants incapable of hexamer formation may be more suitable during experimental validation and could moreover be utilized in analysis of the hexamer's impact on the viral fitness, which could finally clarify Y1's status of a potential drug target.

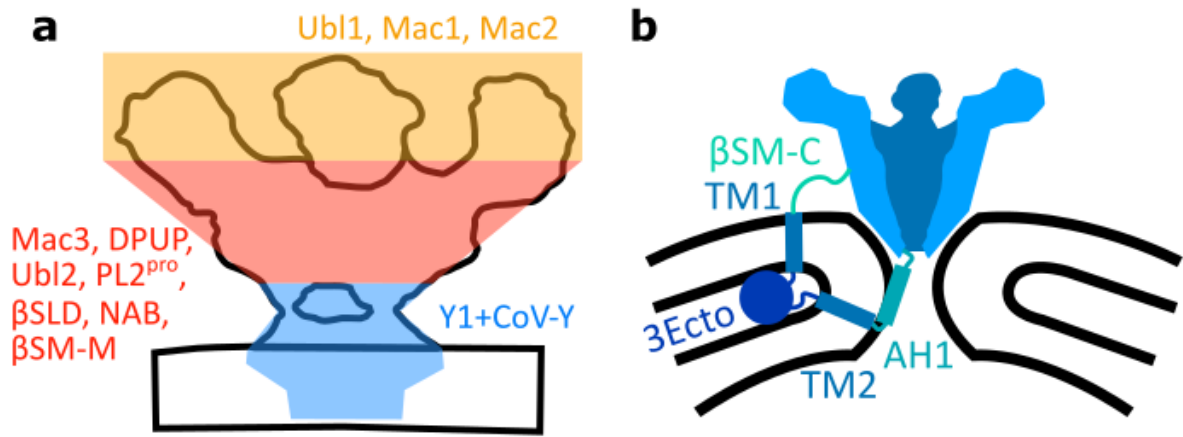
#### 4.3.6 Integrative modeling

A few low resolutions structures of the pore complex are solved and available as subtomogram averaged cryo electron tomography (cryoET) maps. These show either the wildtype complex of MHV from infected cells or the complex with an NSP3 construct N-terminally fused to GFP [19]. The other maps show the complex from cells only transfected by SARS-CoV-2 NSP3 and NSP4 in wildtype form, without the NSP3 domains Ubl1 to Mac1, or without the NSP3 domains Ubl1 to Ubl2 [32]. While the domain composition of NSP3 differs between MHV and SARS-CoV-2 (Figure 2.1.2), the general shape of the complex remains: a short disc embedded in the outer DMV membrane builds the foundation, from which six pillars rise towards the cytosol and hold a massive hexameric ring, from which six prolongations stretch away from the complex (see Figures 20 and 21). Due to the resolution difference of 10 Å, these maps are not directly comparable. The major differences, however, are visible in the prolongations, which were experimentally shown to contain the domains Ubl1 and Mac1 and probably also contain in SARS-CoV-2 Mac2 and Mac3, which are absent in MHV [19, 32]. The remaining important question: where are all the other domains within the assembly?

With the two studies of the complex [19, 32] and additional studies about NSP3 in the past, it is possible to gather strong clues about the position of each domain. The location of Ubl1 was shown in two experiments: via GFP fusion in MHV [19] and via deletion mutants in SARS-CoV-2 [32]. Since the bound GFP added density to the volume and its fluorescence was clearly measurable [19], it cannot float around freely and must be attached to Ubl1, which in turn must be attached to the rest of the complex. However, Ubl1 is linked to Mac1 via the disordered domain HVR and Mac1 is attached to Mac2 via a 34-residue linker. In a deletion experiment where the domains Ubl1, HVR, and Mac1 are removed from NSP3, the cryoTM map shows much shorter prolongations, which speaks for a location of Ubl1 and Mac1 in that segment of the complex [32]. Because Mac1 is connected to Mac2 with a linker, it is possible that Mac1 is not assembled adjacent to Mac2. Since HVR is disordered, it is floating in the cytosol, is wrapped around other domains, or even both as its length ranges from 96 to 158 residues (Table 4). This has

interesting implications for the functionality of the entire complex, which are discussed at the end of this subsection. The next domains of SARS-CoV-2 are Mac2, Mac3, DPUP, and Ubl2. These domains were also removed from NSP3 in a deletion experiment, where the resulting pore complex lost its crown-like structure beyond its foundation [32]. Therefore, PL2<sup>pro</sup>,  $\beta$ SLD, NAB, and  $\beta$ SM are not able to assemble into a stable hexameric ring without the domains Mac2, Mac3, DPUP, or Ubl2. Since all of these domains were shown to possess biological functions, it is astonishing that they additionally provide structural stability. The  $\beta$ SM domain is next to HVR and the Mac1-Mac2 linker the only remaining disordered domain. Since all previous domains are located somewhere in the upper "crown" and Y1 is likely at the base of the complex together with the transmembrane region,  $\beta$ SM could easily function as a linker between the upper ring and the first transmembrane helix. As discussed in section 4.3.4, Y1 and CoV-Y must be located at the base of the pore complex, which is supported by the fit of the predicted structure into the cryoTM maps of both viruses, calculated vacuum electrostatics, close linker to the transmembrane region from the N-terminus of Y1, and additionally due to the proximity of the C-terminus of CoV-Y to the outer DMV membrane. The latter is interesting since NSP3 and NSP4 are connected at this C-terminus and the very first residues of NSP4 are functioning as a transmembrane helix. Zimmermann et al. showed that preventing PL2<sup>pro</sup> from cleaving this connection results in deformed DMVs, tighter spacing of the DMV membranes, and no formation of visible pore complexes [32]. Because DMV-like structures are forming anyways, NSP3 and NSP4 must be embedded in the membrane to allow the interaction between both proteins' ectodomains [27]. Therefore it is possible, that the remaining polyprotein NSP3-NSP4 (potentially with NSP1 and NSP2 still attached) is first inserting itself into the membrane and cleaving NSP4 from NSP3 afterwards. The whole process is likely highly complicated and a dynamic assembly of the complex is expected, but the study from Zimmermann et al. on DMV morphology based on different NSP3 mutants provided a good foundation [32].

The only remaining domains of the complex are the ectodomains of NSP3 and NSP4. To understand their location it is important to note that the DMV interior was previously the cytosol and that the DMV lumen is the space between both membranes, which was previously the lumen of the endoplasmic reticulum [27]. Figures 1 and 2.1.3 illustrate the anatomy of the DMVs. The spacing between these two membranes measures 16 nm [32], which allows to fit easily the large ectodomain of NSP4 inside (Figure 18). Towards the pore complex, the membranes curve and form a pore [32]. While the confidence of the NSP4 ectodomain hexamer is not optimal, it is still promising as for the



**Figure 23:** Schematic illustration of the pore complex and rough regions, where specific domains must be located based on the discussed reasons (a). The hypothesis on the transmembrane region and connection to the complex via  $\beta$ SM-C is shown in (b). The scale of the complex and membranes is based on the subtomogram averaged cryoTM map and images from Zimmermann et al. [32].

second subdomain of the assembly, the pAE values are very low between each monomers and the channel diameter matches the observed channel diameter of the pore complex. However, this domain is located in the lumen between the two membranes of the DMV. Therefore, if the NSP4 ectodomain forms a multimeric ring, it must have a much larger channel diameter, where the curving membrane passes through the ring. Alternatively, NSP4 is not forming a multimer on its own at all and is just interacting with the NSP3 ectodomain to form DMVs [27]. In the latter scenario, however, it could still be possible that NSP4 monomers are interacting. For the ectodomain of NSP3, it is only certain, that it is located in the DMV lumen as well and that it must be interacting with the NSP4 ectodomain. A schematic illustration of the pore complex and regions where all domains must be located after considering all restrictions is shown in Figure 23.

The main role of the complex as a whole is the export of RNA from the DMV interior towards the cytosol. The Y1 hexamer at the base fits well into this function, as the structure is predicted with a positively charged channel which would be able to export negatively charged RNA. RNA export could be supported by the nucleic acid binding domain (NAB) and by Mac3, since both domains are capable of binding RNA [57]. DPUP, the domain immediately following Mac3, potentially supports the RNA binding ability of Mac3 by being positively charged at neutral pH [57]. However, Mac3 is specialized in binding G-quadruplexes, which are specific 3D structures of folded RNA [52, 55, 56] and are not present on the positive strand of viral RNA from SARS-CoV-2 [127]. Therefore, Mac3 either binds only the negative strand, although it is predicted to have only one

stable G-quadruplex, or it interacts with host RNA [127]. Clarifying these roles would give strong clues on whether these domains are oriented towards the channel or towards the cytosol. Vice versa, solving the pore complex structure assembly first could give good indications about the functions of these domains.

Lastly, the pore complex is involved in arranging the RNA and the nucleocapsid protein (N-protein) for packaging of the viral genome. Ubl1 was shown to bind N-protein [33] and electrostatics calculations show that it is strongly negatively charged. Since this interaction between Ubl1 and N-protein was shown to be essential [47] and N-protein changes its conformation upon binding Ubl1 and RNA [33], Ubl1 is likely catching N-protein from the cytosol and bringing it in close contact with the exported RNA. This step is essential, because only RNA packaged in N-protein is binding to the M-protein at the virion assembly site and is therefore only then packed into new virions [39]. Now comes the role of HVR into play. As mentioned above, deletion experiments have show Ubl1 and Mac1 to assembly to prolongations of the complex [32], which means that the large HVR is either wrapping around the complex, floating in the cytosol, or both. This is supported by the fact, that Ubl1 function is independent of HVR [142]. It was shown that HVR interacts with N-protein [49], since HVR is able to mimic RNA [17, 48]. To make use of this function, at least parts of HVR must float in the cytosol. Here, they could catch N-proteins and guide them to the nearby Ubl1 domain. It was also shown, that Ubl1 is stronger interacting with N-protein than HVR [49], which fits well with this hypothesis. To test this hypothesis, an NSP3 construct lacking HVR could be tested against the wildtype for viral fitness. Since it seems like Ubl1 and Mac1 are assembling without HVR to the pore complex (which can be tested as well), the transfected cells should still produce intact pore complexes. However, it could still be required to leave a short linker of about 10-20 residues from HVR, as the relative orientation between Ubl1 and Mac1 is not clear. With this  $\Delta$ HVR-NSP3, one can measure the viral fitness by quantifying viral titers and compare them to the control. Analysing the quantity and shape of double membrane vesicles would ensure, that the mutation has no impact on complex assembly and DMV formation, as Hagemeyer et al. did for examining the interaction between the ectodomain and NSP4 [27].

Finally, in an integrative modelling approach it was tested to fit all predicted domain structures into the cryoEM maps. Since the cryoET map of MHV is at a resolution of 30.5 Å, only the shape and envelope of the folded domains were important, but at the same time it was not easy to fit domains unambiguously into one position and orientation of the

map at that resolution. Unfortunately, this limitation turned out to make it impossible with the given approach to solve the complex's structure. At least, it also fitted the Y1+CoV-Y hexamer at the same position as the UCSF Chimera algorithm. The study by Zimmermann et al [32] was only found during writing of this dissertation. Since no time was left to perform the time and resource intensive model fitting, no additional integrative modeling for the SARS-CoV-2 pore complex was performed. In that study, a manual fitting was presented. However, it left out Y1 and CoV-Y completely, which would therefore contradict the here presented fitting. Furthermore, no publication featuring the Y1+CoV-Y hexamer was found to this day, despite the low effort to generate it. Therefore, no experiments regarding the function of Y1 are available.

Conclusively, NSP3 and the pore complex it is forming are highly complex. Each of its 17 domain fulfills a different role in the viral replication cycle, while also functioning in the assembly of the complex itself, making NSP3 truly to a coronaviral Swiss army knife [13]. Structure prediction made models of each domain available, but the low resolution cryoET maps make it difficult to model the entire complex assembly. Nevertheless, each new experiment about NSP3 provides new insights into NSP3 internal and external interactions, which are slowly completing the big picture of non-structural protein 3.

#### 4.3.7 Y1 as potential drug target

The export of replicated RNA from the DMV interior into the cytosol is an essential step in the viral replication cycle. Mutation experiments on the NSP4 ectodomain in MHV have already shown to interrupt the formation of DMVs and the same study identified Y1 and/or CoV-Y to improve the interaction between NSP3 and NSP4 [27]. From these results, it is clear that Y1 is dispensable for DMV formation, but it could still serve an essential role in the export of RNA and/or formation of the pore complex. While the study showed that NSP3 and NSP4 interact in absence of Y1, it is unclear whether a pore complex was assembled. Considering the previous hypothesis that Y1 forms a hexamer at the lower section of the pore complex and contributes a large number of hydrogen bonds, it is thinkable of Y1 being one major driving force of the assembly. The deletion experiments by Zimmermann et al. identified an important role of Mac2, Mac3, DPUP and Ubl2 in the assembly of the upper ring [32], but Y1 and CoV-Y were present in each mutational variant. With deletion of the whole C-terminus one can test whether these domains are essential for the virus functionally and/or structurally. More targeted experiments would aim at preventing the assembly of a Y1 hexamer and examine its effect, while the general Y1 and CoV-Y structures are still intact. The results of such an experiment will finally



determine whether Y1 is essential and whether its hexamer formation is key to the viral replication cycle.

One practical property of Y1 is its high conservation in *Orthocoronavirinae*. If the structures of different viruses are truly as similar as predicted, Y1 could serve as a potential drug target for broad-spectrum antivirals. These could be applied against human-infecting betacoronaviruses, but also against gammacoronaviruses, which infect various bird species and are therefore a potential threat for livestock. The strong selection pressure against mutations in the loop region would make such drugs also applicable against potential pathogenic coronaviruses in the future and reduce the risk of escape variants.

## 5 Outlook

The domains of non-structural protein 3 are involved in various, sometimes essential parts of the viral life cycle of SARS-CoV-2. These include PL2<sup>pro</sup>, which releases NSP1 to NSP4 from the polyprotein [4, 5, 25]; the ectodomain, which induces the formation of double membrane vesicles [27, 65]; Ubl1, which arranges interaction between exported, replicated RNA and the Nucleocapsid protein [33, 47]; Mac1, which participates in immune response reversal [44–46, 53]; and Mac3, a G-quadruplex-binding domain that is essential for unknown reasons [52, 56, 127]. However, only Mac1 and PL2<sup>pro</sup> are extensively researched drug targets and the results of this work add Y1 as a potential drug target to the list. Formation of the RNA-exporting pore is a crucial step, but whether Y1 is actually the major contributor must be shown experimentally. All in all, NSP3 needs further exploration of domain functions. Validating all structure predictions experimentally is the first step towards completing the full picture of this multidomain protein.

## 6 Materials and Methods

### 6.1 Utilizing AlphaFold2 for domain boundary determination and construct design

#### 6.1.1 Sequence information and preliminary domain ranges

Protein sequences of SARS-CoV-2, SARS-CoV-1, and MHV NSP3 were obtained from NCBI (reference ids YP\_009742610.1, NP\_828862.2, and NC\_048217.1, respectively).

At the start of this work at the end of 2021, experimentally determined structures for the following domains were available in the PDB [82]: For SARS-CoV-2, the domains Ubl1, Mac1, Ubl2, PL2<sup>pro</sup>, and NAB. During the work,  $\beta$ SM-M, and CoV-Yb were solved and deposited. Much after the determination of the final domain ranges, the structure 8F2E was deposited, showing Y1b, CoV-Ya, and CoV-Yb. Due to its late release, it is was not considered for domain determination.

For SARS-CoV-1, PDB structures included the domains Ubl1, Mac1, Mac2, Mac3, DPUP, Ubl2, PL2<sup>pro</sup>, and NAB.

For MHV, structures of the domains Ubl1, the DPUP-like domain, Ubl2, and PL2<sup>pro</sup> were available.

Since SARS-CoV-1 was the virus with the most available experimentally solved domains, preliminary domain ranges were set for SARS-CoV-1 first, which then served as template via sequence alignments for the other two viruses. The goal was to map every residue of NSP3 to exactly one domain or linker. First, domain ranges were defined from the sequences of experimentally solved structures from the PDB. These defined the domains Ubl1 (PDB codes 2GRI, 2IDY), Mac1 (2FAV, 2ACF), Mac2 (6YXJ), Mac3 (2JZE, 2JZD, 2JZF, 2RNK), DPUP (2KAF, 2KQW), Ubl2+PL2<sup>pro</sup> (4M0W, 5TL6, 3E9S, 4OVZ, 3MJ5, 5Y3Q, 2FE8, 4OW0, 5Y3E, 5E6J, 4MM3, 5TL7), and NAB (2K87). Structures 2W2G and 2WCT include the sequence of Mac2 and Mac3.

From the two Mac1 structures it was unclear whether residues 356-358 (SARS-CoV-1 NSP3 numbering) belong to Mac1 or the following linker. Unclear regions were resolved in later steps. For the residues 513-526 it was unclear whether they belong to the C-term of Mac2 or the N-term of Mac3 as the structures 2JZE, 2JZD, 2JZF, 2RNK show both cases. However, 2KQV (from other authors) suggests Mac3 to begin at Gly527. Since the residues Glu517 to Leu526 were not modelled in the structure 6YXJ for Mac2, that segment was kept as a linker for the preliminary domain ranges.

The ranges of the domains HVR and DPUP were defined automatically due to their surrounding domains being experimentally solved. The remaining domains were taken from the most recent NSP3 review by Lei et al. [17].

For MHV, the domains Ubl1 (2M0A), DPUP-like (4YPT), and Ubl2+PL2<sup>pro</sup> (5WFI). For SARS-CoV-2, experimentally determined structures covered the domains Ubl1 (PDB code 7KAG), Mac1 (used PDB codes: 6WEY, 6WOJ, 7CZ4, 6YWL, 7BF5, 7KQP), Ubl2+PL2<sup>pro</sup> (7CMD, 7CJD, 7CJM, 7LLZ), NAB (7LGO), and CoV-Yb (7RQG). A large number of structures was available for Mac1 and PL2<sup>pro</sup>, but only the listed structures were used.

The remaining domain ranges of both viruses were obtained by global sequence alignments between each sequence from the SARS-CoV-1 preliminary domain ranges and full length NSP3 from SARS-CoV-2 or MHV. The sequence alignments were performed with Clustal Omega [74]. Transmembrane domains were defined by prediction via TMHMM 2.0 [143]. Remaining domains of MHV NSP3 were defined with the information from Lei et al. [17]. For SARS-CoV-2, the gene annotations of the NCBI entry with the reference id YP\_009742610.1 were used and regions between domains were designated as linkers.

For local pairwise sequence alignments to calculate the sequence identity between the sequences from preliminary domain ranges, EMBOSS Water version 6.6.0 [74] was used.

### 6.1.2 AlphaFold2 predictions

Sequences from preliminary domain ranges were submitted to AlphaFold2 [15] via Google Colab, ColabFold [128], which uses the MMSeqs2 algorithm [144] for multiple sequence alignment. Default settings without templates and without relaxation were used. The AlphaFold2 version v2.1.1 was used for both sarbecoviruses. Since structures for MHV were predicted later, version v2.1.2 was used for MHV. Towards the end of this work, v2.3.1 was available and capable of taking larger sequences as input. Hence, it was used to predict a full length NSP3 structure, which was not possible prior this version.

Predicted models were loaded with the respective experimental structure of a domain, if available, into PyMOL [131] and were there aligned via PyMOL's build in alignment feature, which calculated the root mean square deviation (RMSD) between the two structures. The same method was used in later steps as well to calculate the RMSD. If the sequences deviated too much, this method did not deliver an optimal superimposition

of the structures. In that case, PyMOL's command "cealign" was used, which uses the CEalign algorithm [133] and also outputs the RMSD. Both results return slightly different RMSD values while giving visually a highly similar result. However, in the tested cases, this difference was below 0.1 Å and was therefore not visible in the listed results.

### 6.1.3 Classification into ordered or disordered regions

The predicted structure models from AlphaFold2 and their respective pLDDT values were used to classify each residue of NSP3 into ordered or disordered. The decision tree is depicted in Figure 5, while the method is explained in section 3.1.4.

In a more detailed and algorithmic way, the whole structure prediction is searched manually for secondary structure elements in PyMOL [131], with the cartoon representation enabled. More specifically, one looks for  $\alpha$ -helices and  $\beta$ -sheets. Starting with regions lacking those structural elements, the pLDDT of the respective residues is checked. If the pLDDT is below 50, it is considered as "definitely disordered". Is the pLDDT above 80, it is considered as "definitely ordered". If it is between 50 and 80, it is only considered "ordered" if it is surrounded by secondary structure elements, otherwise it is considered "disordered".

If a region contains secondary structure elements and a pLDDT above 80, it is considered as "definitely ordered". If the pLDDT is lower, and if the whole structure prediction contains a disordered loop at N-terminus or C-terminus, one must crop the respective disordered region from the sequence and resubmit it to AlphaFold2. After cropping, the new structure should have increased pLDDT values for ordered regions, while decreased values for disordered regions. In any case, if the region of interest contains  $\beta$ -sheets, it is considered as "very likely ordered". In case of  $\alpha$ -helices, there is still a chance that it is a false helix prediction. If possible, crop more disordered regions from termini and resubmit, otherwise consider that secondary structure element with low pLDDT as "probably disordered".

So-called "barbed-wire" conformations, as described by Williams et al. [132], were used as an additional hint to decide whether a region is disordered or not. Multiple iterations of predicting a structure were only applied when disordered regions had to be cropped off the sequence in order to improve the pLDDT values.

### 6.1.4 Structure Similarity Search

To identify structures with a similar fold, a structure similarity search was performed on the whole PDB [82] via the respective PDB search tool [124]. For the domains NAB and

$\beta$ SM, the PDB entries 7LGO and 7T9W were used as target model, respectively. For the domains  $\beta$ SLD, the folded ectodomain core, and Y1, the structure prediction was used. It was searched for "assemblies" and the return format was set to "assembly" as well. Both, "strict" and "relaxed", searches were formed for each case.

## 6.2 Experimental validation of the Betacoronavirus-specific linker domain ( $\beta$ SLD)

### 6.2.1 Structure prediction analysis

The structure was predicted via AlphaFold2 as described in section 6.1.2. The used sequence covers the residues 1057 to 1090 from SARS-CoV-2 NSP3 and the predicted structure was inspected manually and visually in PyMOL [131], where the "technical" preset was used to show all atoms and to display hydrogen bonds. The exact sequence is: TTIKPVTYKLDGVVCTEIDPKLDNYYKKDNSYFT

Ramachandran plots were generated with MolProbity [145].

### 6.2.2 Conservation of domain sequence

To analyse the conservation of the linker domain, its sequence was submitted to BLAST [76] with SARS-CoV-2 sequences excluded. In additional BLAST queries, the whole taxonomy of *Coronaviridae* was excluded to find other homologues.

For a more informative search, sequences of all betacoronaviruses found under the NCBI taxonomy id 694002 were used in local pairwise sequence alignments with the SARS-CoV-2 linker domain sequence. If available, sequences annotated as "NSP3" were used from the betacoronaviruses. Otherwise, the sequences of the whole polyprotein 1ab were used (which contains all non-structural proteins). In the latter case, it was checked that the alignment was located in the region of NSP3 on the polyprotein. Local pairwise sequence alignments were performed via EMBOSS Water version 6.6.0 [74]. Based on the alignment results, the sequences were extended to match the length of 34 residues of SARS-CoV-2 linker domain. These extended sequences were then used in a global multiple sequence via Clustal Omega version 2.1 [74] to identify conserved amino acids. Afterwards, the extended sequences were submitted to ColabFold [128] to predict structures via AlphaFold2 [15].

The whole procedure was repeated for other coronaviruses listed under the NCBI taxonomy id 694002. Specifically, this included viruses with the following NCBI accession

numbers:

24 genomes from *Alphacoronavirus*: NC\_046964, NC\_010437, NC\_022103, NC\_028811, NC\_028833, NC\_028814, NC\_028824, NC\_028752, NC\_002645, NC\_034972, NC\_002306, NC\_028806, NC\_038861, NC\_030292, NC\_005831, NC\_032730, NC\_010438, NC\_023760, NC\_003436, NC\_009988, NC\_018871, NC\_009657, NC\_048211, NC\_035191

17 genomes from *Betacoronavirus*: NC\_025217, NC\_038294, NC\_019843, NC\_039207, NC\_026011, NC\_003045, NC\_006213, NC\_006577, NC\_048217, NC\_012936, NC\_009020, NC\_017083, NC\_030886, NC\_009021, NC\_004718, NC\_045512, NC\_009019

5 genomes from *Gammacoronavirus*: NC\_010646, NC\_046965, NC\_048214, NC\_001451, NC\_010800

10 genomes from *Deltacoronavirus*: NC\_011547, NC\_016996, NC\_016993, NC\_011550, NC\_016994, NC\_039208, NC\_016992, NC\_011549, NC\_016991, NC\_016995

One genome from unclassified Shrew coronavirus: NC\_046955

Sequence alignments with these viruses were also conducted with the sequence of NAB and  $\beta$ SM-M from SARS-CoV-2, as both domains were also assumed to be specific to *Betacoronavirus* [18].

### 6.2.3 Expression of Betacoronavirus-specific linker domain construct

A construct comprising the domains Ubl2, PL2<sup>pro</sup>, and  $\beta$ SLD was designed to validate the presence of the folding linker domain experimentally. The construct consists of the 345 SARS-CoV-2 (NCBI reference YP\_009742610.1) NSP3 residues 746-1090 with the following amino acid sequence:

```
EVRTIKVFVTTVDNINLHTQVVDMSMTYGQQFGPTYLDGADVTKIKPHNSHEGKTFYVLPN
DDTLRVEAFEYYHTTDP SFLGRYMSALNHTKKWKYPQVNGLT SIKWADNNSYLATALLTL
QQIELKFNPPALQDAYRARAGEAANFCALILAYCNKTVGELGDVRETMSYLFQHANLDS
CKRVLNVVCKTCGQQQTTLKGVEAVMYMGTL SYEQFKKGVQIPCTCGKQATKYL VQQESP
FVMSAPPAQYELKHGFTFCASEYTGNYQCGHYKHITSKETLYCIDGALLTKSSEYKGPI
TDVFKENSYTTTTIKPVTYKLDGVVCTEIDPKLDNYYKKDNSYFT
```

*(The following paragraph was primarily written by David Briggs for a publication related to this work, who also conducted the described experiment.)*

The construct contains the C111S mutation, which is commonly used in crystallization of Ubl2+PL2<sup>pro</sup> [25]. It was amplified by PCR from a template kindly provided by David LV Bauer (Francis Crick Institute). The coding sequence was ligated into a

pGEX-6P vector using (5') BamHI and NotI (3') restriction sites introduced during amplification. Correct insertion was confirmed by Sanger sequencing. Recombinant protein was expressed in LB in BL21Gold E.coli as a GST-3C-fusion protein. After induction at OD600 of 0.6, the temperature was reduced to 16°C, and cells were harvested the following morning. Cell pellets were frozen at -80°C until needed.

#### 6.2.4 Purification of Betacoronavirus-specific linker domain construct

*(The following paragraph was primarily written by David Briggs for a publication related to this work, who also conducted the described experiment.)*

Cell pellets were resuspended in lysis buffer (50mM Tris pH 7.5, 300mM NaCl, 5 % (v/v) Glycerol, 0.5mM TCEP, 1µM Zinc Acetate. EDTA-free protease inhibitors (Roche) were added as per manufacturer's instructions. Cells were lysed by sonification, and lysate clarified by centrifugation at 45,000g, 4°C for 45 minutes. Protein was harvested by incubating the lysate with Glutathione Sepharose 4B (Cytiva) for 2 hours at 4°C with constant mixing. The beads were then harvested and washed with 10 bed-volumes of lysis buffer and then 20 bed-volumes of SEC buffer (20mM Tris, pH 7.5, 150mM NaCl, 0.5mM TCEP, 1µM Zinc Acetate). Beads were then resuspended in 5 bead volumes of wash buffer, and then incubated with GST-HRV3C protease overnight at 4°C with constant mixing. The cleaved product was collected from the beads via filtration. This material was then concentrated to 5mg/mL, aliquoted, snap frozen in liquid nitrogen and stored at -80°C until needed.

#### 6.2.5 Crystallization of Betacoronavirus-specific linker domain construct

*(The following paragraphs were primarily written by David Briggs for a publication related to this work, who also conducted the described experiment.)*

Prior to crystallization, the protein was thawed on ice and any remaining aggregates or impurities were removed via a final size exclusion chromatography step, using a Superdex200 increase column, equilibrated in SEC buffer (20mM Tris-HCl pH 7.5, 150mM NaCl, 1mM TCEP, 1µM zinc acetate). This material was diluted 1:2 with milliQ water and concentrated to 7.5mg/mL for crystallization trials. Sitting-drop vapour diffusion crystallisation experiments were setup in MRC 2-well 96 well plates using a Formulatrix NT-8 drop setting robot. Initial microcrystals we obtained using small (200nl protein + 100 mother liquor) drops, which were then used to streak seed into larger 400nl + 200nl drops. Crystals appeared after 4-5 weeks after incubation and streak seeding at 4°C in 0.1M MES pH 6.8, 8.4 % PEG20K. Crystals were cryocooled in liquid nitrogen using



crystallization liquor supplemented with 20 % Ethylene glycol.

Diffraction data were collected at beamline I24 at Diamond Light Source. Auto processing using the AutoProc pipeline [146] indicated that the data extended to 2.1 Å resolution, and that the crystals had space group P1. Molecular replacement using the MR-BUMP pipeline [147] in CCP4Cloud [148] gave a reasonable solution using PDB 4M0W (“Crystal Structure of SARS-CoV papain-like protease C112S mutant in complex with ubiquitin” [149]) with a TFZ of 12.0 and an initial R<sub>free</sub> of 47 %.

Inspection of the electron density map in COOT [150] revealed that the amino-terminal lobe of PL2<sup>pro</sup> was well defined, but the carboxy-terminal lobe was not, with poor electron density and noisy difference maps. The output from this molecular replacement solution was submitted to the Modelcraft auto-building and refinement pipeline [151] which finished with a final R<sub>free</sub> of 26.5 %. Inspection of the results revealed that the crystals in fact contained two copies of the N-terminal lobe (Arg3 to Cys181 (construct numbering, supplementary information section B)) of the PL2<sup>pro</sup> monomer, forming a close and compact dimer. No further refinement was undertaken.

### 6.2.6 SAXS analysis of Betacoronavirus-specific linker domain construct

*(The following paragraph was primarily written by Yunyun Gao for a publication related to this work, who also conducted the described experiment.)*

The purified protein was shipped on dry ice to P12 BioSAXS beamline [152] at PETRA III (DESY, Hamburg, Germany), where the solution small-angle X-ray scattering (SAXS) experiment was conducted. P12 provides monochromatic X-rays with a 0.2 x 0.05 mm<sup>2</sup> beam at the sample capillary. For this experiment, the X-ray wavelength and sample-to-detector distance were 0.124 nm and 3000 mm respectively. The purified monodisperse fraction was prepared as protein stocks with a concentration of 15 mg/mL. The stock buffer is 20 mM Tris (pH 8.5), 150 mM NaCl, 0.5 mM TCEP, 1 μM zinc acetate. The buffer for background subtraction was prepared by first diluting the protein stock 1:1 (v:v) with a close match of the stock buffer then collecting the flow-through of ultrafiltration when the diluted protein solution was concentrated to the original volume. The concentration series (0.65 mg/mL, 0.97 mg/mL, 1.29 mg/mL, 1.61 mg/mL, 1.94 mg/mL, 3.22 mg/mL) is prepared by dilution with the closely matched buffer and careful filtration with Nanosep 100K OMEGA filters (PALL Life Sciences). The concentration series was measured under the standard “batch mode” at P12. The scattering profile was radially integrated, averaged and absolutely scaled from the corresponding detector images of 30 exposures. Each exposure was 95 ms. The model and its theoretical SAXS curve are from

the most fitted output generated by SREFLEX [134], using the background subtracted SAXS profile of the 3.22 mg/mL sample and relaxed Alphafold2 prediction as the input. The ab-initio modeling was carried out as following: 1) generate 20 bead models with DAMMIF [153] using the GNOM [154] output of the 3.22 mg/mL sample; 2) generate a starting search volume by averaging the 20 bead models using DAMAVER [155]; 3) run DAMMIN [156] using the starting search volume and the GNOM output of the 3.22 mg/mL sample.

### 6.3 Hexameric pore complex

*(From here, all following sections and paragraphs were written again by myself, Maximilian Edich, if not stated otherwise.)*

#### 6.3.1 Prediction of multidomain segments

To assess the validity of predicted arrangement between domains, input sequences for AlphaFold2 were prepared, which covered multiple subsequent domains from SARS-CoV-2 NSP3. Due to length limitations for input in early ColabFold versions, these sequences covered only a few domains. A full length NSP3 prediction, however, was conducted towards the end of this work. Since the first domain, Ubl1, is immediately followed by the large disordered domain HVR, the first multidomain input sequence was designed to cover Mac1, Mac2, and the linker in between. Next came the sequences Mac2 to DPUP, Mac1 to DPUP (including the previous two, to explore differences), DPUP to PL2<sup>pro</sup>, Ubl2 to NAB, NAB with the complete  $\beta$ SM, and Y1 to CoV-Y. Since AlphaFold2 does not take membranes into account, the transmembrane region covering TM1 to AH1 was handled separately and no sequences were designed to include these domains, except for the full length NSP3 prediction.

The resulting pAE matrices were utilized to evaluate the arrangement of multiple domains, where rank1 predictions were prioritized. Other ranks were taken into account to observe whether the arrangements of the structures were predicted consistently similar or if they deviated from each prediction.

An assembly of two domains was considered plausible if the respective pAE values were consistently below ten. A custom python script was used to calculate statistics of pAE values and assemblies with matrices having average pAE values above 15 were considered implausible.

### 6.3.2 Conservation of Y1

To investigate the conservation of Y1, a BLAST search [76] and sequence alignments were conducted as already described for the *Betacoronavirus*-specific linker domain in section 6.2.2. For sequence alignments, the same betacoronaviruses as in section 6.2.2 were used. For conservation within *Coronaviridae*, one or two representatives per genus per used. These were NP\_073549.1 from *Alphacoronavirus*, NP\_066134.1 and YP\_001876435.1 from *Gammacoronavirus*, YP\_002308496.1 and YP\_002308505.1 from *Deltacoronavirus*, and NC\_046955.1 from the unclassified Shrew coronavirus.

Global sequence alignments with SARS-CoV-2 Y1 and the polyprotein of the respective viruses were performed via Clustal Omega [74]. Sequences from these alignments were considered as potential Y1 homologs and were submitted to ColabFold [128] to obtain structure predictions from AlphaFold2 [15]. Alignments to the structure prediction of SARS-CoV-2 Y1 were performed in PyMOL [131], where RMSD values were calculated as well.

For conservation among *Nidovirales*, a similar procedure was conducted. Since global sequence alignments were worse, sequences from the alignment were iteratively extended and resubmitted to ColabFold depending on each result individually. A focus for extending the sequence were the alignment itself and the predicted secondary structure elements, as well as the pLDDT scores. If a predicted segment at a terminus seemed to be cropped, the respective terminus was extended and the structure repredicted. The examined viruses sequences and their NCBI codes were NC\_076697.1 from *Pitovirinae*, MZ203498.1 from *Letovirinae*, YP\_008798231.1 and YP\_009665195.1 from *Torovirinae*, YP\_001661453.1 from *Roniviridae*, NC\_015668.1 from *Mesoniviridae*, and NP\_127506.1 from *Arteriviridae*.

In a BLAST search [76], the taxonomy *Coronaviridae* was excluded with the taxonomy id 11118.

### 6.3.3 Multimer predictions of Y1 and Y1+CoV-Y

Prediction of all multimers was accomplished in ColabFold [128] via the multimer feature of AlphaFold2 [15, 103]. Because the prediction did not terminate and save the results on default settings due to runtime limitations in the free of charge version of ColabFold [128], the code was adapted to terminate after the prediction of a single multimer by setting "num\_models" to 1 and "model\_order" to "[1]". To evaluate the pAE matrices, a custom python script was developed, which calculates statistics from the whole pAE matrix and

individual cells of the multimer pAE matrix.

Electrostatics were calculated in PyMOL [131] via "generate vacuum electrostatics". UCSF Chimera [66] was utilized for fitting the predicted hexamer of MHV Y1+CoV-Y and Y1 into the cryo electron tomography density map from Wolff et al. [19]. To accomplish a decent fit, the hexamer was positioned close to the membrane region of the complex and oriented in the correct way according to the fact that the transmembrane region is at the N-terminus of Y1. Rotation and shift were enabled.

#### 6.3.4 Expression of Y1

*(The following paragraph was primarily written by David Briggs for a publication related to this work.)*

The construct of Y1 (supplementary information section C) was amplified by PCR from a template kindly provided by David LV Bauer (Francis Crick Institute). The coding sequence was ligated into a pGEX-6P vector using (5') BamHI and NotI (3') restriction sites introduced during amplification. Correct insertion was confirmed by Sanger sequencing. Recombinant protein was expressed in LB in BL21Gold E.coli as a GST-3C-fusion protein. After induction at OD600 of 0.6, the temperature was reduced to 16°C, and cells were harvested the following morning. Cell pellets were frozen at -80°C until needed.

Cell pellets were resuspended in lysis buffer (50mM Tris pH 7.5, 300mM NaCl, 5 % (v/v) Glycerol, 0.5mM TCEP, 1µM Zinc Acetate). EDTA-free protease inhibitors (Roche) were added as per manufacturer's instructions. Cells were lysed by sonification, and lysate clarified by centrifugation at 45,000g, 4°C for 45 minutes. Protein was harvested by incubating the lysate with Glutathione Sepharose 4B (Cytiva) for 2 hours at 4°C with constant mixing. The beads were then harvested and washed with 10 bed-volumes of lysis buffer and then 20 bed-volumes of SEC buffer (20mM Tris, pH 7.5, 150mM NaCl, 0.5mM TCEP, 1µM Zinc Acetate). Beads were then resuspended in 5 bead volumes of wash buffer, and then incubated with GST-HRV3C protease overnight at 4°C with constant mixing. Analysis by SDS-PAGE of the purification process showed that no detectable Y1 had been expressed.

#### 6.3.5 Integrative modelling

*(From here, all following sections and paragraphs were written again by myself, Maximilian Edich, if not stated otherwise.)*

The structures of the domains from MHV NSP3 predicted by AlphaFold2 [15] were fitted into the cryo electron tomography map from Wolff et al. [19] via Assemblin [119].

Assembleline was installed according to the official installation instruction pages. Two subunits were used, "nsp3" and "nsp3\_Cterm", where "nsp3" covered the MHV NSP3 domains Ubl1, PL1<sup>Pro</sup>, Mac1, DPUP-like, Ubl2, PL2<sup>Pro</sup> with attached  $\beta$ SLD, NAB, and  $\beta$ SM. The subunit "nsp3\_Cterm" was defined separately to fit the prediction of the Y1+CoV-Y hexamer, because no symmetry copies had to be generated for this structure.

First, the Assembleline steps of generating fit libraries and global optimization were run only for Ubl1 domain and in a separate run only for the C-terminus subunit in order to validate via experimental results or to compare to the previous fit of the hexamer into the cryoTM map, where the RMSD between both fitted hexamers (first fit directly in Chimera described in section 6.3.3 and second fit via Assembleline) was calculated in UCSF Chimera [66]. The resulting fits were examined manually and then used in a third run, where all domains were used. The global optimization step is usually followed by the refinement step, but due to unsatisfying results, this step was not executed. All parameters were identical between runs except for domain specific paths and sequence ranges, and for the disabled generation of symmetry copies for the hexameric C-terminal structure.

For the `efitter_params.py` file, the parameters "CA\_only" and "move\_to\_center" were set to "False", while "backbone\_only" was set to "True". Per structure, 100000 placements were probed. Afterwards, p-values were calculated via "genpval.py".

For global optimization, the files "MHV\_pore\_complex.json", "xla\_project.json", and "params.py" were created, where all files are based on the official templates. To obtain domain specific parameters, the Xlink Analyzer [120] plugin for UCSF Chimera [66] was used according to the instructions of the Assembleline [119] manual pages.

For all domains, a series block for sixfold symmetry was added, where the mode was set to "auto" for all structures of subunit 1 and to "input" for the already hexameric structure of C-terminal subunit. To define the symmetry axis, the cryoTM was loaded into UCSF Chimera, where the command "measure symmetry #0" led to the first symmetry information. For obtaining the axis point required by Assembleline, a second copy of the map was loaded to Chimera and rotation was deactivated for the first map. The second map was rotated manually by roughly 60° and fitted via *Tools > Volume Data > Fit in Map* to the first map. Both maps were activated via *Tools > General Controls > IDLE* and the command "measure rotation #1 #0" generated the output containing the axis point, which was used with the type "C6" and axis parameter "[0, 0, 1]" to define the symmetry block of the JSON file.

In the data block of the JSON file, domain name, respective subunit, path to PDB,

and path to the generated fit library were written for each domain. The parameters "foreach\_serie" and "foreach\_copy" were set to "True" in all cases. In the subunits block, chains were set from "A" to "F" for subunit 1 and from "G" to "K" for the C-terminal subunit. Domain ranges were also set here, where for subunit 1 the full MHV NSP3 sequence was given in a FASTA file, while for the C-terminal subunit only the sequence of the structure was given, which required the range to be set to "1, 342". In order to make the Assembline run successful, the chains and residue numbers must be adapted to the parameters in the JSON file. This was accomplished in PyMOL [131] via the commands "alter sele, resi=int(resi)+X" to shift the residue numbering by X and with "alter (chain A),chain='G'" to rename chain labels.

The second JSON file, "xla\_project.json", contained similar information, which could be completely derived from the previous steps. The "params.py" file based on the official template, but because no information from cross-linking experiments was available, "xlink restraints" was removed from the high-res scoring function, while symmetry restraints were added.

The global optimization was run on a single linux machine. Thus, the code for running Assembline on computer clusters was removed. Due to limited computational resources, only a single trajectory was run via "--prefix 0000000". After the global optimization succeeded, the current working directory was changed to the folder "out", from where scores were extracted and plots generated. This was done via these three commands:

```
extract_scores.py
plot_convergence.R total_score_logs.txt 1
plot_scores.R all_scores.csv
```

Models were created via the command

```
rebuild_atomic.py --top 1 --project_dir \
../ MHV_pore_complex.json all_scores_sorted_uniq.csv \
--rmf_auto
```

Unfortunately, this crashed for all models except for the hexameric Y1+CoV-Y structure due to an internal bug of Assembline, which remained unsuccessful. Hence, only that structure was generated as an atomic model, while the other domains could only be loaded as coarse grained models into UCSF Chimera [66].

## **Part II**

### **Beyond the Surface: Exploring the Solvent Region of Macromolecular Crystals with Experimental Phasing**

## 7 Introduction and Objectives

### 7.1 Introduction

Protein structures are modelled and refined with the corresponding electron density map of a macromolecular crystal. Such a map is not measured directly, but is instead calculated from structure factors describing the electron density of a crystal's content. However, the signal collected from the crystal via X-rays relates only to the structure factor amplitudes, while the phase information is not directly measurable.

One common way of overcoming this phase problem is to use the structure model of a similar protein or from structure prediction to estimate phases, which are then combined with the experimentally obtained structure factor amplitudes. While this technique leads indeed to a structure model of the target protein, it may introduce a bias by using other data than the experimentally measured one for the phases. Furthermore, information about the solvent region is usually weighted down via "solvent flattening", as no solvent model is provided. Alternative methods are capable of deriving the phases experimentally from the collected data, leading to electron density maps free from model bias and solvent flattening, but require a more complicated experimental setup. Although protein structures are rarely solved by such techniques, their data could provide new insights about the solvent region and the discrepancy between experimental data and structure models, which is present for all macromolecular crystals. Moreover, it could provide new insights into the R-factor gap problem, where the solvent region is potential contributor to the discrepancy between our structure models and the observed experimental data [157].

In this work, different methods of obtaining phase information and density modification were compared with a focus on the solvent region. Contribution of the solvent region to the structure factors was analysed and the new information allowed to lower the R-values of a PDB structure by building alternative water models.

### 7.2 Objectives

The two major questions with current methods of phasing via molecular replacement and solvent flattening are: 1) Has the introduced model bias a significant impact and 2) Does the solvent region contain valuable information which is lost when applying solvent flattening? To assess these questions, models built and refined with density modification were compared to purely experimentally phased models. Since this required a significant number of structure examples, a semi-automated pipeline had to be developed, which calculates electron density free of model bias and solvent flattening from publicly avail-



able datasets. Furthermore, the resulting electron density maps had to be evaluated for good quality and statistically analysed with additional software developed for this purpose. Specifically, the PDB had to be filtered for high resolution protein structures associated with data sets from MAD-phasing (multiple-wavelength anomalous diffraction phasing). Developing a manual method for calculating the desired electron density maps via SHELXE [158] and SHARP [159] was necessary before automating most of this process. Finally, the gained insights were tested to lower the R-values of PDB structures by refining the model against purely experimentally determined electron density maps without density modification.

## 8 Theoretical Background

### 8.1 Macromolecular Crystallography

Macromolecular crystallography is currently the most commonly used method for the solution of protein structures. While it has certain limitations (which were explored in section 2.2.2 of Part I), it provides high resolution data, from which atom precise structure models can be obtained. In this method, X-rays are directed towards protein crystals, where the X-ray photons interact with the electrons, resulting in scattering of photons and a measurable diffraction pattern. The intensities of these diffraction patterns are proportional to the amplitudes of structure factors, which describe the electron density of the molecule. The phases of these structure factors, however, are not measured, but they can be obtained from the measurable anomalous dispersion, which is explored in the next section. After obtaining structure factors with amplitudes and phases, an electron density map is generated, which is used for model building. From the built model, calculated structure factors ( $F_{\text{calc}}$ ) can be obtained and compared to the observed structure factors ( $F_{\text{obs}}$ ). To assess how close the modelled structure is to the experimentally measured data, the R factor is calculated:

$$R_{\text{work}} = \frac{\sum ||F_{\text{obs}}| - |F_{\text{calc}}||}{\sum |F_{\text{obs}}|} \quad (1)$$

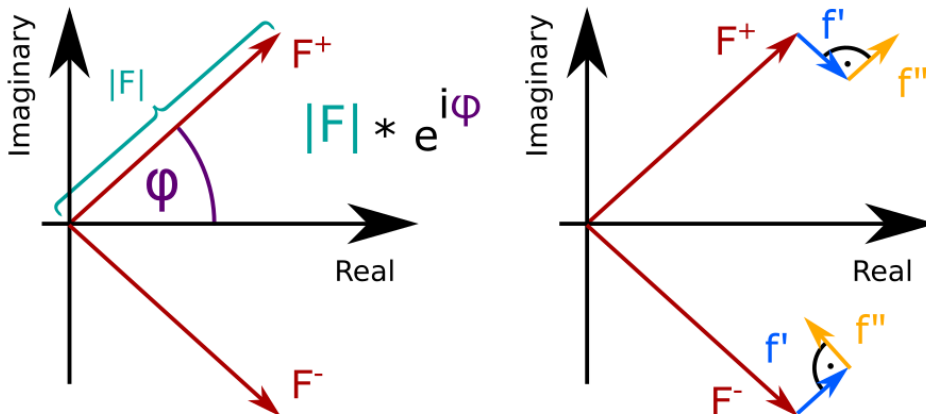
In addition to the  $R_{\text{work}}$ , the  $R_{\text{free}}$  is calculated, where usually 5 % of all reflections are excluded from the refinement to test how well the model predicts unseen data.

### 8.2 Experimental phasing

During a beamline experiment, only amplitudes are measured, while phases are not experimentally obtained. One possibility to calculate those phases is by utilizing anomalous differences in experimental phasing. Alternatively, molecular replacement can be used to obtain phases from an already solved similar structure. Several methods can be described under the term of experimental phasing. They all have in common that they maximize the anomalous differences during data collection, which are later utilized to determine the phases of an atomic substructure.

#### 8.2.1 Anomalous differences

Friedel's law states that



**Figure 24:** Scattering factor  $F^+$  and its centrosymmetrically opposed scattering factor  $F^-$  represented as vectors in the Argand diagram (left). The anomalous contributions  $f'$  and  $f''$  are added to the normal scattering factors  $F^+$  and  $F^-$ , where  $f'$  is mirrored along the real axis, while  $f''$  is mirrored along the imaginary axis (right).

$$|F_{hkl}| = |F_{-h-k-l}| \quad (2)$$

where centrosymmetrically opposed structure factors have the same amplitude. However, this is only true in the absence of anomalous differences [160]. These result from anomalous contributions to the scattering factor of each atom, which add a wavelength dependent complex term to the scattering factor equation:

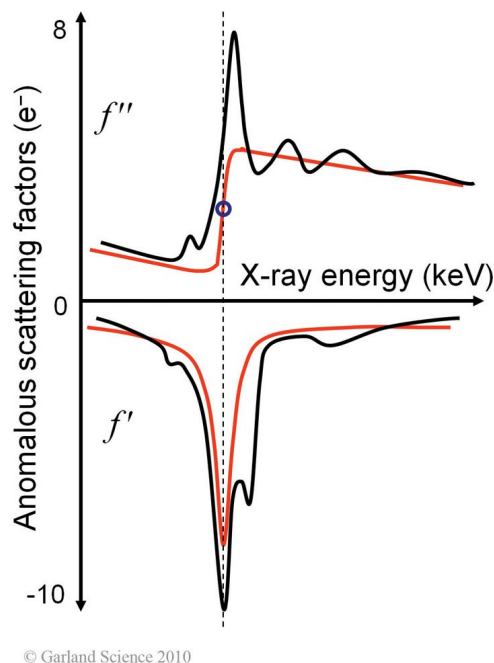
$$f_{ATOM} = f + f'_{(\lambda)} + i * f''_{(\lambda)} \quad (3)$$

where  $f$  is the normal scattering factor. When representing a scattering factor as a vector, it becomes obvious where the anomalous differences come from: While the real contribution  $f'$  is symmetric along the real axis for a pair of centrosymmetrically opposed scattering factors, the complex contribution  $f''$  is symmetric along the imaginary plane and has thus a phase shift of  $90^\circ$  from  $f'$  (Figure 24).

Furthermore,  $f'$  and  $f''$  are wavelength and element dependent (Figure 25), which is the fundamental property that is exploited for experimental phasing.

### 8.2.2 MAD phasing

Several methods allow experimental determination of phases, where MAD (multiple wavelength anomalous diffraction) phasing is one of them. Its key concept is that the macromolecular crystal is measured at three or four different wavelengths, where each is selected



**Figure 25:** Figure is taken from Rupp, *Biomolecular Crystallography* [160]. The plot shows the wavelength dependency of  $f'$  and  $f''$ , where the black curve shows the fine structure of the functions and the red curves show simplified functions.

carefully depending on X-ray absorption scans of the sample or theoretical absorption curves for the anomalous scatterers, which are heaviest element in the unit cell. The wavelengths are either chosen as peak (maximal difference between  $f'$  and  $f''$ ), inflection (minimal  $f'$  values), high energy remote (high  $f''$ ), or low energy remote (low  $f''$ ), according to the anomalous scattering curves of  $f'$  and  $f''$  (Figure 25). From the measurement at different wavelengths,  $f'$  and  $f''$  can be determined, which finally gives phase estimates. These are first made for the substructure, which are the anomalous scatterers of the unit cell. In protein crystals, these are heavy atoms such as zinc, copper, or iron, which are ions bound to the protein, or sulfur from methionines and cysteines, where selen as a strong anomalous scatterer is also possible by incorporating selenomethionine [160]. After solving the substructure via programs such as ANODE [161] or SHELXD [162], the phases of the remaining structure are calculated via programs such as SHELXE [158] or SHARP [159].

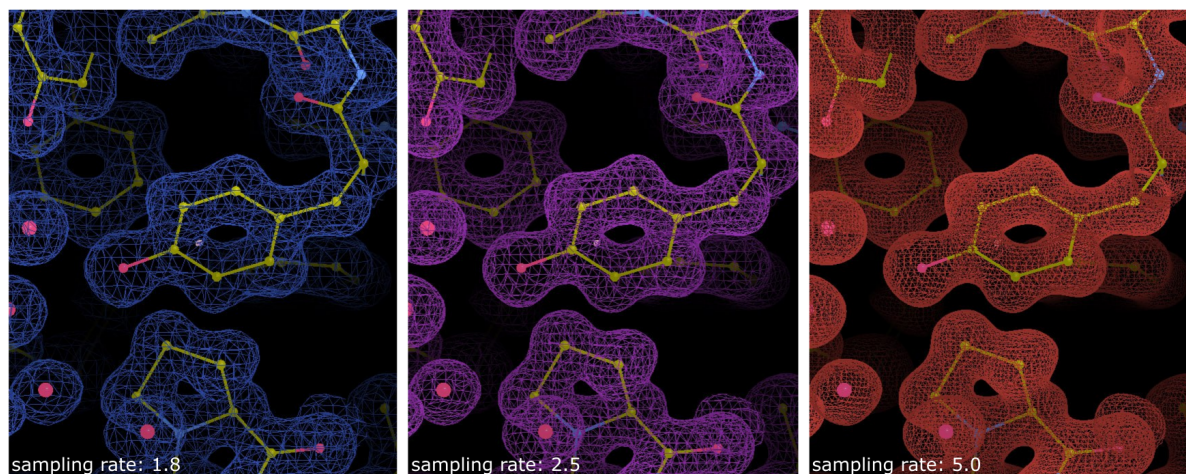
### 8.3 Crystallographic data

Crystallographic experiments collect data in the form of reflections, which can be related to structure factors describing the electron density of a unit cell averaged over all unit cells of a crystal. The geometry of the reflections and experimental setup relates to the miller indices of each reflection, while the measured intensity is proportional to a structure

factor's amplitude. Phases can be estimated and calculated in various ways, as explored before. The interesting data for the analysis of this work are the electron density maps, consisting of grid points with values describing the density, and structure factor tables, which list structure factors by Miller indices and assign next to phase and amplitude various other information. Knowing the format of these data enables understanding of how these can be manipulated or analysed.

### 8.3.1 Electron density maps

Most of the time, electron density maps are generated from structure factors. The Fourier transform generates from these and with the unit cell parameters the density map. The map itself consists of a regular grid of 4-dimensional points, which contain x, y, and z coordinates, as well as a value describing the electron density at that point. The number of points are dependent of the resolution of the data. Since the electron density is calculated as a continuous function, however, the grid points are just discrete samples and the respective sampling rate (Figure 26), which defines the exact number of grid points and their spacing, is set in the program for map visualization, such as COOT [150]. Because this can result in millions of grid points for the entire unit cell calculated from tens of thousands of structure factors, it is cheaper to store the electron density information in form of structure factor tables rather than in .MAP files.



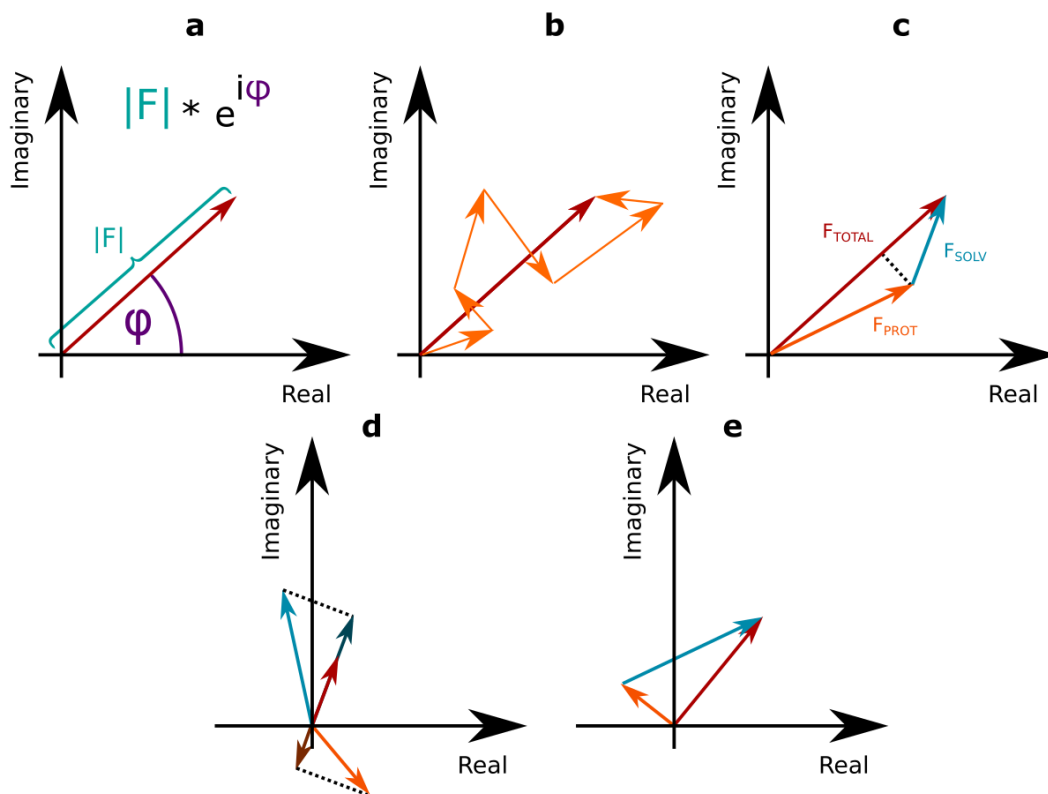
**Figure 26:** Comparison of different sampling rates while loading the same set of amplitudes and phases from the PDB structure 3SEE. The sampling rate of 1.8 (left) is the default setting. Contour level is in all cases at 1.0 RMSD.

To make the 3-dimensional electron density visible, a contour level is defined, which visualizes only the peak values above a certain threshold. Grid point values can be both, positive and negative, hence peaks and valleys exist, but both are usually referred to as positive or negative peaks, respectively. Only the grid points that are above the threshold

and adjacent to grid points below the threshold, i.e. those at the edge of a visible peak, are connected to each other, creating a mesh (Figure 26). To compare electron density maps visually, they must have the same contour level, which is measured in electrons per  $\text{\AA}^3$  or by the RMSD, where higher RMSD or number of electrons per  $\text{\AA}^3$  increases the threshold and shows thus less density.

When building a molecular model into the observed electron density, atoms are placed in real space. From these, one can calculate the structure factors  $F_{\text{calc}}$  which describe the electron density of the molecules. The map from observed structure factors,  $F_{\text{obs}}$ , and  $F_{\text{calc}}$  can be subtracted from each other to obtain a so-called difference map, such as an Fo-Fc map, which provides a more convenient visualization for building and refining a structure model [163]. Ideally, the model is built in such a way, that the difference map has neither strong positive peaks, nor strong negative peaks, as it indicates missing atoms or atoms not supported by electron density. Next to model building, molecular dynamics simulations are another use case where atoms are placed or moved in real space with chemical or physical restraints in order to simulate movement of molecular structures or the distribution of molecules in the solvent region. Calculating the electron density map of these can also provide new insights.

Last but not least, the grid points of a map can be modified. With a molecular model, it is possible to mask certain regions of a map to, for example, obtain only grid points in a defined radius around each atom of the structure. Hereby, the values of all grid points within the selection are set to zero, while all others remain unmodified. The inverse case is of course also possible. While this is a useful for map analysis, as later seen in this work, more practical techniques for the modification of electron density are applied during structure solution. The most interesting here is the solvent flattening, where a flat map for the solvent channel region is expected to be more likely correct, than the purely experimentally based density with high variations [164]. In real-space solvent flattening, the density of that region is lowered and new phases are calculated from the modified map, which are afterwards combined with the experimentally determined phases [164], so that not all of the experimental information is lost. Practically, solvent flattening is applied in reciprocal space, where a likelihood function is maximized for identification of optimal weights for the experimental phases and phases from the modified map. Finding these optimal weights, however, is difficult since both phases are not independent from each other [164–166].



**Figure 27:** Argand diagram visualizing complex numbers in a 2D plane, where structure factors can be represented as vectors with the structure factor amplitude defining the vector's length and the phase defining its angle between the real axis (a). The total structure factor of a unit cell ( $F_{TOTAL}$ ) can be represented as sum of structure factors describing separate components of the unit cell (b). With two components, the respective structure factor vectors can be projected onto  $F_{TOTAL}$  to calculate their contribution (c). The projections can exceed the length of  $F_{TOTAL}$  (d) and are dependent on phase differences to  $F_{TOTAL}$ , which can result in a projection of zero length (e).

### 8.3.2 Structure factor tables and contribution of separate components

Structure factors are stored in tables, most commonly in the .MTZ file format, but also often in the file formats .HKL, .SCA, or .PHS. MTZ files store additional information in a header and can keep track of modifications made to the file. These tables contain at least the Miller indices and intensities or amplitudes, usually a sigma value for the intensities or amplitudes, and often the phases. Additional common columns are  $R_{free}$ -flags, and figures of merit. Which columns are present depends on the programs generating these files and which phasing method was used. For anomalous phases, intensities and amplitudes are present for the positive and negative indexed reflection, from which anomalous differences could be calculated (see previous section 8.2.1).

It is possible to represent structure factors as complex numbers and thus as vectors in the Argand diagram (Figure 27a). The relevant information are the amplitude, which equals the magnitude of the structure factor, and the phase angle  $\phi$  in radians:

$$|F| e^{i\phi} \quad (4)$$

The structure factor for a certain Miller index  $hkl$  of the whole unit cell ( $F_{\text{TOTAL}}$ ) is technically the sum of all structure factors with the same index describing a single atom in the unit cell (Figure 27b). By separating the electron density map into two regions, such as the density of the protein structure and the non-protein density, which is the solvent region, one can calculate the contribution of each partial structure factor to the total structure factor, by projecting the vectors onto  $F_{\text{TOTAL}}$  (Figure 27c). Projections are calculate by

$$V_{\text{SOLV}} = \frac{F_{\text{SOLV}} \overline{F_{\text{TOTAL}}}}{|F_{\text{TOTAL}}|^2} F_{\text{TOTAL}} \quad (5)$$

where  $V_{\text{SOLV}}$  is the projected vector. The projection can then be used to calculate the contribution of  $F_{\text{SOLV}}$  to  $F_{\text{TOTAL}}$  by

$$\text{contribution}_{\text{SOLV}} = \frac{|V_{\text{SOLV}}|}{|V_{\text{SOLV}}| + |V_{\text{PROT}}|} \quad (6)$$

However, such a projection can come with problems. First, the projections can exceed the length of  $F_{\text{TOTAL}}$  (Figure 27d). Second, if the phase difference between  $F_{\text{TOTAL}}$  and  $F_{\text{SOLV}}$  is 90, the contribution from  $F_{\text{SOLV}}$  to  $F_{\text{TOTAL}}$  is zero, independent of the amplitude (27e). An alternative to this phase dependent structure factor contribution is simply diving the amplitudes of  $F_{\text{SOLV}}$  by the amplitude of  $F_{\text{TOTAL}}$ . Since the amplitude of  $F_{\text{SOLV}}$  can be zero, the better approach would be to divide by the sum of the components:

$$\text{contribution}_{\text{SOLV}} = \frac{|F_{\text{SOLV}}|}{|F_{\text{SOLV}}| + |F_{\text{PROT}}|} \quad (7)$$



## 9 Results

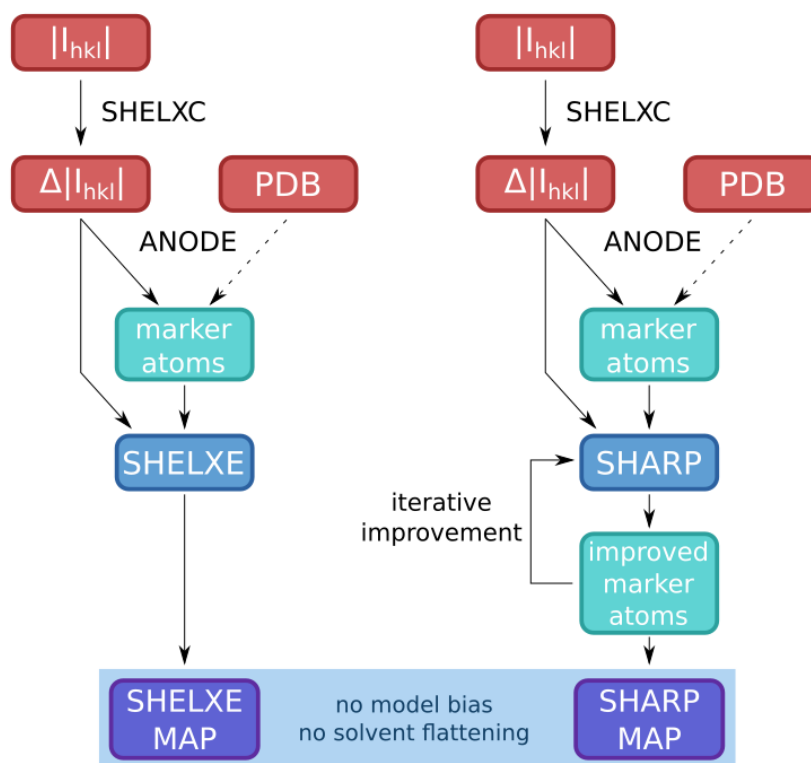
Crystals of proteins contain in each unit cell not only the protein structure, but also large solvent channels. Experimental phasing leads to electron density maps, which show strong peaks in the solvent region, which are usually flattened down via density modification. Here, these regions of the electron density were investigated by generating density maps via different methods and comparing these. Next to this analysis in real space, the contribution of the solvent region to structure factors was examined in reciprocal space.

### 9.1 MAD phasing

Different phasing methods lead to distinct phases for measured reflections and lead thus to visible differences in electron density. In order to evaluate such differences numerically in the solvent region of a unit cell, the RCSB PDB [82] was filtered for datasets from MAD experiments which fulfilled specific criteria. These datasets, containing measured intensities for multiple wavelengths from which anomalous differences could be calculated, were automatically phased within a custom pipeline, where the reflection data was handled via ANODE [161], and SHELX [158, 162, 167], as seen in Figure 28. Automated evaluation by cross-correlation and analysis of average anomalous density peaks provided an overview of each electron density map's quality together with various metrics, which were used to select few candidates for manual phasing with SHARP [159]. To handle the large number of entries, custom python scripts were written to automate most of the processes.

#### 9.1.1 Experimental phasing via SHELXE

A first search from the PDB [82] provided a list of 573 entries associated with MAD data (listed in Appendix A.1). Filtering out entries on various criteria (described in detail in section 12.1.2) left 105 entries, which had datasets measured at three wavelengths to resolutions below 2 Å and a completeness of 97 % or higher. Furthermore, these entries were lacking any potential anomalous scatterers as additives in the crystallization condition and all data per entry was collected from a single crystal. From the entries which were filtered out, 15 were flagged manually due to different reasons listed in section 12.1.2. Among these were, for example, identical values for two wavelengths, incomplete data, and unusual CIF file formats, which were incompatible with subsequent scripts of the pipeline.



**Figure 28:** Schematic workflow of generating electron density maps with SHELXE and SHARP. Intensity differences per reflection per wavelength are calculated from measured intensities via SHELXC. ANODE uses these anomalous differences to calculate the phases and thus the electron density map of the substructure consisting of marker atoms, i.e. anomalous scatterers. The PDB file is only used to provide reference points within the unit cell for manual inspection and does not contribute the calculation of the substructure. Intensities of the entire structure and the substructure are loaded either into SHELXE or into SHARP. While SHELXE calculates the remaining phases of the entire structure, SHARP provides an improved substructure first, which can be iteratively improved. Both methods lead to an experimentally phased electron density map without any contribution of a previous model and without the application of any solvent flattening.

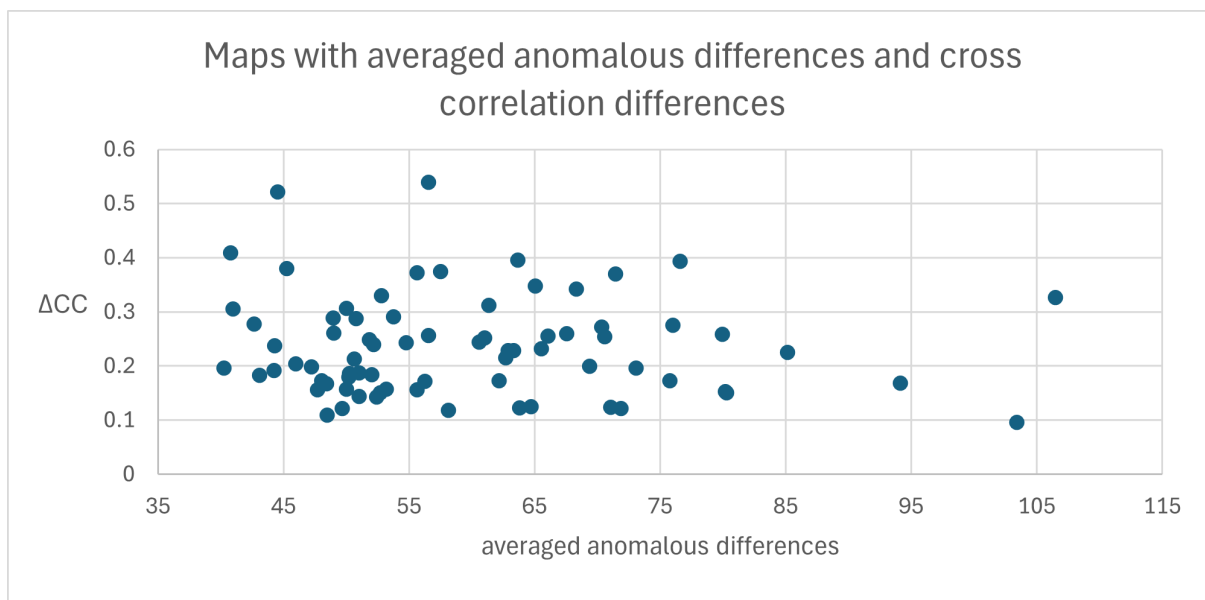
The remaining datasets were automatically converted into different file formats and phased via SHELXE [158] by first finding a substructure of anomalous scatterers via ANODE [161], which took anomalous differences calculated by SHELXC (Figure 28 left). SHELXE was executed with disabled 'free lunch' algorithm and with disabled density modification. The PDB structure was only used to provide atoms and their elements as reference points to evaluate the anomalous difference peaks from the text file alone. However, the automatic pipeline did not consider the elements and only focused on the height of the peaks. Therefore, the resulting electron density maps were free of model bias, solvent flattening and density modification. Because low averaged anomalous densities (AAD) lead to maps of low quality, such structures were filtered out as well, leaving only 74 structures for the further analysis.

To provide a quality indicator for faster evaluation of the electron density maps, the cross correlation between the PDB structure and electron density map from SHELXE was calculated. Since this value varies per structure, the cross correlation between the PDB structure and the electron density map deposited to the PDB was also calculated. The difference between both cross correlation values was then used as an indicator for low quality together with the AAD values calculated from ANODE.

For the 74 structures, the cross correlation values between SHELXE map and the PDB model range from 0.38 to 0.82. Between the deposited map and the model, the values range from 0.87 to 0.95. Only two electron density maps from SHELXE showed AAD values above 90 and cross correlation differences below 0.2 (Figure 29, which showed very high quality. The majority of maps showed cross correlation differences below 0.3, which looked at least decent on manual inspection. Table 8 lists 14 structures, for which good or excellent maps were calculated with the here presented method. Additionally, three examples for low quality maps are listed as well. While the AAD values of these show high values, the cross correlation differences show strong deviations between the deposited map and the map experimentally phased with SHELXE.

### 9.1.2 Experimental phasing via SHARP

From the 74 structures, 14 were selected for phasing via SHARP, which are listed in Table 9. Among them are also two additional data sets (AcNIR) from collaborators, which were kindly provided by Mike Hough and Armin Wagner. Unfortunately, one of these cases and four others from the PDB did not lead to satisfying results, as seen in Figure 30c. While substructure identification with ANODE [161] was successful (Figure 30a), calculating the correct phases for the rest of the structure failed and ripples in the electron density



**Figure 29:** Scatter plot showing the averaged anomalous density values and cross correlation differences ( $\Delta CC$ ) of electron density maps calculated from SHELXE for 74 structures.

**Table 8:** List of entries used for manual phasing in SHARP [159] and three entries with low quality electron density map calculated by SHELXE (highlighted in red). Shown are the averaged anomalous densities (AAD) from ANODE [161], local cross correlation between the PDB structure and the the map from the PDB (CC PDB) or the map calculated from SHELXE [158] (CC SHELXE), the difference between these cross correlation values, and the resolution of the data.

PDB Code	AAD	CC SHELXE	CC PDB	CC difference	Resolution in Å
2QPX	71.05	0.815	0.939	0.124	1.40
3SEE	103.42	0.806	0.902	0.096	1.25
2R01	75.77	0.731	0.904	0.173	1.15
3UE2	51.82	0.649	0.898	0.249	1.23
3NO2	79.93	0.665	0.924	0.259	1.35
2R0X	70.60	0.633	0.887	0.254	1.06
3POH	71.89	0.820	0.941	0.121	1.55
3CJM	73.07	0.740	0.936	0.196	1.50
3N6Z	76.01	0.649	0.924	0.275	1.30
3OHG	50.63	0.734	0.947	0.213	1.80
3LLX	80.19	0.776	0.929	0.153	1.50
2QL8	94.11	0.759	0.927	0.168	1.50
3QC0	80.27	0.785	0.935	0.150	1.45
3NOH	58.14	0.789	0.907	0.118	1.60
2FUP	44.51	0.377	0.899	0.522	1.48
3DCZ	55.64	0.557	0.929	0.372	1.65
4JM1	71.45	0.540	0.910	0.370	1.40

peaks are observable around the heavy atoms (Figure 30b). The map from SHARP in these cases looks at first like it is matching with the structure, especially with oxygen and nitrogen atoms. However, the density structure branches off into multiple chains and the pattern continues into the solvent region (Figure 30c). Another data set of AcNIR at

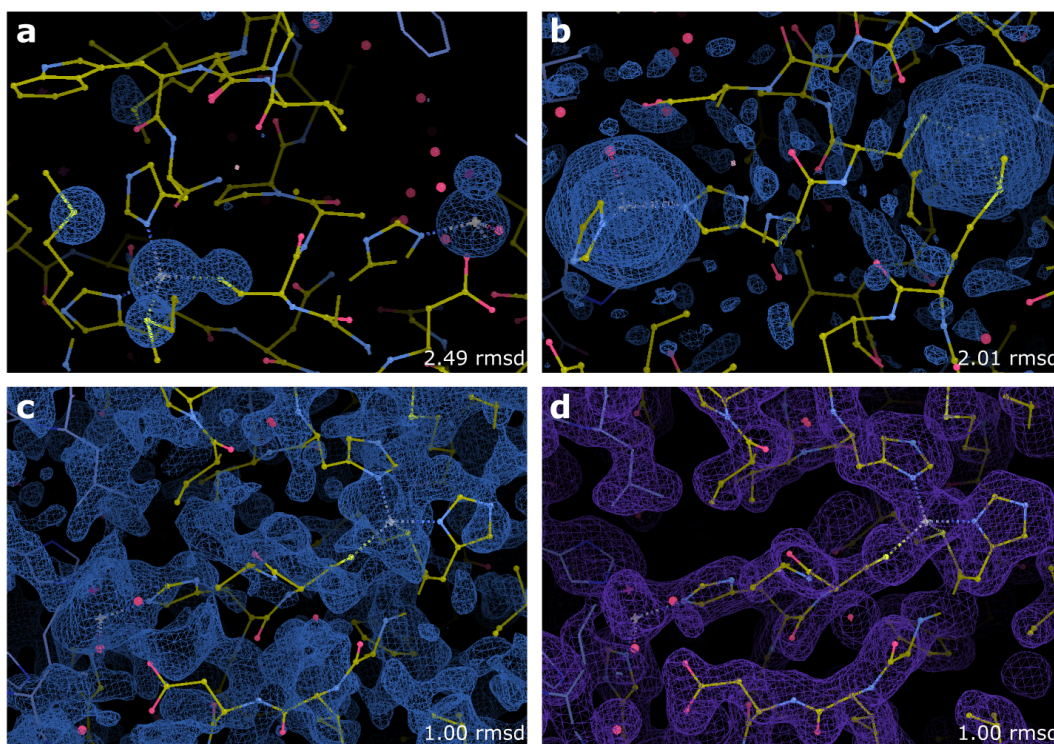
lower resolution did lead to an electron density map that agreed with the PDB structure 5N8I (Figure 30d). This leaves 11 electron density maps for further analysis, which are experimentally phased via SHARP without solvent flattening and without any influence of a pre-existing model which could have biased the map reconstruction.

**Table 9:** *List of structures which were phased with SHARP [159]. Entries are sorted by the crystal solvent content and entries shaded in red were not successfully phased.*

PDB Code	Resolution in Å	Solvent Content in %	PDB Deposition Date	Anomalous Scatterer
3OHG	1.80	68.63	2010-08-17	Se
2QL8	1.50	64.11	2007-07-12	Se
3QC0	1.45	61.81	2011-01-14	Se
3LLX	1.50	60.72	2010-01-29	Se, Zn
3SEE	1.25	58.46	2011-07-10	Se
2QPX	1.40	55.71	2007-07-25	Se, Zn
3NO2	1.35	51.48	2010-07-24	Se
3N6Z	1.30	50.63	2010-05-26	Se
3CJM	1.50	46.77	2008-03-13	Se
3POH	1.55	42.98	2010-11-22	Se
AcNIR(1)	1.42	39.70	-	Cu
AcNIR(2)	0.99	39.70	-	Cu
2R01	1.15	39.29	2007-08-17	Se
3NOH	1.60	37.24	2010-07-25	Se
3UE2	1.23	34.10	2011-10-28	Se
2R0X	1.06	29.7	2007-08-21	Se

## 9.2 Comparison of electron density maps

With the applied methods, electron density maps from SHELXE and SHARP were available in addition to the maps deposited to the PDB. In the case of 2QPX, maps from molecular dynamics (MD) simulations were provided by James Holton. In total, a comparison between the PDB map and SHELXE could be made for 75 structures. For 11 structures, the maps could be compared to the SHARP map additionally. The comparison was made between entire maps, but also between partial maps containing only the protein region or the solvent region. While manual inspection was performed first to identify striking differences, statistical methods were used to quantify those and to identify systematic differences. The increased amount of electron density peaks in maps from SHELXE and SHARP was the most striking of the initial observations, which is depicted in Figure 31.



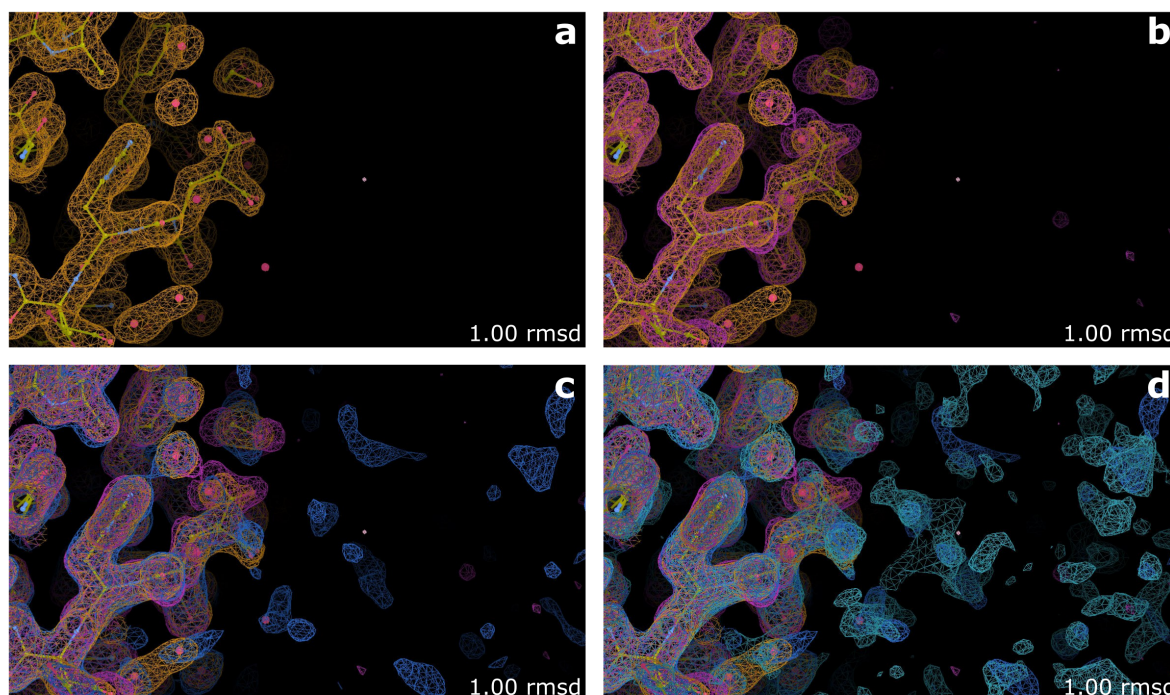
**Figure 30:** *Electron density maps from AcNIR data sets 1 and 2 at different stages. a: map of substructure of AcNIR(2) after ANODE at high threshold, showing strong peaks around copper and sulfur atoms. b: map from PHS file of AcNIR(2) after SHELXE. Ripples in the electron density are observable. c: map of AcNIR(2) after first round of SHARP. While showing partially connected volumes, it does not agree the structure model. d: map of AcNIR(1) after third round of SHARP. The electron density map agrees with the structure model. Same camera angle as in (c). The columns FB and PHIB were used in (c) and (d)*

### 9.2.1 Surface water comparison

As observable from Figure 31, the methods without model bias and solvent flattening, namely the electron density maps calculated with ANODE and SHELXE or SHARP, provide more information about the solvent region far beyond the protein's surface. However, differences are also visible for the surface waters as depicted in Figure 32. Multiple different scenarios are observable: modelled waters positioned in peaks of all three maps, in just two of the maps, in just a peak of single map, and peaks from the SHELXE map overlapping only with a peak from the SHARP map, which look like density of waters, as well as numerous peaks from the SHELXE or SHARP map without any overlap. To quantify the number of modelled waters positioned in peaks common across all maps, a custom python script was written. The resulting counts are listed in Table 10.

First of all, the proportion of waters common across density peaks of all three maps ranges from 10.4 % to 58.9 %. The total number of modelled waters ranges from 189 to 612. The proportion of waters located in peaks of the PDB and SHELXE maps ranges from 39.1 % to 92.4 %, while for waters common in the PDB and SHARP maps this





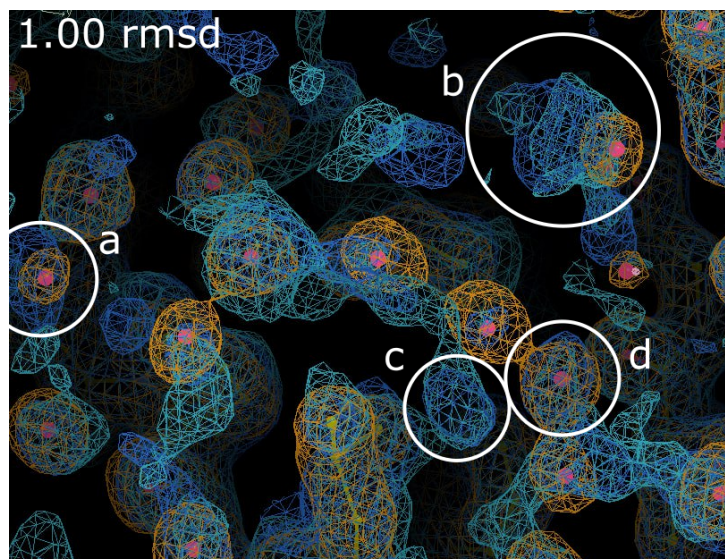
**Figure 31:** Comparison of electron density maps obtained by different methods from the data of PDB structure 2QPX. *a*: model and map from deposited MTZ file in orange, loaded from the columns FP and PHIC. *b*: map from MD simulations in pink. *c*: map from SHARP in blue, from the columns FB and PHIB. *d*: map from SHELXE in cyan, from PHS file (contains only *a*).

number ranges from 21.4 % to 73.7 %. It is noticeable that in all cases, except for 2QPX and AcNIR, more waters are common among the PDB map and SHELXE map than among the PDB map and the SHARP map. For 2QPX and AcNIR, this is vice versa, although the values differ by only 3.2 % and 6 % respectively, while differences in other cases reach up to 36 % as in the case of 3SEE.

### 9.2.2 Similarity of partial maps

While differences in the solvent region were striking, a close look at the electron density in the protein region revealed also some differences. Quantifying the observable differences in certain regions of the electron density maps was therefore the next step. First, partial maps were generated by masking the entire map with the protein model after removing all waters from the structure. The masking resulted in two additional maps per previous map: the protein region only, and the solvent region only, where all grid points which were not belonging to the respective region were set to zero.

After generating all maps, cross correlations were calculated with the results listed in Table 11. In all cases, the cross correlation value for the protein region was higher than for the respective entire map, which is also seen in the minimum and maximum values.



**Figure 32:** Three electron density maps and the PDB structure 2QPX, where differences in peaks for surface waters are visible. The orange map is from the MTZ file deposited to the PDB, loaded from the columns FP and PHIC. The blue map is calculated from SHARP, loaded from the columns FB, PHIB. The cyan map is from SHELXE, loaded from the PHS file. a: this modelled water is only covered by density from the PDB map and SHARP map. b: density peaks from all three maps are close, but vary in size and shape. The water is also only centered in the map from the PDB. c: an overlap of peaks from SHARP and SHELXE maps without a modelled water or peak from the PDB map. d: a modelled water located in peaks of all three maps.

While the protein region comes with values of up to 0.923, the solvent region comes with a maximum correlation of 0.564 and a minimum of 0.120. Splitting the solvent region further into surface solvent region and deep solvent region (Figure 33), leads to increased cross correlation values (Table 12). The deep solvent region shows no correlation, since the PDB map lacks strong density in that region. Furthermore, in all cases, except for 2QPX and AcNIR, the cross correlation value was higher for comparison with the SHELXE map than with the SHARP map, meaning that the SHELXE map has more similarity to the PDB map than the SHARP map. Only for 2QPX and AcNIR it is exactly vice versa. These are also the only cases, where the cross correlation of the solvent region between PDB map and SHARP map is above 0.29, while for the SHELXE map all of the values are above 0.27. The SHARP maps of 3UE2 and 2R01 agree with large portions of the structure, but they show also many gaps in the density for some side chains, while other locations has more density close to the main chain.

Since the deep solvent is only present in maps from SHELXE and SHARP, additional cross correlations were calculated between maps from these two sources (Table 13). Interestingly, the differences in the protein region are relatively high in all cases except for 2QPX. Cross correlation values in the surface solvent are always higher than in the deep solvent, but are lower than the values in the comparison between PDB maps and



**Table 10:** *Thirteen structures for which maps from PDB, SHELXE, and SHARP are available, sorted descending by crystal solvent content. Listed are how many of the modelled waters are covered by density peaks of the respective maps. The column "Common waters" contains counts of waters, which appear in density peaks of all three maps. "in PDB & SHELXE" counts only waters, which appear in the PDB map and the SHELXE map, while "in PDB & SHARP" does the same for peaks from the PDB map and the SHARP map. The number before the slash indicates the respective counts, while the number behind the backslash is the total number of modelled waters in that structure. The respective percentage is given as well.*

PDB	Common waters in all maps	Common waters of PDB & SHELXE maps	Common waters of PDB & SHARP maps
3OHG	49 / 470 (10.4%)	184 / 470 (39.1%)	101 / 470 (21.4%)
2QL8	83 / 431 (19.2%)	224 / 431 (51.9%)	143 / 431 (33.1%)
3QC0	54 / 267 (20.2%)	167 / 267 (62.5%)	91 / 267 (34.0%)
3LLX	143 / 543 (26.3%)	364 / 543 (67.0%)	202 / 543 (37.2%)
3SEE	204 / 397 (51.3%)	367 / 397 (92.4%)	223 / 397 (56.1%)
2QPX	290 / 492 (58.9%)	347 / 492 (70.5%)	363 / 492 (73.7%)
3CJM	73 / 352 (20.7%)	182 / 352 (51.7%)	114 / 352 (32.3%)
3POH	210 / 612 (34.3%)	448 / 612 (73.2%)	272 / 612 (44.4%)
AcNIR(1)	73 / 368 (19.8%)	143 / 368 (38.8 %)	165 / 368 (44.8 %)
2R01	89 / 227 (39.2%)	178 / 227 (78.4%)	107 / 227 (47.1%)
3UE2	33 / 189 (17.4%)	98 / 189 (51.8%)	59 / 189 (31.2%)

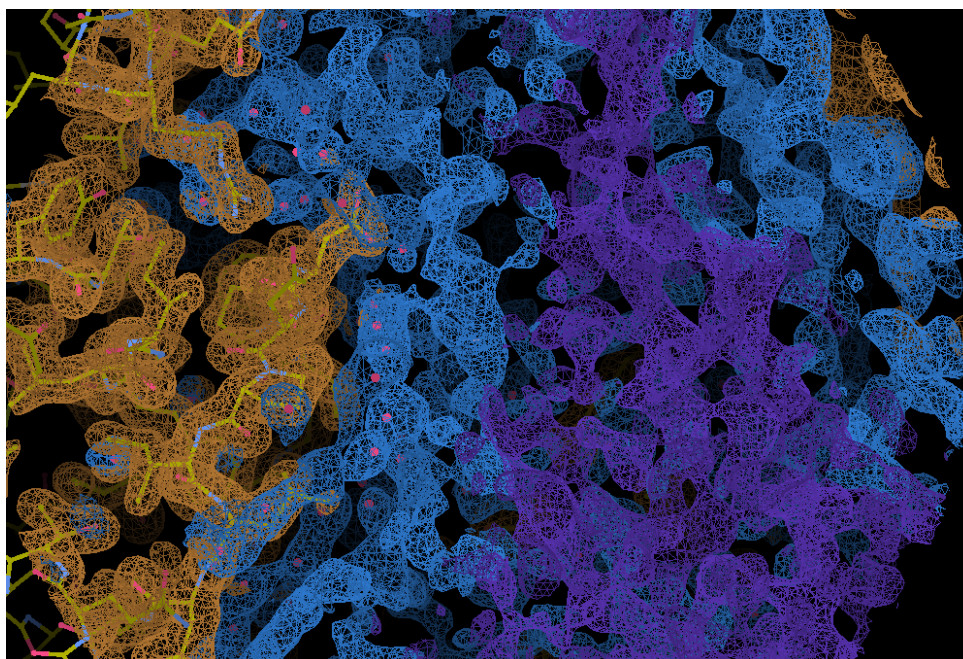
SHELXE maps (listed in Table 12). Interestingly, the cross correlation between the deep solvent maps is relatively high and reaches in the case of 2QPX a value of 0.681. All other correlations between deep solvent regions are much smaller with values below 0.44. The correlation value for the deep solvent is in all cases 0 for 3UE2, since the solvent channels are so narrow, that there is no deep solvent (it has the lowest solvent content, as seen in Table 9).

### 9.2.3 Molecular dynamics simulation

James Holton provided a series of molecular dynamics (MD) simulations for the solvent of the structure 2QPX. From the simulated atoms, an electron density map was calculated and split into a protein-only and solven-only map, as for the other maps in the previous section. Calculation of cross correlation measured the similarity of the simulated maps to the other ones. The results were used to optimize the simulation further, which led to higher cross correlation values. The results of the simulation with the highest values are shown in Table 14. The overall density and density in the protein region are most similar to the PDB map. For the solvent region, no correlation to any of the maps was observed, which was the case for all simulations. Analysis of surface waters (see section 9.2.1) shows that 194 out of 492 modelled waters are in peaks of the PDB map and the MD map. A total of 106 modelled waters is located in electron density peaks of all four

**Table 11:** Cross correlation values between entire maps and partial maps, which contain only the protein region or the solvent region. Waters were removed from the PDB structure before masking the map and a masking radius of 2.5 Å.

PDB	CC SHELXE map to PDB map			CC SHARP map to PDB map		
	entire map	protein region	solvent region	entire map	protein region	solvent region
3OHG	0.674	0.793	0.288	0.488	0.607	0.157
2QL8	0.741	0.844	0.339	0.512	0.662	0.174
3QC0	0.770	0.861	0.349	0.537	0.650	0.172
3LLX	0.794	0.871	0.420	0.558	0.658	0.218
3SEE	0.844	0.904	0.564	0.588	0.694	0.289
2QPX	0.841	0.896	0.459	0.880	0.923	0.530
3CJM	0.752	0.830	0.370	0.548	0.642	0.189
3POH	0.857	0.898	0.563	0.544	0.603	0.252
AcNIR	0.566	0.605	0.270	0.798	0.816	0.404
2R01	0.787	0.833	0.484	0.516	0.575	0.228
3UE2	0.712	0.747	0.482	0.524	0.561	0.282
min.	0.674	0.747	0.288	0.488	0.561	0.157
max.	0.857	0.904	0.564	0.880	0.923	0.530



**Figure 33:** Separation of the entire SHARP map of 2QPX into protein region (orange), surface solvent (blue), and deep solvent (violet). All maps are shown at 1.0 RMSD.

maps, i.e. MD map, PDB map, SHELXE map, and SHARP map, which are 21.5 % of all waters. This is in contrast to the 58.9 % of waters common among the maps from the PDB, SHELXE, and SHARP.

**Table 12:** Cross correlation values between partial maps of the solvent region, which contain only the surface solvent region (generated via a masking radius of 7.5 Å from atoms minus the protein region) or the deep solvent region (entire solvent region minus surface solvent region). Waters were removed from the PDB structure before masking.

PDB	CC SHELXE map to PDB map		CC SHARP map to PDB map	
	surface solvent	deep solvent	surface solvent	deep solvent
3OHG	0.388	0.016	0.220	0.004
2QL8	0.394	0.078	0.210	0.014
3QC0	0.444	0.035	0.230	0.007
3LLX	0.513	0.068	0.273	0.029
3SEE	0.614	0.040	0.327	0.025
2QPX	0.520	0.037	0.591	0.072
3CJM	0.397	0.013	0.206	-0.004
3POH	0.576	0.084	0.262	-0.014
2R01	0.492	0.117	0.233	0.066
3UE2	0.482	0	0.282	0

**Table 13:** Results from cross correlation between partial electron density maps from SHARP and SHELXE.

PDB	entire	protein	solvent	surface	deep
	map	region	region	solvent	solvent
3OHG	0.571	0.668	0.390	0.421	0.351
2QL8	0.650	0.760	0.468	0.478	0.438
3QC0	0.621	0.711	0.409	0.440	0.355
3LLX	0.641	0.728	0.407	0.447	0.327
3SEE	0.636	0.723	0.426	0.452	0.313
2QPX	0.930	0.955	0.783	0.806	0.681
3CJM	0.582	0.662	0.353	0.359	0.317
3POH	0.602	0.649	0.405	0.410	0.323
AcNIR	0.631	0.652	0.383	0.383	0.298
2R01	0.592	0.636	0.417	0.418	0.397
3UE2	0.544	0.567	0.423	0.423	0

### 9.3 Refinement against electron density map from SHARP

The measurable differences in electron density of the protein region and surface waters were leading to the hypothesis that protein models could be improved. To test this hypothesis, structure models were refined against the SHARP map and alternative surface waters were placed, where the refinement started with the PDB model only with the waters common among electron density peaks of PDB, SHELXE, and SHARP maps, while all other waters were removed. The refinement did not use any solvent mask until the very end and after the first automatic placement of waters, all other waters were manually added or deleted. Table 15 compares the  $R_{\text{work}}$  and  $R_{\text{free}}$  from these refinements with the deposited R values from the PDB. For 2QPX, it was possible to lower the R values in comparison to the deposited structure. While the regular refinement and addition of

**Table 14:** *Cross correlation between maps from molecular dynamics simulation (run opt59\_all\_atom) by James Holton and maps from the PDB, SHELXE, and SHARP.*

map origin	entire map	protein region	solvent region
PDB	0.928	0.920	-0.020
SHELXE	0.800	0.840	0.020
SHARP	0.830	0.860	0.010

alternative waters did not get close enough, applying a solvent mask to flatten the deep solvent region was helpful. A similar situation is observable for AcNIR. However, it is important to note, that the AcNIR SHARP map resulted from a different data set than the PDB map. The PDB map contains reflection to a resolution of 1.40 Å, the AcNIR data set to a resolution of 1.42 Å. For 3SEE, the solven mask increased the values and none of the attempts could lower them again.

**Table 15:** *R values from refinement of the PDB model against the SHARP map in comparison to the R values deposited to the PDB. The refinement started with the PDB structure only with waters, which were common among electron density peaks of the PDB, SHELXE, and SHARP map. A solvent mask was avoided as long as possible and was only used in the very last refinement step.*

PDB	PDB		SHARP no solv. mask		SHARP solv. mask	
	R <sub>work</sub>	R <sub>free</sub>	R <sub>work</sub>	R <sub>free</sub>	R <sub>work</sub>	R <sub>free</sub>
2QPX	0.134	0.162	0.164	0.206	0.129	0.158
3SEE	0.130	0.149	0.168	0.177	0.176	0.194
AcNIR	0.153	0.181	0.179	0.208	0.140	0.172

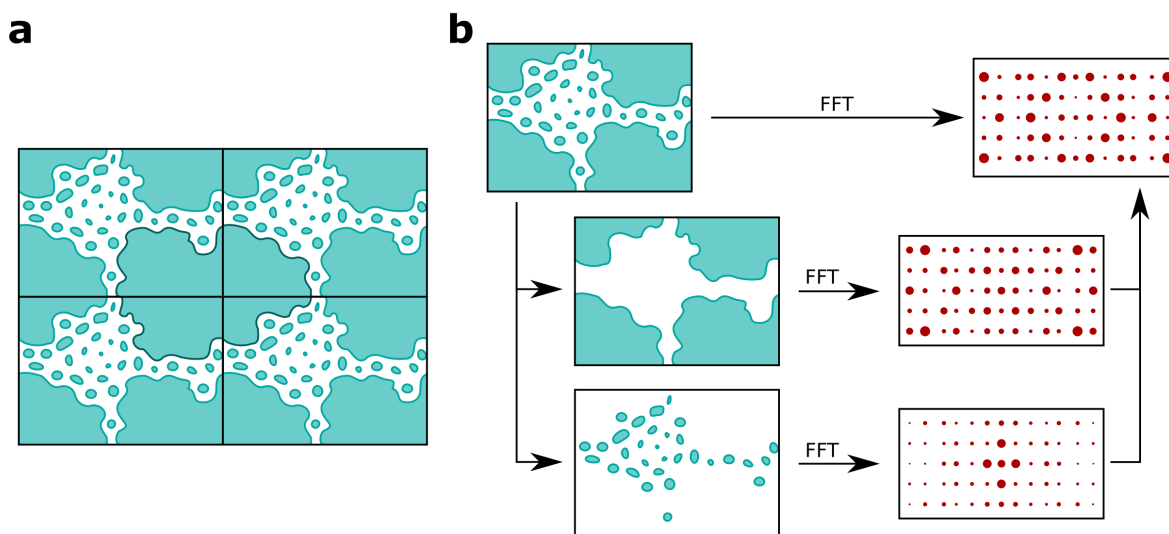
## 9.4 Analysis of solvent region contribution

The differences in the solvent region were measurable, but comparison of electron density in real space did not provide insights on how relevant the observed peaks are and if they result from true solvent contents or noise. More important, however, was the question about how influential the solvent region is on the overall data of the entire unit cell. To address this question, a comparison of data in reciprocal space was conducted.

### 9.4.1 Map separation error

The separation of electron density maps into partial maps, such as protein or solvent region only, potentially has a negative impact on the data represented in reciprocal space, when converting the grid points of the electron density map into structure factors via the Fourier transform. Since the following steps analysed the contribution of the solvent region to each structure factor, it was important to assess the errors, which may have been introduced by map separation.

The idea behind the following procedure is illustrated in Figure 34. The structure factors of the partial maps ( $F_{\text{PROT}}$  for the protein region and  $F_{\text{SOLV}}$  for the solvent region), which are obtained by Fourier transform of a map's grid points, can be represented as complex numbers and vectors in an Argand diagram (Figure 35). The deviation between  $F_{\text{TOTAL}}$  and  $F_{\text{PROT}} + F_{\text{SOLV}}$  quantifies the error, since ideally no deviation is measurable. The analysis was performed for the 74 PDB structures for which PDB maps were available and SHELXE maps were generated. The 11 SHARP maps were also analysed. Statistics about deviations between  $F_{\text{TOTAL}}$  and  $F_{\text{PROT}} + F_{\text{SOLV}}$  for each set of maps are listed in Table 16.

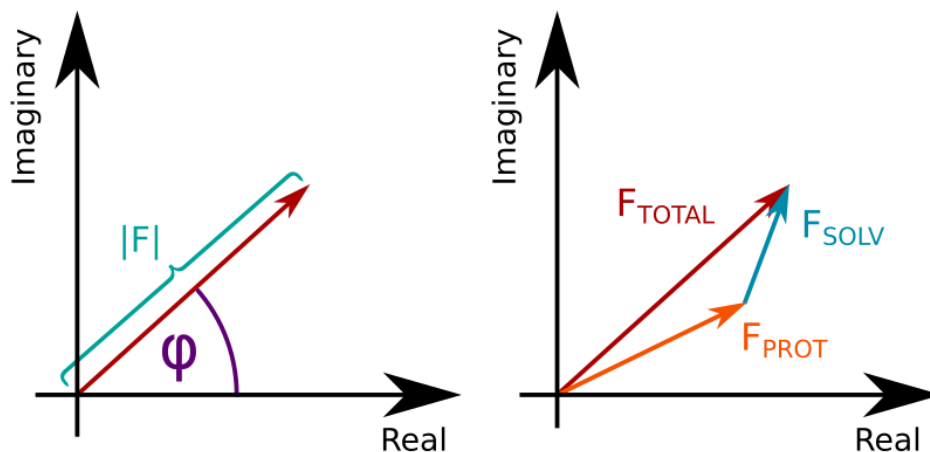


**Figure 34:** *a: Schematic illustration of four units cells with the electron density map of a protein (dense region) and the surrounding solvent channels. b: Separation of the entire density map into protein region only and solvent region only. Fast Fourier transform provides the structure factors of each map. Combining the structure factors of the partial map results in theory in the structure factors of the entire map.*

**Table 16:** *Statistics about percentages of structure factors, which were negatively affected by the map separation. The number in the parentheses indicates the number of maps, that were analysed.*

Map source	minimum	median	mean	maximum
PDB (74)	0.01 %	0.34 %	0.55 %	2.56 %
SHELXE (74)	10.08 %	20.17 %	20.56 %	45.34 %
SHARP (11)	0.04 %	0.30 %	0.50 %	2.50 %

The statistics for the PDB maps and SHARP maps look similar, where the majority of maps has less than 1 % of errors introduced (Table 16). The maximum reaches in both cases around 2.5 %. The map with the minimum error percentage from SHELXE shows four times as many structure factors affected by the map separation, while the median and the maximum have an eight times and 18 times higher proportion of affected structure



**Figure 35:** Argand diagram visualizing complex numbers in a 2D plane, where structure factors can be represented as vectors (left). Structure factors of partial maps can be added like vectors to obtain the structure factor of the entire map.

factors, respectively.

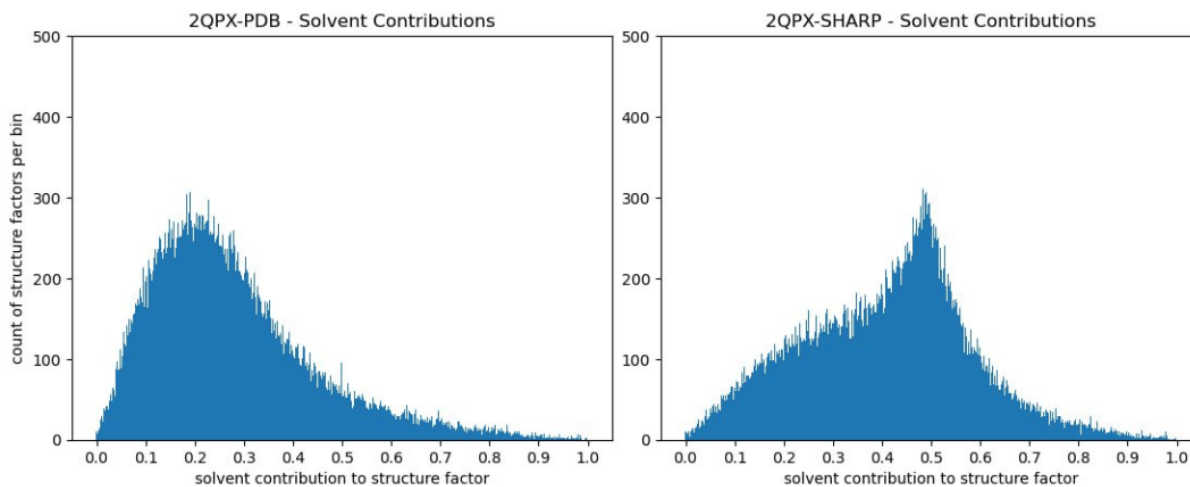
#### 9.4.2 Solvent contribution to structure factors

With the partial maps at hand and a method to evaluate the number of structure factors corrupted by errors, it is now possible to calculate how much of a structure factor's amplitude is determined by electron density of the solvent region. This was accomplished by

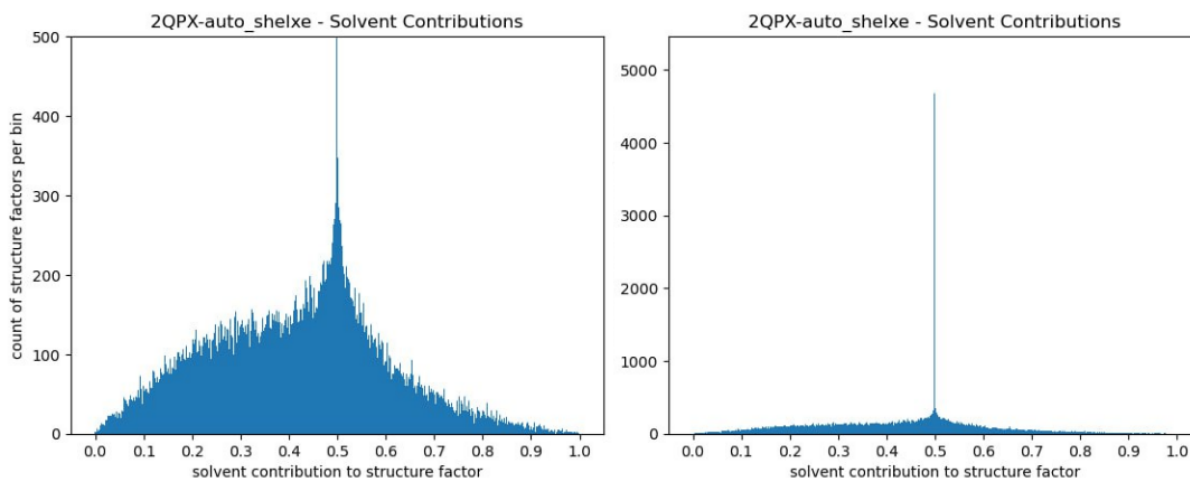
$$\text{contribution}_{SOLV} = \frac{|F_{SOLV}|}{|F_{SOLV}| + |F_{PROT}|} \quad (8)$$

For a set of maps, the contribution from the solvent region, i.e. the solvent contribution to a structure factor, was calculated for each structure factor, which were counted in a histogram. Figure 36 shows the histogram from the PDB map in comparison to that from the SHARP map, both for the structure 2QPX. A smooth distribution is observation for the PDB plot with a peak at  $\tilde{20}$  % solvent contribution, where the majority of counted structure factors has a solvent contribution below 50 %. The plot for the SHARP map, however, has on the first look a very different distribution. The picture becomes clearer when considering also the plot for the SHELXE map of 2QPX, which is depicted in Figure 37. The SHELXE plot shows an extreme peak at exactly 50 % solvent contribution. This anomaly is also present in the plots of the PDB and SHARP maps, but in much weaker form. Where this peak may comes from, is further analyzed in the next section, but here, it is important that around the peak a symmetrical distribution is recognizable, which is also present in the SHARP plot and highlighted in Figure 38. If filtering out the 50 % solvent contribution peaks, a common distribution becomes visible, which shares

similarity with the distribution of the PDB map. However, this distribution is flatter and shifted towards higher solvent contributions compared to the PDB plot.



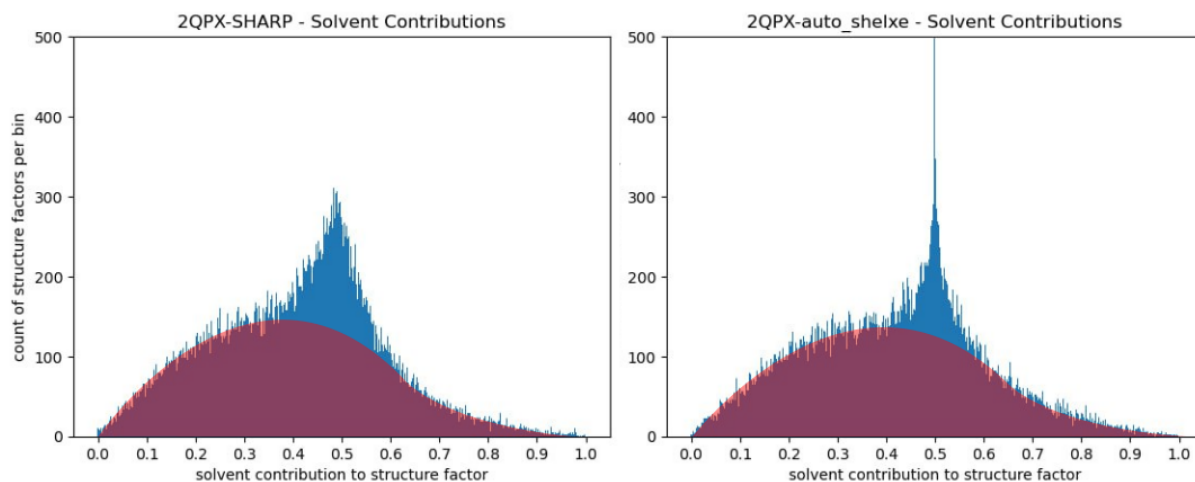
**Figure 36:** Histogram counting the number of structure factors per bin solvent contribution bin. Comparison of maps from PDB (left) and SHARP (right) of the structure 2QPX.



**Figure 37:** Histogram counting the number of structure factors per bin solvent contribution bin. Both histograms show data from SHELXE maps of the structure 2QPX, but the left one does not show the whole plot, as notable by the Y-axis.

While the histograms show a generally higher solvent contribution to structure factors in the SHELXE and SHARP maps, they do not show how much these structure factors contribute truly to the map. The plots depicted in Figure 39 show the normalized amplitudes against the resolution, colored by solvent contribution. The filtered plots (bottom row) show more higher amplitudes, but in these plots alone, it is not easy visible, that more structure factors are in the SHARP plot. In the unfiltered plots (top row), however, higher amplitudes are observable for the PDB plot, which have solvent contribution below 50 %. Otherwise, the distributions look similar.





**Figure 38:** Histograms from Figure 36 and Figure 37 with highlighted similar distribution hidden below the 50 % solvent contribution peaks.

Plots for other structures looked similar, but an analysis over a larger number of structures would provide more insights. Therefore, the 74 structures with PDB maps and SHELXE maps were compared in single metrics. Figure 40 shows all structures sorted by crystal solvent content along the X-axis. The mean solvent contribution across all structure factors of a structure is in all cases higher in the SHELXE map than in the PDB map and within both groups, the mean solvent contribution increases with higher crystal solvent content. Correlation coefficients with the shown regression lines are similarly high, but some outliers and variance are clearly visible.

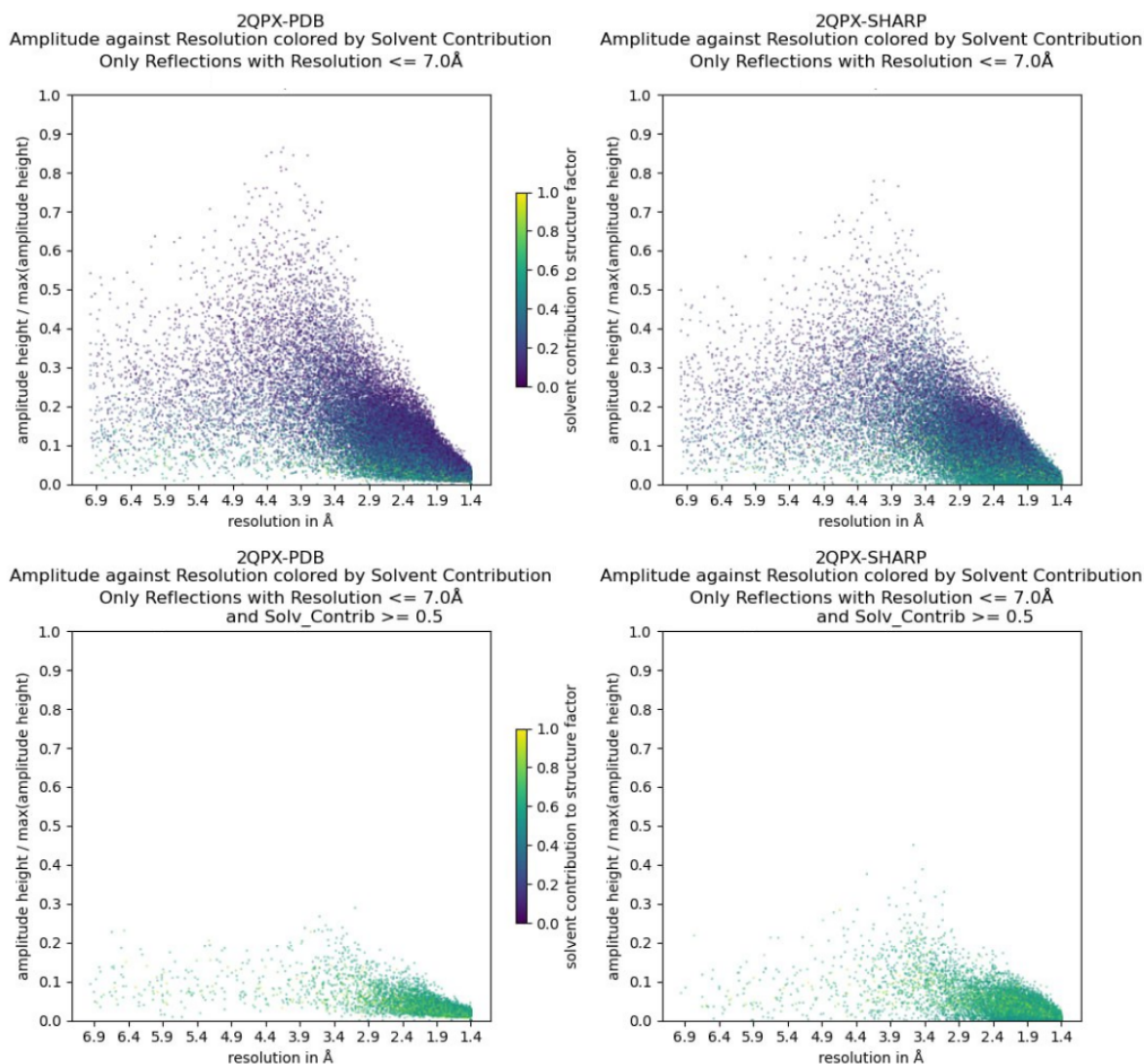
Figure 41 shows how many of the structure factors have a solvent contribution of 51 % or higher. In all PDB maps, this proportion lies below 5 %. For the SHELXE maps, however, the regression line with a correlation coefficient of 0.44 ranges from 10 % above 20 % and in individual cases, more than 25 % of all structure factors have a solvent contribution above 50 %.

Finally, Figure 42 shows how high the impact of the solvent region is on the whole map, by comparing the mean amplitudes multiplied by the respective solvent contribution of that structure factor. Again, SHELXE maps show higher values and a steeper regression curve, although this time, a single PDB map shows high values than the respective SHELXE map.

### 9.4.3 Fifty percent anomaly

The SHELXE maps had a severe influence of errors upon separation of the density map into the protein and solvent regions, as seen in Table 16. Furthermore, suspicious peaks at 50 % solvent contribution are visible in the histograms of Figure 37, where much higher

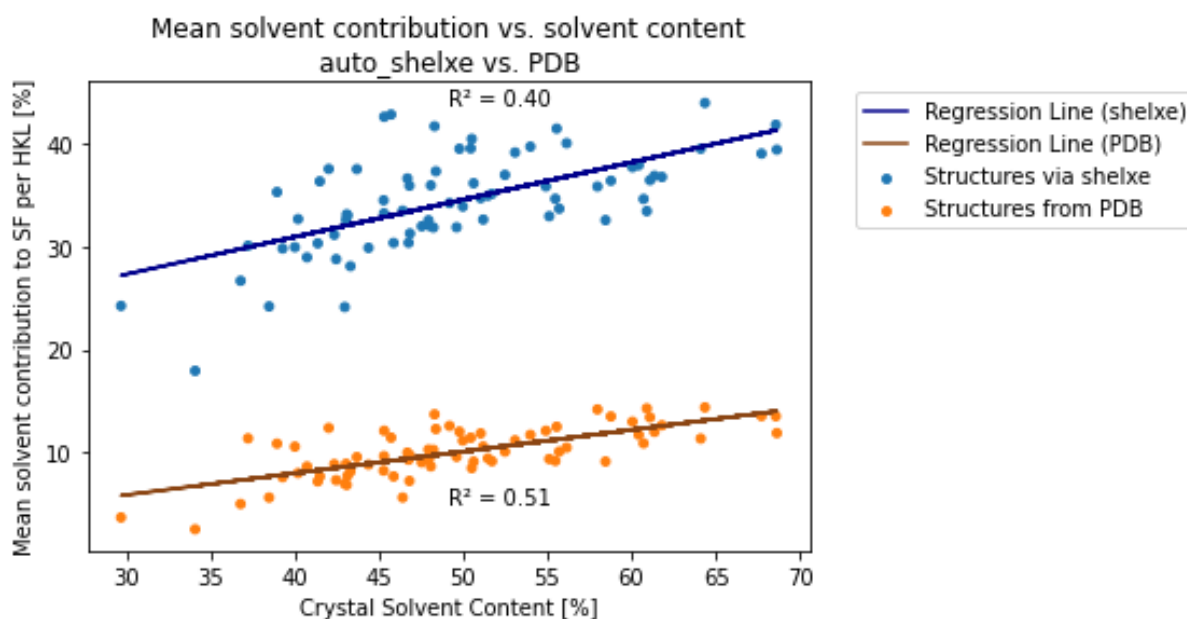




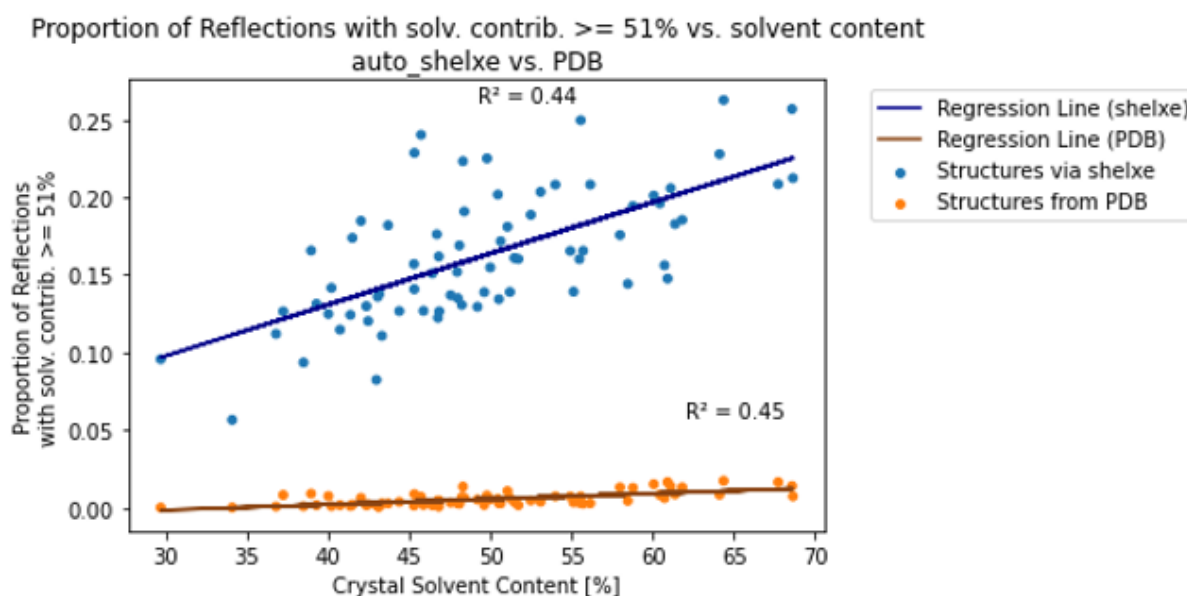
**Figure 39:** Scatter plot showing structure factors plotted by normalized amplitudes against resolution and colored by solvent contribution from PDB map (left column) and SHARP map (right column) of 2QPX. In the bottom row, structure factors with solvent contribution below 50 % are filtered out.

peaks are observable for SHELXE maps than for the other maps.

Deeper investigation shows that across all maps from SHELXE at least 99.2 % of all scaling factor outliers (structure factors in Table 16) have a solvent contribution between 49 % and 51 % and an amplitude value below 0.1. The median of this percentage across all maps is at 99.8 %, meaning that in all SHELXE maps almost all scaling factor outliers are in the peak seen in Figure 37. For the maps from the PDB and SHARP, much fewer scaling factor outliers were observed (Table 16), but in almost all cases more than 99.9 % of these have also amplitudes below 0.1 and a solvent contribution of around 50 %. In fact, most of these amplitudes are smaller than  $10^6$  and only a few are single structure factors show amplitudes above 1.



**Figure 40:** Scatter plot with mean solvent contribution to structure factors plotted against the crystal solvent content. 74 maps from the PDB are plotted (orange) together with 74 SHELXE maps (blue) generated from the same data sets, where two dots resulting of the same data set sharing the same X coordinate.

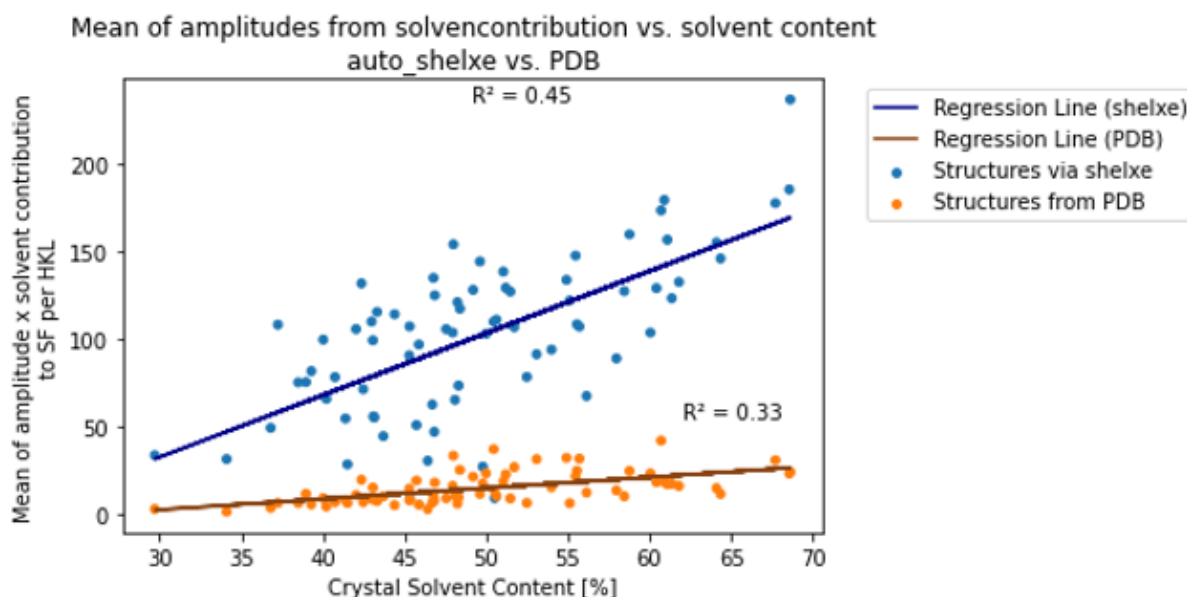


**Figure 41:** Scatter plot with proportion of structure factors per map, that have solvent contribution above 51 %, plotted against the crystal solvent content. Maps from PDB are plotted in orange, SHELXE maps in blue.

## 10 Discussion

### 10.1 MAD phasing

In order to assess the influence of solvent flattening and model bias on electron density maps, a large number of cases must be examined. Here, the analysis started with 573



**Figure 42:** Scatter plot with mean amplitudes times solvent contribution to structure factors plotted against the crystal solvent content. Maps from PDB are plotted in orange, SHELXE maps in blue.

high resolution structures solved from MAD phasing experiments. There were filtered by certain criteria down to 74 structures, for which an electron density map was calculated via ANODE and SHELXE. The quality of those maps was evaluated and the best cases were phased additionally with SHARP.

With 74 structures out of 574, only a fraction of the collected datasets led to MAD maps that were comparable to the deposited maps in the PDB. Because the initial number was so large, manual phasing was not an option. Therefore, an automated process was developed. However, the automatic substructure solution and phasing are potentially prone to errors, which could have rejected entries that could have led to high quality maps. Moreover, datasets missing R values or having a completeness slightly below the threshold were filtered out, although some of these cases could still have led to useful maps as well. Some of the promising cases had even to be filtered out due to submission errors where file formats were unusual or two wavelengths had equal values, as these mistakes resulted in problems within the automated pipeline. Such mistakes show, that the submission process to the PDB did not catch such errors, which persist many years later. After all, none of the remaining structures was deposited later than December 2012, which highlights that MAD phasing is no longer a common method for experimental phasing.

Of the remaining 74 datasets, all had a successful substructure solution, which is also reflected in high averaged anomalous density (AAD) values (Figure 29). These values vary between structures based on the anomalous scatterers. However, as seen for the case of AcNIR, a good substructure solution does not result always in a useful electron density

map (Figure 30). Therefore, the AAD values alone were not sufficient for quick quality estimates. Since the goal was to generate maps comparable to the deposited ones, high cross correlation values between the maps were expected to work better. Because the solvent regions showed many differences, cross correlation values were low. Therefore, the easiest solution was to calculate the correlation against the structure model and evaluate the differences of local cross correlation values of the SHELXE map and the PDB map. Together with the AAD values this provided an overview about structures with strong anomalous differences and for which similar maps from PDB and SHELXE were available.

The best cases were then selected for manual phasing with SHARP, of which all have been inspected manually. Table 9 shows that not all of the datasets which have led to good maps in SHELXE were successfully phased in SHARP. The structure 2R0X, for example, had a high AAD value of 70.6 and a relatively moderate cross correlation difference of 0.254 (Table 8), but did not result in map with SHARP that agreed with the structure, similar to the map shown in Figure 30c. Repeating the SHARP phasing in the hope of avoiding a previous error resulted in the same map. However, since phasing was possible with SHELXE, it should also be possible with SHARP. The correct input parameters are therefore important and it could be possible to improve the quality of the already successful maps by providing the right settings.

Conclusively, 74 maps phased with SHELXE and 11 map phased with SHARP were obtained, of which most showed good to excellent quality. More important, though, was the fact that these maps were free of model bias and that no modification of the electron density in the solvent regions was applied, which provided the basis for the following analysis.

## 10.2 Comparison of electron density maps

The electron density maps calculated with SHELXE and SHARP were compared to the maps generated from the MTZ files deposited to the PDB. Since all three maps were based on the same collected reflection data, this comparison allowed the investigation of differences in maps resulting from the use of different methods for calculating those maps. The most striking difference were the many peaks in the SHELXE and SHARP maps, while the PDB maps had no peaks at all in the respective regions. This was the result of solvent flattening, which flattened the electron density of the solvent region to allow better refinement of the density in the protein region.

In this work, the SHELXE and SHARP maps were explicitly generated without the use of solvent flattening. Furthermore, no extension of reflection data via 'free lunch' algo-

rithms was used, nor any other type of density modification not based on the anomalous differences. While the investigated cases were all from MAD experiments, it is not clear if they relied only experimental phasing information. All structures were solved by Joint Center for Structural Genomics projects, but only for 3OHG a publication is available [168]. The substructure of 3OHG was solved via SHELXD [162] and SOLOMON was used inside autoSHARP [169] for density modification, followed by ARP/wARP [170] for model building. A look into the phasing methods listed in the PDB entries, show that all structures were phased via SHELXD and autoSHARP except for 3SEE and 2R01, which were phased with SHELXD and SHARP, 3UE2, which was phased with SOLVE. In this work, specifically SHARP was used over autoSHARP to gain control over all applied modifications, since in SHARP additional density modification and solvent flattening is optional and was skipped in generating the SHARP maps of this work. In autoSHARP, however, this step is applied automatically and it is likely the case, that all of the examined structures made use of it, which is reflected in the low cross correlation values between the maps. The cases 3SEE and 2R01 were phased via SHARP, but the cross correlation values indicate here as well, that solvent flattening and density modification were applied. Regarding the model bias, it is not possible from the given information which structures made use of molecular replacement and which relied on *ab initio* model building. It is also important to note that none of the structures from the PDB made use of ANODE, since this tool was not available before the year 2011. If there is measurable difference between using SHELXD or ANODE was not examined in this work.

It was also observable that cross correlation values between the PDB maps and SHELXE maps were always higher than those between PDB maps and SHARP maps, except for 2QPX and AcNIR. These were also the only exceptional cases for common surface waters. 2QPX was extensively refined with SHARP, where several iterations were applied and alternative trajectories were tested. So it is possible that in this case the SHARP parameters were just optimized. For AcNIR, however, it is unclear since only two additional refinement runs were applied. With the given data, it is not possible to explain this phenomenon. It is clear, however, that SHELXE is much easier to use since it requires only a single command, while SHARP requires expertise to set the input parameters correctly.

Another important aspect regarding the methodology is the masking of maps. While splitting the map of grid points into two partial maps works fine in real space as adding both partial maps results again in the original map, it causes problems in the reciprocal space, which are discussed late in section 10.4. For the cross correlation analysis in real

space, however, it was sufficient and did not introduce any artifacts.

The differences in the solvent region can be classified into two categories: differences at the surface level and differences in the deep solvent channels. Close to the protein's surface, well-defined electron density peaks typical for waters are observable for all three maps. However, closer investigation revealed that only a fraction of these are common among all maps. The highest fraction of modelled waters which were in peaks of all three maps was 59 %. Between a pair of maps, this number was 92 % and the second highest fraction was already 74 %. These strong differences allow potentially alternative models of surface waters, since only one of the maps could describe the experimental reality. Which one has the best water model, however, is not easy to identify. The numbers of total modelled waters ranged strongly from 189 to 612 waters and the fraction of waters common across all maps ranged from 10 % to 59 %. Potential factors contributing to these vastly different numbers are the crystal solvent content, the depth of solvent channels, protein surface area intersecting with the solvent region, and most importantly, the method behind placing the water molecules, i.e. the experience of the scientist and if density peaks which could result from waters were over- or underfitted. Surface residues may also have an influence on ordered waters. Cross correlation between the surface solvent region of the PDB maps and the maps from SHELXE show also moderate to high similarity with values ranging from 0.388 to 0.614. For SHARP, these values are lower except for 2QPX, which shows a relatively high correlation value of 0.591. In general, the cross correlation values indicate more similarity than the common waters for some structures such as 3OHG, 3CJM, or 3UE2, which could hint at too few modelled waters. All in all, the results from this work indicate the possibility for improved, alternative surface water models, but how this information can be beneficially and reliably applied during model building remains to be elucidated.

Regarding the deep solvent region beyond the surface waters, numerous electron density peaks are observable. However, these are mostly weak, strongly anisotropic, and/or form branched blobs. The arising questions are: do these density peaks result from noise or do they actually contain information which can be utilized in modelling of solvent channel contents? The analysis of solvent region contribution tried to answer these questions, as discussed in section 10.4. While cross correlation showed large differences in the solvent regions, minor differences were notable for the protein regions. However, even the highest value of 0.92 is large enough for alternative protein structure models with minor tweaks in the backbone and side chains. If such alternative models can be used to lower

R values, is discussed in the next section.

The maps from molecular dynamics (MD) simulation were most similar to the PDB map, which was expected since the MD simulation used the model and introduced therefore a bias that is absent in the other two maps. The complete lack of correlation in the solvent region, however, is surprising. Because 39 % of the modelled solvent waters were in peaks of the PDB map and the MD map, one would expect similar density at the surface region. On the other hand, the solvent region of 2QPX makes up 55 % of the entire unit cell, which provides a large volume that can counter the correlation of the surface waters. All in all, the MD simulations did not provide new insights into the deep solvent region, but the surface waters could contribute to improved alternative surface water models. For general statements, however, simulations for more structures must be conducted and analysed.

### 10.3 Refinement against electron density map from SHARP

The refinement examined whether it was possible to obtain a structure model with lower R values, if refining against a purely experimentally determined map without any model bias and without applying density modification, especially without solvent flattening. In the previous section it was already discussed that especially the surface solvent provides electron density for alternative water models, which hold the potential of improving the structure towards the true data.

Starting the refinement with all waters removed, except for those common among the PDB, SHELXE, and SHARP maps, was a good decision, as it lowered the number of waters which had to be removed manually due to poor coverage by electron density while keeping waters which were certainly supported by density of three maps. In all cases, careful placement and removal of waters in COOT during refinement was what lowered R values most. However, at some point it was no longer possible to get lower, unless applying a solvent mask. This reduced the R values drastically and allowed at the same time to model more waters, as some peaks became much clearer water-like density peaks. Both,  $R_{\text{work}}$  and  $R_{\text{free}}$ , were lowered below the deposited R values of 2QPX, which can be seen as a successful refinement. Since rotamer and Ramachandran outliers were not optimized, even lower R values are potentially possible. For AcNIR, also both values were lowered after solvent flattening. However, the SHARP map was here from a different source, making this case not directly comparable. In the case of 3SEE, R values were not lowered, but it could be the case, that the refinement was executed optimally.

In any case, it was shown that structure improvement is technically possible when refining against a purely experimentally determined electron density map and building an alternative solvent model, while applying density modification only after the regular refinement no longer provides any improvements.

#### 10.4 Analysis of solvent region contribution

Separation of maps into different regions provides not only the possibility for real space cross correlation, but also for calculation of the solvent region contribution to the total structure factors. Projecting the partial structure factors onto  $F_{\text{TOTAL}}$  was not suitable, as this was dependent on phases and would cause problems where the projected vector has a length of zero, thus practically negating any contribution of that structure factor, while making the other one look like it is contributing 100 %, although it is not identical to  $F_{\text{TOTAL}}$ . Therefore, a phase-independent way of calculating the contribution was chosen.

The generated plots provided valuable insights into the differences between the SHELXE and SHARP maps and the PDB maps. First, a much higher solvent contribution is observable among all plots, which is no surprise considering the flat solvent of the PDB maps. The similar distribution between SHELXE and SHARP, if ignoring the 50 % solvent contribution peak density, agrees also with the high cross correlation of the solvent region between both maps for 2QPX. Moreover, the increased contribution from the solvent region comes also with high amplitudes from the that region, where the PDB maps are missing this potentially valuable information. It was expected that with higher crystal solvent contents more contribution from the solvent region is introduced, since the solvent region becomes larger and the protein region smaller. However, it was interesting to see a much steeper regression line for the increase of mean amplitudes from the solvent region compared between SHELXE and PDB maps (Figure 42), although the amplitude values were in general rather low.

Regarding the 50 % solvent contribution anomaly, which is especially prominent in the SHELXE plot of Figure 37, it seems at first interesting, that almost all of the scaling factor outliers show such a contribution. However, they also show extremely low amplitude values close to zero. A closer look into the data showed numerous cases across all maps, where  $F_{\text{SOLV}}$  has a very high amplitude.  $F_{\text{PROT}}$  has a very similar amplitude, but a phase angle difference of  $180^\circ$  to  $F_{\text{SOLV}}$ . Combining both would lead to an  $F_{\text{TOTAL}}$  with an amplitude that is practically zero. Due to numerical errors, however, their scaling factor is beyond the set threshold. In many other cases,  $F_{\text{SOLV}}$  and  $F_{\text{PROT}}$  have amplitudes near



zero in the first place, where the phase difference is not restricted to  $180^\circ$  to result in a similar outcome. In any case,  $F_{\text{SOLV}}$  and  $F_{\text{PROT}}$  must have a similar amplitude in order to result in an amplitude near 0 for  $F_{\text{TOTAL}}$ . It is not clear why SHELXE maps have much more structure factors with amplitudes near zero than the PDB or SHARP maps, which could be the result of an underlying process in SHELXE or ANODE. Maybe even multiple sources contributed to this result, since SHARP maps also show higher peaks than the PDB maps. It is, however, unlikely that it is solely the result from the map separation, as it is not equally strong among all cases. Nevertheless, these structure factors do not contribute significantly to the map, since their amplitudes are effectively zero.

Conclusively, the solvent contribution analysis provided some insights about information that was lost in the PDB maps, but the respective structure factors were generally low in amplitudes, thus limiting their impact on the overall map. Also, the solvent was only analysed as a whole portion of the map. Utilizing the surface solvent and deep solvent maps could provide more valuable insights, especially since the alternative surface models have been shown to be able to lower R values.

## 11 Outlook

Experimental phasing methods come with additional efforts and requirements for data collection, processing of reflection data, and phase calculation, but they provide interesting information that is otherwise not present. Here, it was shown that experimentally phased electron density maps without any model bias and density modification can lead to improved structures during refinement, by providing additional information about the surface solvent region and thus enabling alternative water models. However, no general methods of how to build such alternative water models could be developed from the few examples. A much larger study would be required to prove this method to be generally valid. Furthermore, the solvent contribution analysis can be expanded by separating the solvent region into surface and deep solvent and analysing the contributions of the respective regions. Surface residue analysis could also provide valuable information for better modelling of surface waters, since not all amino acids show the same interactions. Finally, R values could successfully be lowered, but the lower limit is still not explored. Pushing the boundary as far as possible can provide more insights on the role of the solvent region in regard to the R-factor gap.

All in all, the deep solvent of macromolecular crystals contains information that is usually lost with conventional methods of phase determination and holds the potential of improving protein structure models. In at least one case, the structure could be improved during refinement and got closer to the experimentally observed data. If this can be applied to any structure, however, must be shown in a large scale study.

## 12 Materials and Methods

### 12.1 MAD phasing

Datasets from MAD phasing were downloaded from the RCSB PDB [82] and filtered based on different criteria. An automatic pipeline was developed to generate electron density maps from these data sets via ANODE [161], and SHELX [158, 162, 167]. The best cases were also phased via SHARP [159].

#### 12.1.1 List of python scripts

A series of python scripts were written during this work to automate the following processes, since the number of structures was too large for manual processing. The scripts for this section of the entire pipeline are:

1. *MAD\_data\_sets\_filter.py*
2. *MAD\_pdb\_stats\_filter.py*
3. *autoMAD\_sca\_converter.py*
4. *autoMAD\_shelxc.py*
5. *autoMAD\_anode.py*
6. *autoMAD\_shelxe.py*
7. *autoMAD\_prep\_for\_SHARP.py*

After downloading all data into the folder "data" and providing a text file with all PDB codes of interest (codes separated by a comma), *MAD\_data\_sets\_filter.py* runs over all entries and keeps only those which have data measured at three or more wavelengths. The PDB codes of kept entries are written into an output text file.

After reducing the initially large number of entries, a more sophisticated filtering following the same base structure was performed by *MAD\_pdb\_stats\_filter.py*. This script filtered out entries that were manually flagged (see next subsection), contained flagged additives in their crystallization conditions, had a resolution above 2 Å, contained more than one crystal, had a completeness below 97 %, or had no given  $R_{\text{work}}$  or had it set to zero. In addition to a text file, it outputs a table, which contains the columns PDB code, Resolution,  $R_{\text{work}}$ ,  $R_{\text{free}}$ , Completeness, Cell Dimensions, Spacegroup, Deposition Date, DOI, Alternative AAs, Solvent Molecules, Crystal Solvent Content in %, Wavelengths,

and Crystallization Conditions. The column Alternative AAs lists anything that does not belong to the canonical 20 amino acids, which included in the given PDB entries non-canonical amino acids and ions.

The script *autoMAD\_sca\_converter.py* identifies the peak-, inflection-, low energy remote-, and high energy remote wavelengths from the pure wavelength values by identifying peak and inflections as the wavelength pair with the smallest difference, where peak has the shorter wavelength and high- and low energy remote have a shorter or longer wavelength than that, respectively. Afterwards, the wavelength-specific datasets are loaded and converted from the .CIF file to .MTZ and .SCA files via the PHENIX [171].

Input files for SHELXC [162] are automatically generated and executed by the script *autoMAD\_shelxc.py*. The resulting output files are used by ANODE [161], which is executed by *autoMAD\_anode.py*. The script handles the modification of the generated output and prepares it for SHELXE [158], which is run by *autoMAD\_shelxe.py* with disabled density modification and 'free lunch'. The same script converts the resulting .PHS files to .MTZ via the CCP4 tool F2MTZ and calculates the overall and local cross correlation with the PDB model via PHENIX.

Finally, *autoMAD\_prep\_for\_SHARP.py* prepares the data to be ready for loading manually into SHARP and it writes additionally custom instructions for each structure in a text file, which contains required data values for copy-pasting.

### 12.1.2 Filtering of datasets

First of all, the RSCB PDB [82] was searched for entries with a maximum resolution of 2.0 Å, solved via multiple wavelength anomalous diffraction (MAD) phasing, and containing only proteins as a polymer. From those, entries with less than three wavelengths, two identical wavelengths, measurements on more than one crystal, no given R-value, incomplete reflection data, minimal given experimental conditions, completeness below 97 %, and additives beyond salt and PEG in their crystallization conditions which could contribute to the signal from the solvent region, were rejected by filtering via the scripts *MAD\_data\_sets\_filter.py* and *MAD\_pdb\_stats\_filter.py*. The table output of the latter one was examined to manually identify and check entries which looked suspicious or contained undesired additives to the crystallization conditions. The respective PDB codes were added to the script as flagged entries, which were filtered out immediately and the script was run again. The next step for conversion spotted also entries with defect data sets, which were flagged as well. The flagged PDB codes contain:

- 1B9M (Crystallization conditions seemed incomplete.)
- 1UV7 (Two wavelengths are identical.)
- 1VME (Two wavelengths are identical.)
- 1VP8 (Has two low energy remote wavelengths.)
- 2I51 (Not possible to clearly identify peak and inflection wavelength.)
- 3BEM (Ni bound to surface, potential additive.)
- 3C0V (Not possible to clearly identify peak and inflection wavelength.)
- 3HTN (Ni bound to surface, potential additive.)
- 1VKN (Datasets per wavelength are not unique.)
- 1XY7 (Three wavelengths in a single dataset, which caused problems with GEMMI.)
- 2CXY (Datasets per wavelength are incomplete.)
- 4F98 (Has unusual CIF format, which is incompatible with GEMMI.)
- 4FSV (Has unusual CIF format, which is incompatible with GEMMI.)
- 5XVL (Only data for just one wavelength included.)
- 7T1M (Only data for just one wavelength included.)

Entries were also filtered if they contained flagged additives in their crystallization conditions. These flagged additives are Ni, Zn, Cacodylate, Ca, CaCl<sub>2</sub>, Mg, MgSO<sub>4</sub>, MgNO<sub>3</sub>, MgCl<sub>2</sub>, Mg formate, Mg acetate, iodide, Cd, phosphate, and the text "additive".

During the filtering, the crystallization conditions were simplified to check all entries faster manually. This involved removing the following strings from the conditions: "TEMPERATURE", "SITTING DROP", "SITTINGDROP", "NANO DROP", "NANODROP", "VAPOR DIFFUSION", "VAPORDIFFUSION", "HANGING DROP", "HANGINGDROP", and "EVAPORATION".

### 12.1.3 Conversion of data and labeling of columns

The downloaded -SF.CIF files contained a list of wavelengths used during the MAD experiment, each with a wavelength-index. Furthermore, the files contain a table with miller indices (H, K, L), intensity (I), and sigma of intensity (SIGI) for each respective wavelength-id. In order to handle this data with SHELX [167] or SHARP [159], each table must be converted into a .SCA or .MTZ file. A simple conversion from .CIF to .MTZ is accomplished by GEMMI [172] via

```
gemmi cif2mtz PDB_CODE-sf.cif --dir=OUT-DIR
```

creating an .MTZ file for table, named according to the dataset-identifier. The MTZ files can then be converted to .SCA files via PHENIX [171]. While this approach is sufficient for SHELX, it does not work for SHARP, since the columns DANO and SIGDANO must be generate, which contain anomalous differences. Therefore, PHENIX is used for data conversion. Also, all columns require the wavelength-label as suffix.

Since only wavelength-ids are stored in the .CIF files, the respective wavelength-label, i.e. “PEAK”, “INFL”, “HREM”, and “LREM”, must be identified first to make sense of the converted data. This is achieved by calculating the differences between each pair of wavelengths. The pair with the shortest difference is identified as containing “PEAK” and “INFL”. Since “PEAK” has always lower wavelength than “INFL”, both labels can be assigned to the correct wavelength id. The wavelengths for “HREM” and “LREM” are identified by their difference to “PEAK”, since “HREM” has higher energy and thus lower wavelength than “PEAK”. The opposite applies to “INFL”.

With the correct wavelength labels at hand, PHENIX is called for file conversion:

```
phenix.reflection_file_converter data/PDB_CODE-sf.cif \  
--mtz PDB_CODE_WVL_LABEL.mtz \  
--label="PDB_CODE-sf.cif:DATA_SET_NAME,wavelength_id=WVL_ID,  
_refln.intensity_meas,_refln.intensity_sigma" \  
--space_group="SPCG"
```

with PDB\_CODE as the respective PDB code, WVL\_LABEL as the identified label, DATA\_SET\_NAME as the identifier of the converted data table, WVL\_ID as the wavelength id, and SPCG as the spacegroup of the data. The commands are filled and executed automatically via a custom script as part of the pipeline. The columns of the resulting .MTZ file look like this:

```
H K L r1vpmAsf(+) SIGr1vpmAsf(+) r1vpmAsf(-) SIGr1vpmAsf(-)
```

Renaming these columns to more convenient names (I(+), SIGI(+), I(-), SIGI(-)) is achieved via the python library "gemmi" [172] as part of the custom script. As final preparation for SHARP, the columns DANO and SIGDANO are generated via CTRUNCATE [173]. The wavelength label is added as prefix to each column, in order to identify its origin after merging all columns in a future step:

```
ctruncate -mtzin PDB_CODE_WVL_LABEL.mtz \
-mtzout PDB_CODE_WVL_LABEL.mtz \
-colin '[H,K,L,I(+),SIGI(+),I(-),SIGI(-)]' \
-colout '_WVL_LABEL' \
-colano '/*/*/[I(+),SIGI(+),I(-),SIGI(-)]'
```

Since SHELXC requires .SCA files, all generated (and relabeled) .MTZ files are converted:

```
phenix.reflection_file_converter PDB_CODE_WVL_LABEL.mtz \
--sca PDB_CODE_WVL_LABEL.sca --label="IMEAN,SIGIMEAN"
```

This is done for each wavelength. All of these steps are performed within the custom script *autoMAD\_sca\_converter.py*.

### 12.1.4 Preparation via SHELXC

The SHELXC part of the pipeline generates automatically an input file for each structure within the script *autoMAD\_shelxc.py*. The content of all input files follows the same scheme, where the number of wavelengths varies between three and four. Specifically, the input files contain unit cell dimensions, the spacegroup, and a file path to the .SCA file with the respective wavelength label for each wavelength. Additionally, it contains the number of iterations and number of expected strong anomalous scatterers, which are set in all input files to 100 and 10, respectively. The exact numbers by SHELXC, but are not important here since ANODE is used to identify the strong anomalous scatterers. Below is an example of such an auto-generated input file:

```
cell 49.231 79.449 106.339 90.00 90.00 90.00
spag P 21 21 21
find 10
ntry 100
PEAK ../../reflection_data/1VJV/1VJV_PEAK.sca
```

```
INFL ../../reflection_data/1VJV/1VJV_INFL.sca
HREM ../../reflection_data/1VJV/1VJV_HREM.sca
LREM ../../reflection_data/1VJV/1VJV_LREM.sca
```

Next, the script calls SHELXC and utilized the input file, as well as the .PDB file of each structure:

```
shelxc PDB_CODE < auto-gen_shelxc_input.txt
```

The standard output is written into the textfile "PDB\_CODE\_shelxc\_out.txt" into the same folder and is checked manually in case of unexpected results. It contains statistics about the reflection data of each data set. Additionally, SHELXC generates the files "PDB\_CODE.hkl", "PDB\_CODE\_fa.hkl" and "PDB\_CODE\_fa.ins", which are used by SHELXD or ANODE. In this work, SHELXD was not used, only ANODE was used.

### 12.1.5 Substructure identification via ANODE

The PDB files are copied over to the folder of the SHELX results. Then, ANODE (version 2013/1) [161] is executed:

```
anode PDB_CODE
```

Since it may be required to rerun ANODE with the additional option "-i" to reindex for certain space groups, the standard output of anode is read to identify the run with the highest averaged anomalous density (AAD) value. That run is then executed again to generate the correct files. If the highest AAD value is below 40, a warning about the low value is given, since this likely results in low quality electron density. An electron density map is still generated for manual inspection, but the entry is filtered out and ignored in subsequent steps.

ANODE writes the standard output automatically into the "PFB\_CODE.lsa" file, the substructure into the "PDB\_CODE\_fa.res" file, and an electron density map based on anomalous differences into the "PDB\_CODE.pha" file.

Since the substructure is modified based on the AAD peaks, a copy for manual inspection is created under the name "copy\_PDB\_CODE\_fa.res". From the table of the .lsa file, substructure atoms with peak heights below 6.0 are identified and removed from the list in the "PDB\_CODE\_fa.res" file via a custom python script. The element of the nearest element is not considered. All steps are executed and handled within *autoMAD\_anode.py*, which generates also a table listing each entry with its respective AAD value, which is later used as a quality indicator for identifying best cases.



### 12.1.6 Phasing via SHELXE

After preparing the substructure, phasing of the complete structure is performed via SHELXE [158]:

```
shelxe PDB_CODE PDB_CODE_fa -m0 -e999 -h
```

where `-h` is added, if the native structure contains heavy atoms, which is true for all examined structures. The option `-mN` forces `N` iterations of density modification, where `-m0` turns density modification off. The option `-eX` augments the data via a "free lunch" approach up to a resolution of `X` Å, where `-e999` switches off the default filling.

The results are the desired phases for the entire structure stored as structure factors in the "PDB\_CODE.phs" file. The console print is stored in "PDB\_CODE.lst", while the heavy atom substructure becomes available with the "PDB\_CODE.hat" file.

Since some of the following scripts along the pipeline are not able to read .PHS files, a conversion from .PHS to .MTZ is performed immediately afterwards. It is accomplished by the CCP4 command

```
f2mtz HKLIN 'PDB_CODE.phs' HKLOUT 'PDB_CODE_from_phs.mtz' <<EOF
CELL UNIT_CELL_PARAMS
ABOUT H K L FOBS FOM PHIB SIGOBS
CTYPE H H H F W P Q
FORMAT '*'
SYMM SPACE_GROUP
END
EOF
```

where `PDB_CODE` is the entry's PDB identifier, `UNIT_CELL_PARAMS` are the unit cell parameters in the format "a b c alpha beta gamma", and `SPACE_GROUP` is the space group, such as "P 21 21 21". The rows `ABOUT` and `CTYPE` define the column labels and type, respectively, and the row `FORMAT` is always kept with "\*".

All of these steps are performed by the script *autoMAD\_shelxe.py*

### 12.1.7 Quality assessment via cross correlation

The script *autoMAD\_shelxe.py* performs after the conversion from .PHS to .MTZ a realspace cross correlation between the PDB structure and the electron density calculated from the .MTZ file. The same is repeated with the same PDB structure and the map

from the deposited .MTZ file of the PDB entry and the difference between those cross correlation values is calculated in EXCEL.

For calculation of the cross correlation, PHENIX is used with the command

```
phenix.get_cc_mtz_pdb MTZ_FILE.mtz PDB_CODE.pdb labin="LABELS"
```

where MTZ\_FILE is the path to the respective MTZ file and LABELS defines the column labels for structure factor amplitudes and phases. For the SHELXE .MTZ file these labels are "FP=FOBS PHIB=PHIB", while for the deposited .MTZ file from the PDB the labels are "FP=FP PHIB=PHIC".

From the command line output, the overall and the local cross correlation values are extracted. After running this for both .MTZ files, the values are saved to a table file. The differences between cross correlation from SHELXE and deposited data were used together with the AAD values from ANODE as indicators for bad quality in order to identify high quality electron density maps.

### 12.1.8 Automated preparation for SHARP

SHARP was operated manually, but the preparation of files was included into the automatic pipeline by writing the script *autoMAD\_prep\_for\_SHARP.py*. One requirement for SHARP is having the datasets of all wavelengths in a single MTZ file. To still identify each wavelength-specific data, the wavelength label was already attached to each column by CTRUNCATE as described in section 12.1.3. The MTZ files from that step were merged pairwise in two steps, until all data was in a single MTZ file, which was accomplished by

```
mtzutils hklin1 PDB_CODE_WVL_LABEL.mtz \  
hclin2 PDB_CODE_WVL_LABEL.mtz hklout OUT_MTZ.mtz
```

where the final MTZ was named "PDB\_CODE\_merged\_all\_mtzutils.mtz". Additionally, a text file for each structure is written, which contains relevant information such as unit cell parameters, count of atoms, wavelengths and calculated  $f'$  and  $f''$ , or atom coordinates and elements of the heavy atom substructure from ANODE. Some of this had to be generated first. A .HATOM file contains element and fractional coordinates of each heavy atom, where a single line looks like this:

```
ATOM Se 0.28356 0.14581 0.04351
```

This file is generated automatically based on atoms listed in the .LSA file generated by ANODE. However, manual inspection and confirmation is required. For atom count estimates, all atoms from the PDB file are counted inside the script, where the structure is accessed via the python library of GEMMI. Calculation of  $f'$  and  $f''$  is done per wavelength and depends on the anomalous scatterer, which is accomplished via this shell command:

```
gemmi fprime -w WAVELENGTH ELEMENT
```

where WAVELENGTH is the wavelength in Å and ELEMENT the element of the heavy atom. If multiple elements were present, the command was repeated for each.

### 12.1.9 Phasing via SHARP

Since SHARP was operated manually, only 14 structures were selected for this procedure. Criteria for the selection were high AAD values and a high quality electron density map, where the SHARP map had comparable quality to the map from the PDB on manual inspection. After selecting the best eleven cases fulfilling these criteria, additional structures were selected with the goal of covering the given spectrum of crystal solvent content. Higher resolutions were always preferred.

SHARP version 2.8.12 was used for phasing via the locally hosted web browser graphical user interface. In order to launch SHARP, CCP4 must be sourced. Here, CCP4 version 8.0 was used. Sourcing and launching SHARP is done via the commands:

```
~/sharp-2.8/setup.sh  
~/sharp-2.8/start.sh
```

Unfortunately, this crashes the terminal, although the server still runs in the background. Logging in to "http://inf2021-dub01:8000" (or another address that was given during the SHARP installation) gives access to the GUI via a web browser. The MTZ and HATOM files must be in the sharpfiles folder before launching the project to be selectable.

On the global information editor page, the calculation options were left at default and it was ensured, that  $f'$  and  $f''$  are refined in the last cycle. Both, residual maps and centroid electron density maps, were enabled to be generated. Unit cell parameters and space group were automatically imported from the MTZ file and were manually confirmed. Counts of atoms per element in the chemical composition section were calculated in the previous step from the PDB and are entered here. Heavy atoms were not entered here. Approximate numbers were not entered as these override specific numbers.

On the G-sites page, the heavy atoms were added from the .HATOM file. All coordinates were open for refinement. In the compound editor (C-sites) the elements were assigned to the previously added atom coordinates. No changes were made in the crystal editor, as only data from single crystals was used. By selecting W-1 on the left and clicking on "New", two more wavelength editors were added. These were modified according to the wavelengths and resolution limits from the experiment. Resolution limits were taken via the shell command "mtzdmp" from each individual MTZ file per wavelength. The reference wavelength, which is the first entered wavelength, was always picked as that with the highest resolution.

On the batch editor page of each wavelength, the respective columns were chosen and  $f'$  and  $f''$  were entered. The calculation of  $f'$  and  $f''$  was described in the previous step. For the start, refinement of  $f'$  and  $f''$  was disabled. For the reference wavelength, the scaling parameters K and B6 were enabled for estimates, and the global non-isomorphism and model imperfection parameters on anomalous differences were enabled for refinement. For the other wavelengths, the scaling parameters K and B6 were enabled for refinement and estimates, while the isotropic B-factor was disabled for refinement and estimates. Global non-isomorphism and model imperfection parameters were enabled for refinement on isomorphous and anomalous differences. Values were in all cases set to the default values.

After setting all values, the job was submitted and results were accessed from the results page reachable from the home page. Phase improvement was skipped. Residual map analysis was performed for each case, where plausible suggestions for refinement and modification of parameters for additional runs were accepted. If  $f'$  and  $f''$  were not calculated but obtained from experimental measurement, they were not refined. Otherwise, refinement of these values was enabled if suggested so. Switching to anisotropic B factor refinement was also accepted. New C-sites were most of the time not close to heavy atoms in the PDB file and where thus rejected most of the time. Resulting electron density maps were inspected in Coot [150]. In any case, at least a second run was performed for refinement. If the result improved, additional runs were conducted, each after refinement suggestions from residual map analysis.

## 12.2 Comparison of electron density maps

Modelled waters from the PDB structures were examined automatically for peaks in electron density of multiple maps. For the comparison of the solvent region, each electron density map was separated into protein region only and deep solvent only. The separated

maps, as well as the complete map, were transformed into reciprocal space where analysis and comparison of various properties was performed.

In addition to the electron density map from the PDB and those generated via SHELXE and SHARP, extreme cases with flat solvent and crystalline solvent were generated for each dataset. In the case of 2QPX, maps obtained from molecular dynamics simulations were considered in the comparison as well.

The following python scripts were programmed and used for this part of the work:

1. *modelled\_water\_peak\_density.py*
2. *water\_removal\_from\_pdb.py*
3. *autoMAD\_coot\_map\_separator.py*
4. *autoMAD\_mtz\_cc.py*

The functionality of these scripts is explained in the next subsections.

### 12.2.1 Surface water comparison

The script *modelled\_water\_peak\_density.py* was used to analyze the modelled surface waters. This script used the python library of GEMMI [172] to iterate over each modelled water and compare the electron density values of three maps at that position. In the first run, the maximum value per map was identified (maximum among the values at water positions, not of the entire map). In a second run, the normalized values between the PDB map, SHELXE map, and SHARP map were compared. Here, the difference between two maps was calculated and if it was above a threshold of 0.1, the respective water was counted as an outlier for this pair of maps. This was repeated for all map-pairs and that way, water common in all three maps and in a pair of maps were identified. A new PDB file was saved for each structure, where these waters were kept, while all others were deleted.

### 12.2.2 Separation of protein and solvent region

The PDB maps, SHELXE maps, and SHARP maps were compared as entire maps and additionally as partial maps. In the script *water\_removal\_from\_pdb.py*, all modelled waters and alternative conformations were removed from the structure and all occupancies were set to 1, all via GEMMI [172].

In *autoMAD\_coot\_map\_separator.py*, a coot script was generated with all required paths and settings. This script was launched without graphics via the python library "subprocess", which executed the command line command:

```
coot --no-graphics --script SCRIPT.py
```

where "SCRIPT" is the generated coot script written in python. The coot script loaded the PDB file via

```
molecule = read_pdb(PDB-FILE-PATH)
```

For SHELXE, the following code was used to load the map from PHS file:

```
eden_map = read_phs_and_make_map_using_cell_symm_from_mol(PHS, 0)
```

where the value 0 is the molecule ID providing the unit cell parameters. For loading the SHARP and PDB map from MTZ, the following code was used:

```
set_auto_read_column_labels(COL_F, COL_PHI, 0)
set_auto_read_do_difference_map_too(0)
eden_map = auto_read_make_and_draw_maps_from_mtz(MTZ)
```

where COL\_F and COL\_PHI are amplitude column label and phase column label, respectively, and MTZ is the path to the MTZ file. By defining the labels, the correct columns were loaded with the third command. The second command disabled loading a difference map. Because this did not work for the MTZ deposited to the PDB, the map ID for subsequent steps had to be increased by 2 to mask the correct map. The two zero values are the molecule ID. The loaded maps were then masked with the PDB model:

```
# map_mol_no, coords_mol_no, mmdb_atom_selection, invert_flag
# solvents map
set_map_mask_atom_radius(2.5)
mask_map_by_atom_selection(map_id, 0, "", 0)

# protein map
set_map_mask_atom_radius(2.5)
mask_map_by_atom_selection(map_id, 0, "", 1)

# empty solvents map
set_map_mask_atom_radius(2.5)
mask_map_by_atom_selection(map_id + 2, 0, "", 0)
```

where a masking radius of 2.5 Å was used. The atom selection was set to empty, which selected all atoms of the model. Since only one model was loaded, the molecule ID was 0. To gain an empty map with the correct unit cell parameters, the protein map was masked again without inverted mask, resulting in a map with all gridpoints of the unit cell set to zero. Lastly, the maps were exported:

```
export_map(map_id + 0, "full.map")
export_map(map_id + 1, "solv.map")
export_map(map_id + 2, "prot.map")
export_map(map_id + 3, "solv_empty.map")

coot_real_exit(0)
```

where the map ID was adjusted. Exported file names contained the PDB code, a suffix containing the masking radius, and an indication about which map it is.

Additionally, the solvent region was split into surface solvent region and deep solvent region. This was accomplished via the following masking:

```
# map_mol_no, coords_mol_no, mmdb_atom_selection, invert_flag
# protein + surface solvent region
set_map_mask_atom_radius(7.5)
mask_map_by_atom_selection(map_id, 0, "", 1)

# only surface solvent region
set_map_mask_atom_radius(2.5)
mask_map_by_atom_selection(map_id + 4, 0, "", 0)

# deep solvent region
set_map_mask_atom_radius(7.5)
mask_map_by_atom_selection(map_id, 0, "", 0)
```

where first a larger region around the protein was kept. This map with protein and surface solvent was masked again removing the protein region. To create the surface solvent and deep solvent maps, a masking radius of 7.5 Å was used.

### 12.2.3 Cross correlation of partial maps

Cross correlation between partial maps was calculated via PHENIX [171] within the script *autoMAD\_mtz\_cc.py*, with the command:

```
phenix.get_cc_map_map MAP1 MAP2
```

where MAP1 and MAP2 are the two maps which are compared. The script extracted the cross correlation values and stored them in a table. Cross correlations were calculated between all full maps, between all protein-only maps, and between all solvent-only maps.

### 12.3 Refinement against electron density map from SHARP

The refinement was performed in CCP4i2 [174] and involved these steps:

1. import PDB file (no waters or only common waters)
2. import MTZ from PDB, load the Free R set
3. import MTZ from SHARP, no Free R set
4. REFMAC5, refine PDB against SHARP map
5. REFMAC5, refine PDB against SHARP map + add waters
6. manual inspection and addition of waters in COOT
7. additional refinement cycles with REFMAC5

Per structure, two initial conditions were refined: The PDB structure with all waters deleted or with all waters deleted except those, which were common among the PDB, SHELXE, and SHARP map (see section 32). The PDB MTZ was only loaded to provide the same Free R set, to make valid comparison. Resolution cutoff was disabled. The SHARP MTZ was used to generate the map against which the structure model was refined.

All runs of REFMAC5 [175] used the initial (or in later cycles refined) PDB model, the reflections and phases from SHARP MTZ (columns FB and PHIB on initial imports), and the Free R set from PDB MTZ, and were run in restrained refinement mode, used anisotropic B-factors, simple solvent scaling with no solvent mask, and no translation-libration-screw parameters.

The first run of REFMAC5 (step 4) ran 10 cycles, used hydrogens during refinement with generation of riding hydrogens, no waters were added, and as the only restraint a weight restraint against the experimental data, which was the automatically calculated one. It was rerun with suggested weight restraint parameter and cycle number and the better run was kept.



The second run of REMAC5 (step 5) used the refined phases, ran 10 cycles, used hydrogens during refinement with generation of riding hydrogens, added waters if R was 0.4 or lower, and used the suggested weight restraint. It was again rerun with suggested weight restraint and cycle number, where the best case was kept.

Afterwards, model and map were inspected manually in COOT [150], where automatically placed water not agreeing with the map were deleted. The functions "find unmodelled blobs" and "difference map peaks" (with RMSD above 8.0) were used to eventually add more waters manually. Suspicious waters were removed and suggested changes from COOT were considered. The following refinement steps depended on the results and aimed at lowering the R values, while also increasing the verdict score and also trying not to overfit. Waters were only added manually in COOT. Calculated phases from refinements were not used in any case. Only phases from the SHARP MTZ were used. After no further improvements were observable, additional refinement with the use of explicit solvent masks was applied, followed by manual water adding in COOT and a last round of REFMAC5.

## 12.4 Analysis of solvent region contribution

The entire electron density map, as well as all partial maps, were converted to MTZ files via GEMMI [172] for each structure, which was automatically done with the script *hkl\_stats\_generator.py* within the script *reflection\_filter.py*. Then, the script iterates over each reflection, where the structure factors are represented as complex numbers with the amplitude as the real part and phase angle in radians as the complex part. The complex numbers of partial maps were added and the magnitude of their sum was divided by the magnitude of the structure factor from the entire map to obtain the scaling factor of that reflection. If it was beyond a tolerance of 0.001, it was counted as a scaling factor outlier. Reflections with magnitudes below 1 were also counted.

The solvent contribution of a structure factor was calculated by dividing the magnitude of the solvent region structure factor by the sum of magnitudes of the solvent region structure factor and the protein region structure factor. Plots were generated via the scripts *reflection\_filter.py* and *autoMAD\_plot\_gen.py*, by collecting various statistics per reflection.

## References

1. Lamers, M. M. & Haagmans, B. L. SARS-CoV-2 pathogenesis. *Nature reviews microbiology* **20**, 270–284 (2022).
2. Huang, C. *et al.* A new generation Mpro inhibitor with potent activity against SARS-CoV-2 Omicron variants. *Signal Transduction and Targeted Therapy* **8**, 128 (2023).
3. She, Z. *et al.* Mpro-targeted anti-SARS-CoV-2 inhibitor-based drugs. *Journal of Chemical Research* **47**, 17475198231184799 (2023).
4. Jiang, H., Yang, P. & Zhang, J. Potential inhibitors targeting papain-like protease of SARS-CoV-2: two birds with one stone. *Frontiers in chemistry* **10**, 822785 (2022).
5. Calleja, D. J., Lessene, G. & Komander, D. Inhibitors of sars-cov-2 plpro. *Frontiers in Chemistry* **10**, 876212 (2022).
6. Chakraborty, C., Sharma, A. R., Bhattacharya, M. & Lee, S.-S. A detailed overview of immune escape, antibody escape, partial vaccine escape of SARS-CoV-2 and their emerging variants with escape mutations. *Frontiers in immunology* **13**, 801522 (2022).
7. Harvey, W. T. *et al.* SARS-CoV-2 variants, spike mutations and immune escape. *Nature Reviews Microbiology* **19**, 409–424 (2021).
8. Mengist, H. M. *et al.* Mutations of SARS-CoV-2 spike protein: Implications on immune evasion and vaccine-induced immunity in *Seminars in immunology* **55** (2021), 101533.
9. Gonzalez Lomeli, F., Elmaraghy, N., Castro, A., Osuna Guerrero, C. V. & Newcomb, L. L. Conserved targets to Prevent emerging coronaviruses. *Viruses* **14**, 563 (2022).
10. Melo-Filho, C. C. *et al.* Conserved coronavirus proteins as targets of broad-spectrum antivirals. *Antiviral Research* **204**, 105360 (2022).
11. Percha, B. & Altman, R. B. Informatics confronts drug–drug interactions. *Trends in pharmacological sciences* **34**, 178–184 (2013).
12. Marsilio, N. R., Silva, D. d. & Bueno, D. Drug incompatibilities in the adult intensive care unit of a university hospital. *Revista Brasileira de terapia intensiva* **28**, 147–153 (2016).

13. Von Soosten, L. C. *et al.* The Swiss army knife of SARS-CoV-2: the structures and functions of NSP3. en. *Crystallography Reviews* **28**, 39–61. ISSN: 0889-311X, 1476-3508 (Jan. 2, 2022).
14. Croll, T. *et al.* Making the invisible enemy visible. en. *bioRxiv*, 2020.10.07.307546 (Oct. 7, 2020).
15. Jumper, J. *et al.* Highly accurate protein structure prediction with AlphaFold. en. *Nature* **596**, 583–589. ISSN: 0028-0836, 1476-4687 (Aug. 26, 2021).
16. Jumper, J. *et al.* Applying and improving AlphaFold at CASP14. *Proteins: Structure, Function, and Bioinformatics* **89**, 1711–1721 (2021).
17. Lei, J., Kusov, Y. & Hilgenfeld, R. Nsp3 of coronaviruses: Structures and functions of a large multi-domain protein. en. *Antiviral Research* **149**, 58 (Jan. 2018).
18. Neuman, B. W. Bioinformatics and functional analyses of coronavirus nonstructural proteins involved in the formation of replicative organelles. en. *Antiviral Research* **135**, 97–107. ISSN: 01663542 (Nov. 2016).
19. Wolff, G. *et al.* A molecular pore spans the double membrane of the coronavirus replication organelle. en. *Science* **369**, 1395–1398. ISSN: 0036-8075, 1095-9203 (Sept. 11, 2020).
20. Piovesan, D., Monzon, A. M. & Tosatto, S. C. *Intrinsic Protein Disorder, Conditional Folding and AlphaFold2* en. Tech. rep. DOI: 10.1101/2022.03.03.482768 (Mar. 3, 2022). <http://biorxiv.org/lookup/doi/10.1101/2022.03.03.482768>.
21. Grellet, E., Goulet, A. & Imbert, I. Replication of the coronavirus genome: A paradox among positive-strand RNA viruses. *Journal of Biological Chemistry* **298** (2022).
22. Oliveira, G. P. & Kroon, E. G. Mouse hepatitis virus: A betacoronavirus model to study the virucidal activity of air disinfection equipment on surface contamination. *Journal of Virological Methods* **297**, 114274 (2021).
23. Pizzato, M. *et al.* SARS-CoV-2 and the host cell: a tale of interactions. *Frontiers in Virology* **1** (2022).
24. Emrani, J. *et al.* SARS-COV-2, infection, transmission, transcription, translation, proteins, and treatment: A review. *International journal of biological macromolecules* **193**, 1249–1273 (2021).

25. Osipiuk, J. *et al.* Structure of papain-like protease from SARS-CoV-2 and its complexes with non-covalent inhibitors. en. *Nature Communications* **12**, 743. ISSN: 2041-1723 (Feb. 2, 2021).
26. Lee, J. *et al.* X-ray crystallographic characterization of the SARS-CoV-2 main protease polyprotein cleavage sites essential for viral processing and maturation. *Nature Communications* **13**, 5196 (2022).
27. Hagemeyer, M. C. *et al.* Membrane rearrangements mediated by coronavirus non-structural proteins 3 and 4. en. *Virology* **458-459**, 125–135. ISSN: 00426822 (June 2014).
28. Kirchdoerfer, R. N. & Ward, A. B. Structure of the SARS-CoV nsp12 polymerase bound to nsp7 and nsp8 co-factors. *Nature communications* **10**, 2342 (2019).
29. Malone, B., Urakova, N., Snijder, E. J. & Campbell, E. A. Structures and functions of coronavirus replication–transcription complexes and their relevance for SARS-CoV-2 drug design. *Nature Reviews Molecular Cell Biology* **23**, 21–39 (2022).
30. Van Hemert, M. J. *et al.* SARS-Coronavirus Replication/Transcription Complexes Are Membrane-Protected and Need a Host Factor for Activity In Vitro. en. *PLoS Pathogens* **4** (ed Baric, R. S.) e1000054. ISSN: 1553-7374 (May 2, 2008).
31. De Haan, C. A. & Rottier, P. J. en. in *Advances in Virus Research* DOI: 10.1016/S0065-3527(05)64006-7, 165–230 (Elsevier, 2005). ISBN: 978-0-12-039863-8. <https://linkinghub.elsevier.com/retrieve/pii/S0065352705640067>.
32. Zimmermann, L. *et al.* SARS-CoV-2 nsp3 and nsp4 are minimal constituents of a pore spanning replication organelle. *Nature Communications* **14**, 7894 (2023).
33. Bessa, L. M. *et al.* The intrinsically disordered SARS-CoV-2 nucleoprotein in dynamic complex with its viral partner nsp3a. en. *Science Advances* **8**, eabm4034. ISSN: 2375-2548 (Jan. 21, 2022).
34. Chang, C.-k., Hou, M.-H., Chang, C.-F., Hsiao, C.-D. & Huang, T.-h. The SARS coronavirus nucleocapsid protein–forms and functions. *Antiviral research* **103**, 39–50 (2014).
35. McBride, R., Van Zyl, M. & Fielding, B. C. The coronavirus nucleocapsid is a multifunctional protein. *Viruses* **6**, 2991–3018 (2014).
36. Scherer, K. M. *et al.* SARS-CoV-2 nucleocapsid protein adheres to replication organelles before viral assembly at the Golgi/ERGIC and lysosome-mediated egress. *Science advances* **8**, eabl4895 (2022).

37. Vennema, H. *et al.* Nucleocapsid-independent assembly of coronavirus-like particles by co-expression of viral envelope protein genes. *The EMBO journal* **15**, 2020–2028 (1996).
38. Zhang, Z. *et al.* Structure of SARS-CoV-2 membrane protein essential for virus assembly. *Nature communications* **13**, 4399 (2022).
39. Lu, S. *et al.* The SARS-CoV-2 nucleocapsid phosphoprotein forms mutually exclusive condensates with RNA and the membrane-associated M protein. *Nature communications* **12**, 502 (2021).
40. Westerbeck, J. W. & Machamer, C. E. A coronavirus E protein is present in two distinct pools with different effects on assembly and the secretory pathway. *Journal of virology* **89**, 9313–9323 (2015).
41. Ghosh, S. *et al.*  $\beta$ -Coronaviruses use lysosomes for egress instead of the biosynthetic secretory pathway. *Cell* **183**, 1520–1535 (2020).
42. Oostra, M. *et al.* Topology and Membrane Anchoring of the Coronavirus Replication Complex: Not All Hydrophobic Domains of nsp3 and nsp6 Are Membrane Spanning. en. *Journal of Virology* **82**, 12392–12405. ISSN: 0022-538X, 1098-5514 (Dec. 15, 2008).
43. Shin, D. *et al.* Papain-like protease regulates SARS-CoV-2 viral spread and innate immunity. *Nature* **587**, 657–662 (2020).
44. Alhammad, Y. M. & Fehr, A. R. The viral macrodomain counters host antiviral ADP-ribosylation. *Viruses* **12**, 384 (2020).
45. Rack, J. G. M. *et al.* Viral macrodomains: a structural and evolutionary assessment of the pharmacological potential. *Open biology* **10**, 200237 (2020).
46. Michalska, K. *et al.* Crystal structures of SARS-CoV-2 ADP-ribose phosphatase: from the apo form to ligand complexes. *IUCrJ* **7**, 814–824 (2020).
47. Hurst, K. R., Koetzner, C. A. & Masters, P. S. Characterization of a critical interaction between the coronavirus nucleocapsid protein and nonstructural protein 3 of the viral replicase-transcriptase complex. *Journal of virology* **87**, 9159–9172 (2013).
48. Chou, C.-C. & Wang, A. H.-J. Structural D/E-rich repeats play multiple roles especially in gene regulation through DNA/RNA mimicry. *Molecular BioSystems* **11**, 2144–2151 (2015).

49. Ni, X., Han, Y., Zhou, R., Zhou, Y. & Lei, J. Structural insights into ribonucleoprotein dissociation by nucleocapsid protein interacting with non-structural protein 3 in SARS-CoV-2. *Communications Biology* **6**, 193 (2023).
50. Rack, J. G. M. *et al.* Viral macrodomains: a structural and evolutionary assessment of the pharmacological potential. en. *Open Biology* **10**, 200237. ISSN: 2046-2441 (Nov. 2020).
51. Michalska, K. *et al.* Crystal structures of SARS-CoV-2 ADP-ribose phosphatase: from the apo form to ligand complexes. *IUCrJ* **7**, 814–824. ISSN: 2052-2525 (Sept. 1, 2020).
52. Kusov, Y., Tan, J., Alvarez, E., Enjuanes, L. & Hilgenfeld, R. A G-quadruplex-binding macrodomain within the “SARS-unique domain” is essential for the activity of the SARS-coronavirus replication–transcription complex. *Virology* **484**, 313–322 (2015).
53. Taha, T. Y. *et al.* A single inactivating amino acid change in the SARS-CoV-2 NSP3 Mac1 domain attenuates viral replication in vivo. *PLoS Pathogens* **19**, e1011614 (2023).
54. Neuman, B. W., Chamberlain, P., Bowden, F. & Joseph, J. Atlas of coronavirus replicase structure. *Virus research* **194**, 49–66 (2014).
55. Chen, Y. *et al.* X-ray Structural and Functional Studies of the Three Tandemly Linked Domains of Non-structural Protein 3 (nsp3) from Murine Hepatitis Virus Reveal Conserved Functions. en. *Journal of Biological Chemistry* **290**, 25293–25306. ISSN: 00219258 (Oct. 2015).
56. Tan, J. *et al.* The SARS-unique domain (SUD) of SARS coronavirus contains two macrodomains that bind G-quadruplexes. *PLoS pathogens* **5**, e1000428 (2009).
57. Neuman, B. W. *et al.* Proteomics analysis unravels the functional repertoire of coronavirus nonstructural protein 3. *Journal of virology* **82**, 5279–5294 (2008).
58. Stogios, P. *et al.* Crystal structure of the Nsp3 bSM (Betacoronavirus-Specific Marker) domain from SARS-CoV-2 Dec. 20, 2021. <https://doi.org/10.2210/pdb7t9w/pdb>.
59. Stogios, P. *et al.* Crystal structure of the Nsp3 Y3 domain from SARS-CoV-2 Aug. 6, 2021. <https://doi.org/10.2210/pdb7rqg/pdb>.
60. Li, Y. *et al.* Crystal structure of the CoV-Y domain of SARS-CoV-2 nonstructural protein 3. en. *Scientific Reports* **13**, 2890. ISSN: 2045-2322 (Feb. 18, 2023).

61. Harak, C. & Lohmann, V. Ultrastructure of the replication sites of positive-strand RNA viruses. *Virology* **479**, 418–433 (2015).
62. Wolff, G., Melia, C. E., Snijder, E. J. & Bárcena, M. Double-membrane vesicles as platforms for viral replication. *Trends in microbiology* **28**, 1022–1033 (2020).
63. Ujike, M. & Taguchi, F. Recent Progress in Torovirus Molecular Biology. en. *Viruses* **13**, 435. ISSN: 1999-4915 (Mar. 8, 2021).
64. Ertel, K. J. *et al.* Cryo-electron tomography reveals novel features of a viral RNA replication compartment. *Elife* **6**, e25940 (2017).
65. Gadlage, M. J. *et al.* Murine hepatitis virus nonstructural protein 4 regulates virus-induced membrane modifications and replication complex function. *Journal of virology* **84**, 280–290 (2010).
66. Pettersen, E. F. *et al.* UCSF Chimera—A visualization system for exploratory research and analysis. en. *Journal of Computational Chemistry* **25**, 1605–1612. ISSN: 0192-8651, 1096-987X (Oct. 2004).
67. Kandler, L. *et al.* SARS-CoV-2 envelope protein and its relationship to the membrane protein. *Crystallography Reviews* **29**, 128–150 (2023).
68. De Haan, C. A., Vennema, H. & Rottier, P. J. Assembly of the coronavirus envelope: homotypic interactions between the M proteins. *Journal of virology* **74**, 4967–4978 (2000).
69. Gu, J. & Bourne, P. E. *Structural bioinformatics* (John Wiley & Sons, 2009).
70. Grigoriev, I. V. & Kim, S.-H. Detection of protein fold similarity based on correlation of amino acid properties. *Proceedings of the National Academy of Sciences* **96**, 14318–14323 (1999).
71. Laurents, D., Subbiah, S. & Levitt, M. Different protein sequences can give rise to highly similar folds through different stabilizing interactions. *Protein Science* **3**, 1938–1944 (1994).
72. Venkateswaran, J., Song, B., Kahveci, T. & Jermaine, C. Trial: a tool for finding distant structural similarities. *IEEE/ACM Transactions on Computational Biology and Bioinformatics* **8**, 819–831 (2009).
73. Marsden, R. L. *et al.* Exploiting protein structure data to explore the evolution of protein function and biological complexity. *Philosophical Transactions of the Royal Society B: Biological Sciences* **361**, 425–440 (2006).

74. Madeira, F. *et al.* Search and sequence analysis tools services from EMBL-EBI in 2022. en. *Nucleic Acids Research* **50**, W276–W279. ISSN: 0305-1048, 1362-4962 (July 5, 2022).
75. Polyanovsky, V. O., Roytberg, M. A. & Tumanyan, V. G. Comparative analysis of the quality of a global algorithm and a local algorithm for alignment of two sequences. *Algorithms for molecular biology* **6**, 1–12 (2011).
76. Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local alignment search tool. en. *Journal of Molecular Biology* **215**, 403–410. ISSN: 00222836 (Oct. 1990).
77. Sayers, E. W. *et al.* Database resources of the national center for biotechnology information. *Nucleic acids research* **49**, D10 (2021).
78. Brown, C. J., Johnson, A. K., Dunker, A. K. & Daughdrill, G. W. Evolution and disorder. en. *Current Opinion in Structural Biology* **21**, 441–446. ISSN: 0959440X (June 2011).
79. Dunker, A. K., Silman, I., Uversky, V. N. & Sussman, J. L. Function and structure of inherently disordered proteins. *Current opinion in structural biology* **18**, 756–764 (2008).
80. Schlessinger, A. *et al.* Protein disorder—a breakthrough invention of evolution? *Current opinion in structural biology* **21**, 412–418 (2011).
81. Wang, Y., Zhang, H., Zhong, H. & Xue, Z. Protein domain identification methods and online resources. *Computational and structural biotechnology journal* **19**, 1145–1153 (2021).
82. Berman, H. M. The Protein Data Bank. *Nucleic Acids Research* **28**, 235–242. ISSN: 13624962 (Jan. 1, 2000).
83. Ruggiero, A., Smaldone, G., Squeglia, F. & Berisio, R. Enhanced crystallizability by protein engineering approaches: a general overview. *Protein and peptide letters* **19**, 732–742 (2012).
84. Quevillon-Cheruel, S., Leulliot, N., Gentils, L., van Tilbeurgh, H. & Poupon, A. Production and crystallization of protein domains: how useful are disorder predictions? *Current Protein and Peptide Science* **8**, 151–160 (2007).
85. Schafer, J. W. & Porter, L. L. Evolutionary selection of proteins with two folds. *Biophysical Journal* **122**, 474a (2023).



86. Gräwert, M. & Svergun, D. A beginner's guide to solution small-angle X-ray scattering (SAXS). *The Biochemist* **42**, 36–42 (2020).
87. Anfinsen, C. B. Principles that govern the folding of protein chains. *Science* **181**, 223–230 (1973).
88. Dill, K. A., Ozkan, S. B., Shell, M. S. & Weikl, T. R. The protein folding problem. *Annu. Rev. Biophys.* **37**, 289–316 (2008).
89. Swope, W. C., Pitera, J. W. & Suits, F. Describing protein folding kinetics by molecular dynamics simulations. 1. Theory. *The Journal of Physical Chemistry B* **108**, 6571–6581 (2004).
90. Piana, S., Klepeis, J. L. & Shaw, D. E. Assessing the accuracy of physical models used in protein-folding simulations: quantitative evidence from long molecular dynamics simulations. *Current opinion in structural biology* **24**, 98–105 (2014).
91. Miao, Y., Feixas, F., Eun, C. & McCammon, J. A. Accelerated molecular dynamics simulations of protein folding. *Journal of computational chemistry* **36**, 1536–1549 (2015).
92. Sippl, M. J. Boltzmann's principle, knowledge-based mean fields and protein folding. An approach to the computational determination of protein structures. *Journal of computer-aided molecular design* **7**, 473–501 (1993).
93. Mackenzie, C. O., Zhou, J. & Grigoryan, G. Tertiary alphabet for the observable protein structural universe. *Proceedings of the National Academy of Sciences* **113**, E7438–E7447 (2016).
94. Frappier, V., Jenson, J. M., Zhou, J., Grigoryan, G. & Keating, A. E. Tertiary structural motif sequence statistics enable facile prediction and design of peptides that bind anti-apoptotic Bfl-1 and Mcl-1. *Structure* **27**, 606–617 (2019).
95. Kuhlman, B. & Bradley, P. Advances in protein structure prediction and design. *Nature reviews molecular cell biology* **20**, 681–697 (2019).
96. Gu, J. *et al.* Recent advances in convolutional neural networks. *Pattern recognition* **77**, 354–377 (2018).
97. Li, Z., Liu, F., Yang, W., Peng, S. & Zhou, J. A survey of convolutional neural networks: analysis, applications, and prospects. *IEEE transactions on neural networks and learning systems* **33**, 6999–7019 (2021).
98. Wei, G.-W. Protein structure prediction beyond AlphaFold. *Nature Machine Intelligence* **1**, 336–337 (2019).

99. AlQuraishi, M. AlphaFold at CASP13. *Bioinformatics* **35**, 4862–4865 (2019).
100. Baek, M. *et al.* Accurate prediction of protein structures and interactions using a three-track neural network. *Science* **373**, 871–876 (2021).
101. Lin, Z. *et al.* Language models of protein sequences at the scale of evolution enable accurate structure prediction. *BioRxiv* **2022**, 500902 (2022).
102. Varadi, M. *et al.* AlphaFold Protein Structure Database: massively expanding the structural coverage of protein-sequence space with high-accuracy models. *Nucleic acids research* **50**, D439–D444 (2022).
103. Evans, R. *et al.* Protein complex prediction with AlphaFold-Multimer. en. [Online; accessed 2022-04-05]. <http://biorxiv.org/lookup/doi/10.1101/2021.10.04.463034> (Oct. 4, 2021).
104. Vaswani, A. *et al.* Attention is all you need. *Advances in neural information processing systems* **30** (2017).
105. Akdel, M. *et al.* A structural biology community assessment of AlphaFold 2 applications. en. [Online; accessed 2022-04-05]. <http://biorxiv.org/lookup/doi/10.1101/2021.09.26.461876> (Sept. 26, 2021).
106. Palmer III, A. G. Enzyme dynamics from NMR spectroscopy. *Accounts of Chemical Research* **48**, 457–465 (2015).
107. Costa, R. G. L. & Fushman, D. Reweighting methods for elucidation of conformation ensembles of proteins. *Current opinion in structural biology* **77**, 102470 (2022).
108. Camacho-Zarco, A. R. *et al.* NMR provides unique insight into the functional dynamics and interactions of intrinsically disordered proteins. *Chemical Reviews* **122**, 9331–9356 (2022).
109. Sala, D., Engelberger, F., Mchaourab, H. & Meiler, J. Modeling conformational states of proteins with AlphaFold. *Current Opinion in Structural Biology* **81**, 102645 (2023).
110. Saldaño, T. *et al.* Impact of protein conformational diversity on AlphaFold predictions. *Bioinformatics* **38**, 2742–2748 (2022).
111. Wayment-Steele, H. K. *et al.* Predicting multiple conformations via sequence clustering and AlphaFold2. *Nature* **625**, 832–839 (2024).
112. Lomize, A. L. *et al.* Membranome 3.0: Database of single-pass membrane proteins with AlphaFold models. *Protein Science* **31**, e4318 (2022).

113. Minor, W., Dauter, Z. & Jaskolski, M. The young person's guide to the PDB. *Postepy biochemii* **62**, 242 (2016).
114. Edich, M., Briggs, D. C., Kippes, O., Gao, Y. & Thorn, A. The impact of AlphaFold on experimental structure solution. *Faraday Letters*, invited, --in review.
115. Terwilliger, T. C. *et al.* AlphaFold predictions are valuable hypotheses and accelerate but do not replace experimental structure determination. en. *Nature Methods* **21**, 110–116. ISSN: 1548-7091, 1548-7105 (Jan. 2024).
116. Lasker, K. *et al.* Integrative structure modeling of macromolecular assemblies from proteomics data. *Molecular & Cellular Proteomics* **9**, 1689–1702 (2010).
117. Lasker, K. *et al.* Integrative structure modeling of macromolecular assemblies from proteomics data. *Molecular & Cellular Proteomics* **9**, 1689–1702 (2010).
118. Lasker, K., Topf, M., Sali, A. & Wolfson, H. J. Inferential optimization for simultaneous fitting of multiple components into a CryoEM map of their assembly. *Journal of molecular biology* **388**, 180–194 (2009).
119. Rantos, V., Karius, K. & Kosinski, J. Integrative structural modeling of macromolecular complexes using Assemblin. *Nature Protocols* **17**, 152–176 (2022).
120. Kosinski, J. *et al.* Xlink Analyzer: software for analysis and visualization of cross-linking data in the context of three-dimensional structures. *Journal of structural biology* **189**, 177–183 (2015).
121. O'Reilly, F. J. & Rappsilber, J. Cross-linking mass spectrometry: methods and applications in structural, molecular and systems biology. *Nature structural & molecular biology* **25**, 1000–1008 (2018).
122. Sinz, A. Cross-linking/mass spectrometry for studying protein structures and protein–protein interactions: where are we now and where should we go from here? *Angewandte Chemie International Edition* **57**, 6390–6396 (2018).
123. Zvelebil, M. J. & Sternberg, M. J. Analysis and prediction of the location of catalytic residues in enzymes. *Protein Engineering, Design and Selection* **2**, 127–138 (1988).
124. Guzenko, D., Burley, S. K. & Duarte, J. M. Real time structural search of the Protein Data Bank. *PLoS computational biology* **16**, e1007970 (2020).
125. Novotni, M. & Klein, R. Shape retrieval using 3D Zernike descriptors. *Computer-Aided Design* **36**, 1047–1062 (2004).
126. Izidoro, S. C., de Melo-Minardi, R. C. & Pappa, G. L. GASS: identifying enzyme active sites with genetic algorithms. *Bioinformatics* **31**, 864–870 (2015).

127. Lavigne, M. *et al.* SARS-CoV-2 Nsp3 unique domain SUD interacts with guanine quadruplexes and G4-ligands inhibit this interaction. *Nucleic acids research* **49**, 7695–7712 (2021).
128. Mirdita, M. *et al.* ColabFold - Making Protein folding accessible to all. *bioRxiv* (2021).
129. Keane, S. C. & Giedroc, D. P. Solution Structure of Mouse Hepatitis Virus (MHV) nsp3a and Determinants of the Interaction with MHV Nucleocapsid (N) Protein. en. *Journal of Virology* **87**, 3502–3515. ISSN: 0022-538X, 1098-5514 (Mar. 15, 2013).
130. Serrano, P. *et al.* Nuclear Magnetic Resonance Structure of the N-Terminal Domain of Nonstructural Protein 3 from the Severe Acute Respiratory Syndrome Coronavirus. en. *Journal of Virology* **81**, 12049–12060. ISSN: 0022-538X, 1098-5514 (Nov. 2007).
131. Schrödinger, L. L. C. "The PyMOL Molecular Graphics System, Version 2.3. 2019."
132. Williams, C., Richardson, D. C. & Richardson, J. S. Extreme backbone outlier patterns when AlphaFold gives up. *Computational Crystallography Newsletter* **13**, 7–12 (2022).
133. Shindyalov, I. N. & Bourne, P. E. Protein structure alignment by incremental combinatorial extension (CE) of the optimal path. *Protein engineering* **11**, 739–747 (1998).
134. Panjkovich, A. & Svergun, D. I. Deciphering conformational transitions of proteins by small angle X-ray scattering and normal mode analysis. en. *Physical Chemistry Chemical Physics* **18**, 5707–5719. ISSN: 1463-9076, 1463-9084 (2016).
135. Andreini, C., Arnesano, F. & Rosato, A. The zinc proteome of SARS-CoV-2. en. *Metallomics* **14**, mfac047. ISSN: 1756-591X (July 25, 2022).
136. Wilbur, W. J. On the PAM matrix model of protein evolution. *Molecular biology and evolution* **2**, 434–447 (1985).
137. Brandolini, M. *et al.* Correlating qRT-PCR, dPCR and viral titration for the identification and quantification of SARS-CoV-2: a new approach for infection management. *Viruses* **13**, 1022 (2021).
138. Lapetina, S. & Gil-Henn, H. A guide to simple, direct, and quantitative in vitro binding assays. *Journal of biological methods* **4** (2017).

139. Zhang, Y. *et al.* *Applications of AlphaFold beyond Protein Structure Prediction* en. Tech. rep. DOI: 10.1101/2021.11.03.467194 (Nov. 4, 2021). <http://biorxiv.org/lookup/doi/10.1101/2021.11.03.467194>.
140. Buel, G. R. & Walters, K. J. Can AlphaFold2 predict the impact of missense mutations on structure? en. *Nature Structural & Molecular Biology* **29**, 1–2. ISSN: 1545-9993, 1545-9985 (Jan. 2022).
141. Pak, M. A. *et al.* Using AlphaFold to predict the impact of single mutations on protein stability and function. en. *PLOS ONE* **18** (ed Budisa, N.) e0282689. ISSN: 1932-6203 (Mar. 16, 2023).
142. Serrano, P. *et al.* Nuclear magnetic resonance structure of the N-terminal domain of nonstructural protein 3 from the severe acute respiratory syndrome coronavirus. *Journal of virology* **81**, 12049–12060 (2007).
143. Krogh, A., Larsson, B., von Heijne, G. & Sonnhammer, E. L. Predicting transmembrane protein topology with a hidden markov model: application to complete genomes<sup>11</sup>Edited by F. Cohen. en. *Journal of Molecular Biology* **305**, 567–580. ISSN: 00222836 (Jan. 2001).
144. Mirdita, M., Steinegger, M. & Soding, J. MMseqs2 desktop and local web server app for fast, interactive sequence searches. *Bioinformatics* **35**, 2856–2858 (2019).
145. Williams, C. J. *et al.* MolProbity: More and better reference data for improved all-atom structure validation. en. *Protein Science* **27**, 293–315. ISSN: 0961-8368, 1469-896X (Jan. 2018).
146. Vonrhein, C. *et al.* Data processing and analysis with the *autoPROC* toolbox. *Acta Crystallographica Section D Biological Crystallography* **67**, 293–302. ISSN: 0907-4449 (Apr. 1, 2011).
147. Keegan, R. M. & Winn, M. D. *MrBUMP* : an automated pipeline for molecular replacement. *Acta Crystallographica Section D Biological Crystallography* **64**, 119–124. ISSN: 0907-4449 (Jan. 1, 2008).
148. Krissinel, E. *et al.* CCP 4 Cloud for structure determination and project management in macromolecular crystallography. *Acta Crystallographica Section D Structural Biology* **78**, 1079–1089. ISSN: 2059-7983 (Sept. 1, 2022).

149. Chou, C.-Y. *et al.* Structural basis for catalysis and ubiquitin recognition by the *Severe acute respiratory syndrome coronavirus* papain-like protease. *Acta Crystallographica Section D Biological Crystallography* **70**, 572–581. ISSN: 1399-0047 (Feb. 1, 2014).
150. Emsley, P. & Cowtan, K. Coot : model-building tools for molecular graphics. en. *Acta Crystallographica Section D Biological Crystallography* **60**, 2126–2132. ISSN: 0907-4449 (Dec. 1, 2004).
151. Bond, P. S. & Cowtan, K. D. *ModelCraft* : an advanced automated model-building pipeline using *Buccaneer*. *Acta Crystallographica Section D Structural Biology* **78**, 1090–1098. ISSN: 2059-7983 (Sept. 1, 2022).
152. Hajizadeh, N. R., Franke, D. & Svergun, D. I. Integrated beamline control and data acquisition for small-angle X-ray scattering at the P12 BioSAXS beamline at PETRAIII storage ring DESY. *Journal of Synchrotron Radiation* **25**, 906–914. ISSN: 1600-5775 (May 1, 2018).
153. Franke, D. & Svergun, D. I. *DAMMIF* , a program for rapid *ab-initio* shape determination in small-angle scattering. *Journal of Applied Crystallography* **42**, 342–346. ISSN: 0021-8898 (Apr. 1, 2009).
154. Semenyuk, A. V. & Svergun, D. I. GNOM – a program package for small-angle scattering data processing. *Journal of Applied Crystallography* **24**, 537–540. ISSN: 00218898 (Oct. 1, 1991).
155. Volkov, V. V. & Svergun, D. I. Uniqueness of *ab initio* shape determination in small-angle scattering. *Journal of Applied Crystallography* **36**, 860–864. ISSN: 0021-8898 (June 1, 2003).
156. Svergun, D. Restoring Low Resolution Structure of Biological Macromolecules from Solution Scattering Using Simulated Annealing. en. *Biophysical Journal* **76**, 2879–2886. ISSN: 00063495 (June 1999).
157. Holton, J. M., Classen, S., Frankel, K. A. & Tainer, J. A. The R-factor gap in macromolecular crystallography: an untapped potential for insights on accurate structures. *The FEBS journal* **281**, 4046–4060 (2014).
158. Sheldrick, G. M. Macromolecular phasing with SHELXE. *Zeitschrift für Kristallographie-Crystalline Materials* **217**, 644–650 (2002).

159. Bricogne, G., Vornrhein, C., Flensburg, C., Schiltz, M. & Paciorek, W. Generation, representation and flow of phase information in structure determination: recent developments in and around SHARP 2.0. *Acta Crystallographica Section D: Biological Crystallography* **59**, 2023–2030 (2003).
160. Rupp, B. *Biomolecular crystallography: principles, practice, and application to structural biology* (Garland Science, 2009).
161. Thorn, A. & Sheldrick, G. M. ANODE: anomalous and heavy-atom density calculation. *Journal of applied crystallography* **44**, 1285–1287 (2011).
162. Sheldrick, G. M. Experimental phasing with SHELXC/D/E: combining chain tracing with density modification. *Acta Crystallographica Section D: Biological Crystallography* **66**, 479–485 (2010).
163. Lamb, A. L., Kappock, T. J. & Silvaggi, N. R. You are lost without a map: Navigating the sea of protein structures. *Biochimica et Biophysica Acta (BBA)-Proteins and Proteomics* **1854**, 258–268 (2015).
164. Terwilliger, T. C. Reciprocal-space solvent flattening. *Acta Crystallographica Section D: Biological Crystallography* **55**, 1863–1871 (1999).
165. Cowtan, K. & Main, P. Phase combination and cross validation in iterated density-modification calculations. *Acta Crystallographica Section D: Biological Crystallography* **52**, 43–48 (1996).
166. Roberts, A. & Brünger, A. Phase improvement by cross-validated density modification. *Acta Crystallographica Section D: Biological Crystallography* **51**, 990–1002 (1995).
167. Sheldrick, G. *et al.* Shelx (2012).
168. Das, D. *et al.* Structure and function of the DUF2233 domain in bacteria and in the human mannose 6-phosphate uncovering enzyme. *Journal of Biological Chemistry* **288**, 16789–16799 (2013).
169. Vornrhein, C., Blanc, E., Roversi, P. & Bricogne, G. Automated structure solution with autoSHARP. *Macromolecular Crystallography Protocols: Volume 2: Structure Determination*, 215–230 (2007).
170. Perrakis, A., Harkiolaki, M., Wilson, K. S. & Lamzin, V. S. ARP/wARP and molecular replacement. *Acta Crystallographica Section D: Biological Crystallography* **57**, 1445–1450 (2001).

171. Adams, P. D. *et al.* PHENIX: a comprehensive Python-based system for macromolecular structure solution. *Acta Crystallographica Section D: Biological Crystallography* **66**, 213–221 (2010).
172. Wojdyr, M. GEMMI: A library for structural biology. *Journal of Open Source Software* **7**, 4200 (2022).
173. Stein, N. & Ballard, C. Intensity to amplitude conversion using CTRUNCATE. *Foundations of Crystallography* **65**, 161–161 (2009).
174. Potterton, L. *et al.* CCP4i2: the new graphical user interface to the CCP4 program suite. *Acta Crystallographica Section D: Structural Biology* **74**, 68–84 (2018).
175. Murshudov, G. N. *et al.* REFMAC5 for the refinement of macromolecular crystal structures. *Acta Crystallographica Section D: Biological Crystallography* **67**, 355–367 (2011).



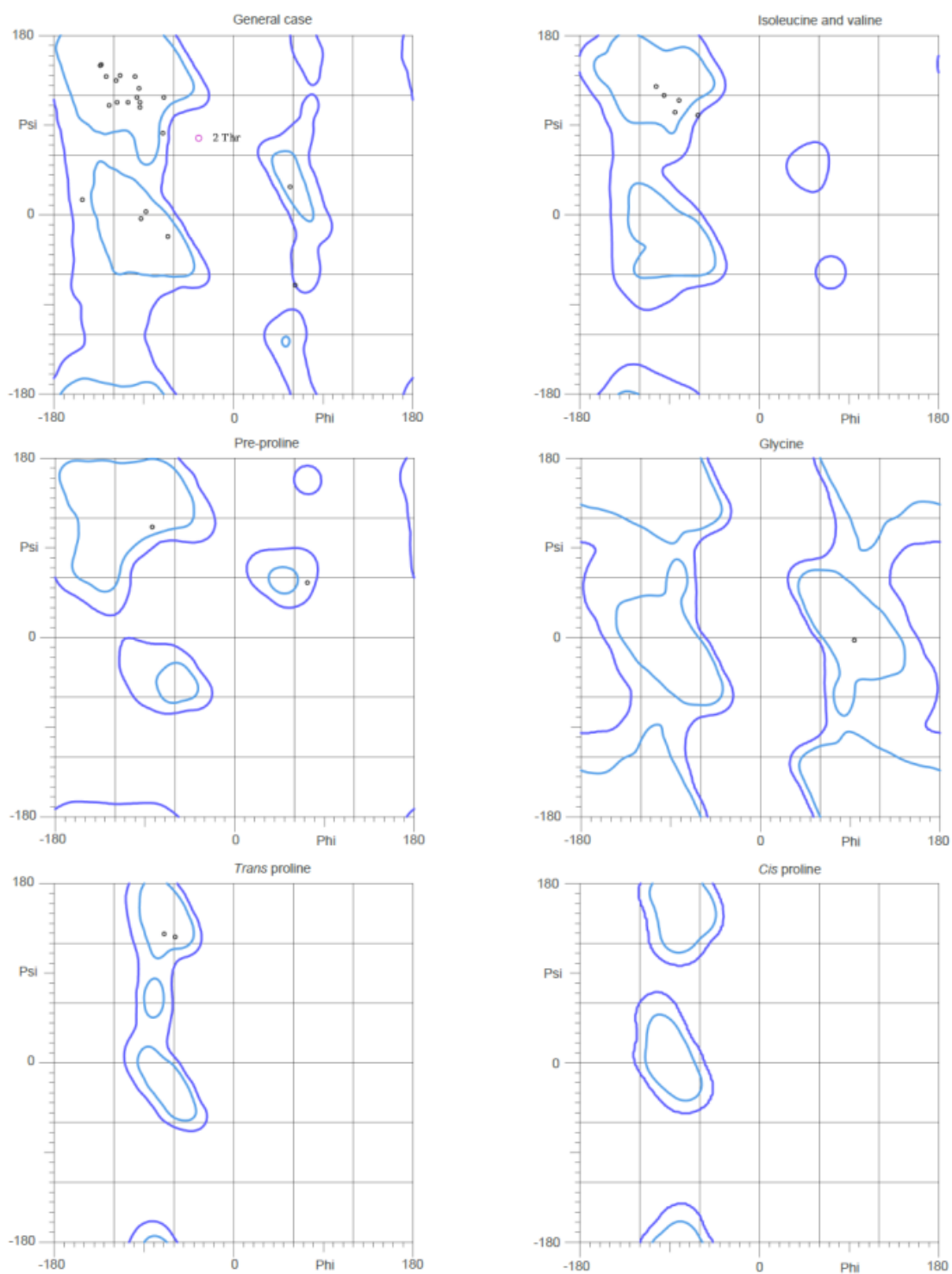
# A Appendix

**Table 17:** Preliminary residue ranges of *nsp3* domains from SARS-CoV-1 and MHV in comparison to ranges defined in previous literature. Previously defined ranges are from Lei et al. [17]. a: These ranges were not stated explicitly and are derived from the surrounding ranges. b: These ranges had no defined start/end. c: These ranges were predicted by TMHMM 2.0. d: This domain was not determined in the literature. Ranges with a orange shade are identical between the preliminary and the previous ranges. Ranges with a blue shade differ by a maximum of 5 residues into each direction.

Complete Name	Preliminary residue ranges for SARS-CoV-1	Previously defined residue ranges for SARS-CoV-1	Preliminary residue ranges for MHV	Previously defined residue ranges for MHV
Ubiquitin-like domain 1 (Ubl1)	1-112	1-112	1-115	20-109
Hypervariable region (HVR)	113-181	113-183	116-252	113-243
Papain-like protease 1 (PL1 <sup>pro</sup> )	-	-	253-480	252-500
Macrodomain 1 (Mac1)	182-358	184-365	481-655	509-633
Linker Mac1-Mac2/Linker Mac1-DPUlike	359-388	366-388 <sup>a</sup>	656-693	634-701 <sup>a</sup>
Macrodomain 2 (Mac2)	389-516	389-524	-	-
Linker Mac2-Mac3	517-526	-	-	-
Macrodomain 3 (Mac3)	527-651	525-652	-	-
Domain preceding Ubl2 and PL2 <sup>pro</sup> (DPUP)	652-722	653-720	694-777	702-776
Ubiquitin-like domain 2 (Ubl2)	723-781	723-??? <sup>b</sup>	778-837	777-??? <sup>b</sup>
Papain-like protease 2 (PL2 <sup>pro</sup> )	782-1040	???-1036 <sup>b</sup>	838-1080	???-1074 <sup>b</sup>
Linker PL2 <sup>pro</sup> -NAB	1041-1065	1037-1065 <sup>a</sup>	1081-1116	1075-1110 <sup>a</sup>
Nucleic-acidic-binding domain (NAB)	1066-1180	1066-1180	1117-1232	1111-1229
Betacoronavirus-specific marker domain ( $\beta$ SM)	1181-1390	1203-1318	1233-1448	1277-1401
Transmembrane domain 1 (TM1)	1391-1413 <sup>c</sup>	1391-1413 <sup>c</sup>	1449-1471	1450-1471
Nsp3 ectodomain	1414-1495 <sup>c</sup>	1414-1495 <sup>c</sup>	1498-1564	1472-1577 <sup>a</sup>
Transmembrane domain 2 (TM2)	1496-1518 <sup>c</sup>	1496-1518 <sup>c</sup>	1565-1587	1578-1600
Linker TM2-AH1	1519-1522 <sup>c</sup>	1519-1522 <sup>a, c</sup>	1588-1607	??? <sup>d</sup>
Amphipathic helix 1 (AH1)	1523-1545 <sup>c</sup>	1523-1545 <sup>c</sup>	1608-1630	??? <sup>d</sup>
Nidovirus-conserved domain of unknown function (Y1)	1546-??? <sup>b</sup>	1546-??? <sup>b</sup>	1631-???	??? <sup>d</sup>
Coronavirus-specific C-terminal domain (Y2)	???-1820	???-??? <sup>b</sup>	???-1904	??? <sup>d</sup>
Coronavirus-specific C-terminal domain (Y3)	1821-1922	???-1922 <sup>b</sup>	1905-2006	???-2004

**Table 18:** Sequence similarities and RMSD values between predicted folds of domains from SARS-CoV-1 and MHV. Only domains predicted to fold into a defined structure and large regions of disorder are listed. RMSD values are calculated with PyMOL [131] for folded domains. Results are sorted by decreasing sequence similarity. Domains of the transmembrane region are listed below and are not sorted, since only short alignments were found. For these cases, the alignment-length is given in parentheses. Sequence similarity was used over sequence identity due to also listing disordered domains, which show high similarity and low identities. Domains marked with an asterisk consist of subdomains, which are also listed individually.

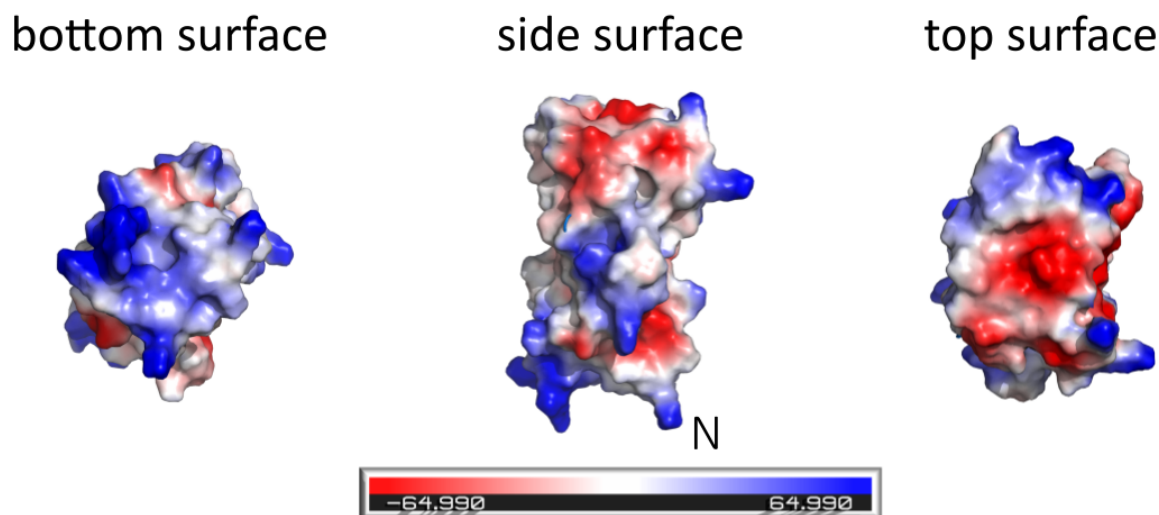
Domain	SARS-CoV-1 to MHV	
	Sequence similarity	RMSD
Y1a	73.0%	0.3 Å
Y1*	65.0%	0.7 Å
Y1b	60.0%	0.6 Å
CoV-Ya	60.0%	0.7 Å
NAB	42.6%	0.9 Å
DPUP	60.0%	3.2 Å
βSM-N	60.0%	-
CoV-Yb	58.0%	0.5 Å
CoV-Y*	57.3%	2.4 Å
Ubl1	55.0%	0.6 Å
βSLD	54.0%	0.4 Å
Ubl2	52.0%	0.6 Å
Mac1	51.7%	0.8 Å
PL2 <sup>pro</sup>	48.2%	1.0 Å
EctoC	48.0%	8.6 Å
βSM-C	46.0%	-
βSM-M	39.0%	2.9 Å
HVR	35.0%	-
TM1	40.0% (5)	0.2 Å
EctoL + EctoTM	69.0% (13)	-
EctoTM	-	-
TM2	75.0% (8)	0.3 Å
AH1	53.0% (17)	0.5 Å

**Ramachandran Analysis for betacoronavirus-specific linker domain ( $\beta$ SLD)**

**Figure 43:** Ramachandran plots for the structure prediction of the Betacoronavirus-specific linker domain. Plot was generated via MolProbity [145].

**Table 19:** Results from local pairwise sequence alignment between the sequences of several beta-coronaviruses and the sequence of SARS-CoV-2 linker domain. The alignment and calculations were performed via EMBOSS Water version 6.6.0 [74]. The alignment was performed either with the sequence from an *nsp3* equivalent or with the whole sequence of the *pp1ab* from the respective virus, if no *NSP3* annotation was present.

Virus	NCBI reference-id	Sequence identity	Sequence similarity
SARS-CoV-1	NC_004718	79 %	91 %
Murine hepatitis virus strain A59	NC_048217.1	43 %	54 %
Rousettus bat coronavirus HKU9	NC_009021.1	52 %	58 %
Rousettus bat coronavirus isolate GCCDC1 356	NC_030886.1	53 %	68 %
Bat Hp-betacoronavirus	NC_025217.1	55 %	69 %
Human coronavirus HKU1	NC_006577.2	37 %	48 %
Rabbit coronavirus HKU14	NC_017083.1	43 %	54 %
Human coronavirus OC43	NC_006213.1	43 %	46 %
Bovine coronavirus isolate BCoV-ENT	NC_003045.1	43 %	49 %
Betacoronavirus HKU24 strain HKU24-R05005I	NC_026011.1	43 %	50 %
Rat coronavirus Parker	NC_012936.1	46 %	54 %
Betacoronavirus Erinaceus	NC_039207.1	45 %	48 %
Tylonycteris bat coronavirus HKU4	NC_009019.1	55 %	66 %
Pipistrellus bat coronavirus HKU5	NC_009020.1	43 %	59 %
MERS-CoV	NC_019843.3	55 %	62 %
Betacoronavirus England 1 isolate H123990006	NC_038294.1	55 %	62 %



**Figure 44:** *Vacuum electrostatics calculated for AlphaFold2 [15] structure prediction of SARS-CoV-2 Y1 domain. Calculation and visualization was made in PyMOL [131]. Shown are the three surfaces, where the bottom surface is positively charged and the top surface negatively.*

## A.1 Appendix of Part I: Beyond the Surface

The initial search query on the PDB for structures associated with MAD data, which fulfilled certain criteria, provided 573 structures. These were filtered down to 105 structures which were handled automatically by SHELXC, ANODE, and SHELXE, listed here:

1VPM, 1VR5, 2FUP, 2GVK, 2HTD, 2HUH, 2HYT, 2ISB, 2IT9, 2NLV, 2NYH, 2OC5, 2OPL, 2OZH, 2OZJ, 2PQ7, 2PYQ, 2QEU, 2QIW, 2QL8, 2QPX, 2QTD, 2QZC, 2R01, 2R0X, 2RA9, 2RAF, 2RB7, 2RBD, 2RDC, 3B8L, 3BDI, 3BHN, 3BN7, 3C8L, 3CGH, 3CJM, 3CK1, 3CVO, 3D4E, 3D5P, 3DCZ, 3DI4, 3E0Z, 3E10, 3EBT, 3EBY, 3E08, 3EQX, 3ETN, 3F14, 3FDH, 3FFR, 3GE5, 3GYD, 3GZA, 3HN5, 3HN7, 3HZP, 3I09, 3IHV, 3IMK, 3IRB, 3K11, 3KE7, 3KGY, 3KS6, 3LHN, 3LHO, 3LLC, 3LLX, 3LWC, 3LYG, 3M1T, 3MCW, 3MZ2, 3N6Z, 3NO2, 3NOH, 3NPD, 3NUF, 3OHG, 3OSD, 3OYV, 3OZ2, 3POH, 3Q1N, 3QC0, 3RJV, 3SD2, 3SEE, 3SGG, 3UE2, 4DWF, 4FS7, 4H08, 4H17, 4HLB, 4IAB, 4ICI, 4IPB, 4IRT, 4JM1, 4KH8, 4LER

Additional filtering based on the average anomalous density left 74 structures:

1VPM, 1VR5, 2FUP, 2GVK, 2HUH, 2NLV, 2OC5, 2OZH, 2PQ7, 2QEU, 2QIW, 2QL8, 2QPX, 2QZC, 2R01, 2R0X, 2RA9, 2RB7, 2RBD, 3BDI, 3BHN, 3CGH, 3CJM, 3CK1, 3D4E, 3D5P, 3DCZ, 3DI4, 3EBT, 3EBY, 3EQX, 3F14, 3FDH, 3FFR, 3GE5, 3GZA, 3HN7, 3IHV, 3IMK, 3IRB, 3K11, 3KGY, 3LHN, 3LLC, 3LLX, 3LWC, 3LYG, 3M1T, 3MCW, 3N6Z, 3NO2, 3NOH, 3NPD, 3NUF, 3OHG, 3OSD, 3OYV, 3OZ2, 3POH, 3Q1N, 3QC0, 3RJV, 3SD2, 3SEE, 3SGG, 3UE2, 4FS7, 4H08, 4HLB, 4IAB, 4ICI, 4JM1, 4KH8, 4LER



## A Acknowledgements

First and foremost I thank my supervisor Andrea Thorn, for making all of this possible in the first place. Thank you very much for initiating the project and giving me the chance to complete my PhD projects here in Hamburg! A huge "Thank you!" goes also to Arwen Pearson for creating such a wonderful scientific community at HARBOR, which made working at such a place very enjoyable! I also thank you for the support in supervision, with another great thanks going to Andrew Torda helping in this regard.

Special thanks goes to David Briggs, Yunyun Gao, Armin Wagner, Mike Hough, and to James Holton, who collaborated in the projects in the one or the other way. Without you, many results wouldn't have been possible! I would like to thank Gianluca Santona, Clemens Vornrhein, Dale Tronrud and all the members of Coronavirus Structural Taskforce for the great, helpful, and interesting discussions! These pushed the projects into new directions and it was a pleasure to learn from all of you.

Great thanks goes also to Torben Steenbock, Michael Rübhausen, and Gabriel Bester, who made it possible to teach python. I had great time during the course, especially when teaching with Torben, thanks for that!

Of course I have to thank my great colleagues at HARBOR, with a special thanks to Sabrina, Lea, Erik, Philip, Kristopher, Jan, Yunyun and Pairoh! It was fun to share the working space with you and I got a lot of support from you during some difficult times, thank you! Great thanks goes also to Irene's group for all the nice invitations to their group events and of course for organizing the H-Bar. All of the other groups at HARBOR receive also a great thank you, without you the HARBOR would just be a boring office, but you made it fun and living place, had great time with all of you!

Last but not least, I have to leave a huge THANK YOU! to my parents, my brothers Alex and Nik, and all of my room mates who always supported me through anything of this wonderful stage of my life. Special thanks goes out to Arne, Amjad, Jan, Silva, Janek, and Julia (also for infinite supply of Mate). You made the time in Hamburg great, but I also have to thank all of my close friends from Bielefeld. Despite the distance, you were always here :)

The dissertation was written with L<sup>A</sup>T<sub>E</sub>X via the Overleaf editor. Images were edited or created via Inkscape. DeepL Write was used to suggest synonyms and to correct English grammar mistakes.



## Eigenständigkeitserklärung

Hiermit erkläre ich an Eides statt, dass ich die vorliegende Dissertationsschrift selbst verfasst und keine anderen als die angegebenen Quellen und Hilfsmittel benutzt habe.

---

Ort, Datum

---

Maximilian Edich

## **Auflistung verwendeter Gefahrenstoffe nach GHS**

Im Rahmen dieser Arbeit wurden keine Gefahrenstoffe nach GHS verwendet.