

UC Davis

UC Davis Previously Published Works

Title

Computational models in the service of X-ray and cryo-EM structure determination

Permalink

<https://escholarship.org/uc/item/9rp0b02z>

Authors

Kryshtafovych, Andriy

Moult, John

Albrecht, Reinhard

et al.

Publication Date

2021-07-22

DOI

10.22541/au.162696116.65229185/v1

Copyright Information

This work is made available under the terms of a Creative Commons Attribution License, available at <https://creativecommons.org/licenses/by/4.0/>



Computational models in the service of X-ray and cryo-EM structure determination

Journal:	<i>PROTEINS: Structure, Function, and Bioinformatics</i>
Manuscript ID	Draft
Wiley - Manuscript type:	Research Article
Date Submitted by the Author:	n/a
Complete List of Authors:	<p>Kryshtafovych, Andriy; University of California Davis, Genome Center Moult, John; University of Maryland Biotechnology Institute, Institute for Bioscience and Biotechnology Research and Department of Cell Biology and Molecular Genetics Albrecht, Reinhard; Max Planck Institute for Developmental Biology, Department of Protein Evolution Chang, Geoffrey; University of California San Diego, Department of Pharmacology Chao, Kinlin; University of Maryland Biotechnology Institute, Institute for Bioscience and Biotechnology Research Fraser, Alec; University of Texas Medical Branch at Galveston, Department of Biochemistry and Molecular Biology, Sealy Center for Structural Biology and Molecular Biophysics (SCSB) Greenfield, Julia; University of Maryland Biotechnology Institute, Institute for Bioscience and Biotechnology Research Hartmann, Marcus; Max Planck Institute for Developmental Biology, Department of Protein Evolution Herzberg, Oznat; University of Maryland, Institute of Bioscience and Biotechnology Research; University of Maryland, Chemistry & Biochemistry Josts, Inokentij; University of Hamburg, The Hamburg Advanced Research Center for Bioorganic Chemistry (HARBOR) and Department of Chemistry, Institute for Biochemistry and Molecular Biology Leiman, Petr; University of Texas Medical Branch at Galveston, Department of Biochemistry and Molecular Biology, Linden, Sara; University of Maryland Biotechnology Institute, Institute for Bioscience and Biotechnology Research Lupas, Andrei N.; Max-Planck-Inst. for Developmental Biology, Dept. of Protein Evolution Nelson, Daniel; University of Maryland, Institute for Bioscience and Biotechnology Research and Department of Veterinary Medicine; University of Maryland Biotechnology Institute, Institute for Bioscience and Biotechnology Research Rees, Steven; University of California San Diego, School of Pharmacy and Pharmaceutical Science Shang, Xiaoran; University of Maryland Biotechnology Institute, Institute for Bioscience and Biotechnology Research Sokolova, Maria; Skolkovo Institute of Science and Technology, Center of Life Sciences Tidow, Henning; University of Hamburg, The Hamburg Advanced</p>

	Research Center for Bioorganic Chemistry (HARBOR) and Department of Chemistry, Institute for Biochemistry and Molecular Biology
Key Words:	X-ray crystallography; cryo-EM; CASP, protein structure prediction

SCHOLARONE™
Manuscripts

1
2
3 **Computational models in the service of X-ray and cryo-EM structure**
4
5
6 **determination**
7
8

9 Running title: Models in service
10
11
12
13

14 **Andriy Kryshtafovych**, ¹ - Genome Center, University of California, Davis, California
15
16 95616, USA
17
18

19 **John Moutl**, ² - Institute for Bioscience and Biotechnology Research, Department of Cell
20
21 Biology and Molecular genetics, University of Maryland, 9600 Gudelsky Drive, Rockville,
22
23 MD 20850, USA
24
25

26 **Reinhard Albrecht**, ³ - Department of Protein Evolution, Max Planck Institute for
27
28 Developmental Biology, 72076 Tübingen, Germany
29
30

31 **Geoffrey A. Chang**, ⁴ - Department of Pharmacology, University of California-San Diego, La
32
33 Jolla, CA, 92093, USA
34
35

36 **Kinlin Chao** ⁵, Institute for Bioscience and Biotechnology Research, University of Maryland,
37
38 Rockville, MD 20850, USA
39
40

41 **Alec Fraser**, ⁶ - Department of Biochemistry and Molecular Biology, Sealy Center for
42
43 Structural Biology and Molecular Biophysics (SCSB), The University of Texas Medical Branch
44
45 at Galveston, TX 77555, USA
46
47

48 **Julia Greenfield** ⁵, Institute for Bioscience and Biotechnology Research, University of
49
50 Maryland, Rockville, MD 20850, USA
51
52

53 **Marcus D. Hartmann**, ³ - Department of Protein Evolution, Max Planck Institute for
54
55 Developmental Biology, 72076 Tübingen, Germany
56
57
58
59
60

1
2
3 **Osnat Herzberg**, ^{5,7} – ⁵ Institute for Bioscience and Biotechnology Research, University of
4 Maryland, Rockville, MD 20850, USA; ⁷ Department of Chemistry and Biochemistry,
5 University of Maryland, College Park, MD 20742, USA
6
7
8

9
10 **Inokentij Josts**, ⁸ - The Hamburg Advanced Research Center for Bioorganic Chemistry
11 (HARBOR) & Department of Chemistry, Institute for Biochemistry and Molecular Biology,
12 University of Hamburg, Luruper Chaussee 149, 22761 Hamburg, Germany
13
14
15

16
17
18 **Petr G. Leiman**, ⁶ - Department of Biochemistry and Molecular Biology, Sealy Center for
19 Structural Biology and Molecular Biophysics (SCSB), The University of Texas Medical Branch
20 at Galveston, TX 77555, USA
21
22
23

24
25 **Sara B. Linden** ⁵ - Institute for Bioscience and Biotechnology Research, University of
26 Maryland, Rockville, MD 20850, USA
27
28

29
30 **Andrei N. Lupas**, ³ - Department of Protein Evolution, Max Planck Institute for Developmental
31 Biology, 72076 Tübingen, Germany
32
33
34

35
36 **Daniel C. Nelson** ^{5,9} - ⁵ Institute for Bioscience and Biotechnology Research, University of
37 Maryland, Rockville, MD 20850, USA; ⁹ Department of Veterinary Medicine, University of
38 Maryland, College Park, MD 20742, USA
39
40

41
42 **Steven D. Rees**, ⁴ - Department of Pharmacology, University of California-San Diego, La
43 Jolla, CA, 92093, USA
44
45

46
47 **Xiaoran Shang** ⁵ - Institute for Bioscience and Biotechnology Research, University of
48 Maryland, Rockville, MD 20850, USA
49
50

51
52 **Maria L. Sokolova**, ¹⁰ - Center of Life Sciences, Skolkovo Institute of Science and Technology,
53 Moscow, 121205, Russia
54
55
56
57
58
59
60

1
2
3 **Henning Tidow**,⁸ - The Hamburg Advanced Research Center for Bioorganic Chemistry
4
5 (HARBOR) & Department of Chemistry, Institute for Biochemistry and Molecular Biology,
6
7 University of Hamburg, Luruper Chaussee 149, 22761 Hamburg, Germany
8
9

10 and

11
12
13 **AlphaFold2 team**,¹¹ - DeepMind, London, EC4A 3TW, UK
14
15
16
17
18

19 **Keywords:** X-ray crystallography; cryo-EM; CASP, Protein Structure Prediction.
20
21

22 **Abbreviations:** **CASP:** community wide experiment on the Critical Assessment of
23
24 Techniques for Protein Structure Prediction; **MR:** molecular replacement; **cryo-EM:** cryo-
25
26 electron microscopy; **NMR:** nuclear magnetic resonance
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Abstract

CASP (Critical Assessment of Structure prediction) conducts community experiments to determine the state of the art in computing protein structure from amino acid sequence. The process relies on the experimental community providing information about not yet public or about to be solved structures, for use as targets. For some targets, the experimental structure is not solved in time for use in CASP. Calculated structure accuracy improved dramatically in this round, implying that models should now be much more useful for resolving many sorts of experimental difficulty. To test this, selected models for seven unsolved targets were provided to the experimental groups. These models were from the AlphaFold2 group, who overall submitted the most accurate predictions in CASP14. Four targets were solved with the aid of the models, and, additionally, the structure of an already solved target was improved. An a-posteriori analysis showed that in some cases models from other groups would also be effective. This paper provides accounts of the successful application of models to structure determination, including molecular replacement for X-ray crystallography, backbone tracing and sequence positioning in a Cryo-EM structure, and correction of local features. The results suggest that in future there will be greatly increased synergy between computational and experimental approaches to structure determination.

Introduction

Besides being a challenging and interesting problem in itself, computational modeling of protein structure has significant practical impact on the biomedical field¹⁻³. The most direct application is in structural biology where models are used to help determine protein structures by experimental methods including X-ray crystallography, cryo-electron microscopy and NMR spectroscopy. In X-ray crystallography, models are often used to solve the phase problem by molecular replacement (MR), which relies on the existence of similar protein structures or accurate models that serve as templates to be placed in the crystal cell, consistent with the diffraction data⁴. In NMR, models can assist with the prediction of chemical shifts and NMR spectra, or the interpretation of real spectra (i.e., chemical shift assignments and then NOE assignments) and in building structures that satisfy experimentally derived distance and angle restraints^{5,6}. In cryo-EM, models are of value for backbone tracing and fitting sequence into a map, especially at low and moderate resolution (3.5-5.0 Å)^{7,8}. Regardless of the structure determination technique, models can be used to identify and sometimes fix problematic regions in experimental structures⁹. With the recent major advances in protein structure modeling¹⁰⁻¹³, it is clear that in future models will play a substantially larger role in determining and validating experimental structures.

In CASP, not-yet solved or not-yet released structures are solicited from the experimental community as modeling targets. The suitability of a structure as a target is largely determined by three factors: estimated modeling difficulty (some may be too easy), whether there is sufficient time available before experimental structure release, and conversely, whether the experimental structure will be solved in time for model assessment. Inevitably, some targets will encounter problems, and normally have to be abandoned. There were eleven such targets in CASP14, including seven where experimental data have been collected, but, nevertheless, the structure could not be determined. Because of the very high accuracy of many submitted

1
2
3 models on other targets, especially those from the AlphaFold2 group ^{11,12}, the organizers
4 decided to see how many of the challenging structures could be resolved with the aid of models.
5
6 In previous CASPs, generated models have occasionally helped solve structures. For example,
7
8 the crystal structure of Sla2 ANTH domain of *Chaetomium thermophilum* (CASP11 target
9
10 T0839) was determined by molecular replacement using CASP models¹⁴, but these have been
11
12 exceptions. In CASP14, four structures were solved with the aid of AlphaFold2 models. A post-
13
14 CASP analysis has shown¹⁵ that models from other groups would also have been effective in
15
16 some cases. In the three remaining unsolved cases, poor data quality appears to have been the
17
18 issue. These are all ‘hard’ targets with limited or no homology information available for at least
19
20 some domains, demonstrating the power of the new methods for all classes of modeling
21
22 difficulty. For one other target, provision of the models resulted in correction of a local
23
24 experimental error. The paper discusses these success stories, with content for each target
25
26 provided by the corresponding experimental group.
27
28
29
30
31
32
33

34 **Results**

35 **1. AlphaFold2 models help solve crystal structure of the inner membrane reductase**

36 **FoxB (CASP: T1058) by molecular replacement – by IJ and HT.**

37
38 From email to the CASP Prediction Center: *The model you sent me (from the leading group)*
39
40 *worked for molecular replacement and we finally solved the structure by MR-SAD. I am still*
41
42 *astonished that the human expert model worked, while none of the server models we tried did*
43
44 *(as they were rather similar). Henning Tidow*
45
46
47
48
49
50

51 **1.1. Brief description of the target**

52
53 Most microorganisms rely on the bioavailability of iron for their survival. Due to the
54
55 low solubility of ferric iron, they often use secreted siderophores for the chelation and uptake
56
57 of iron. In Gram-negative bacteria, siderophores are usually taken up by TonB-dependent
58
59 transporters (TBDTs) located in the bacterial outer membrane. The route of ferric-siderophores
60

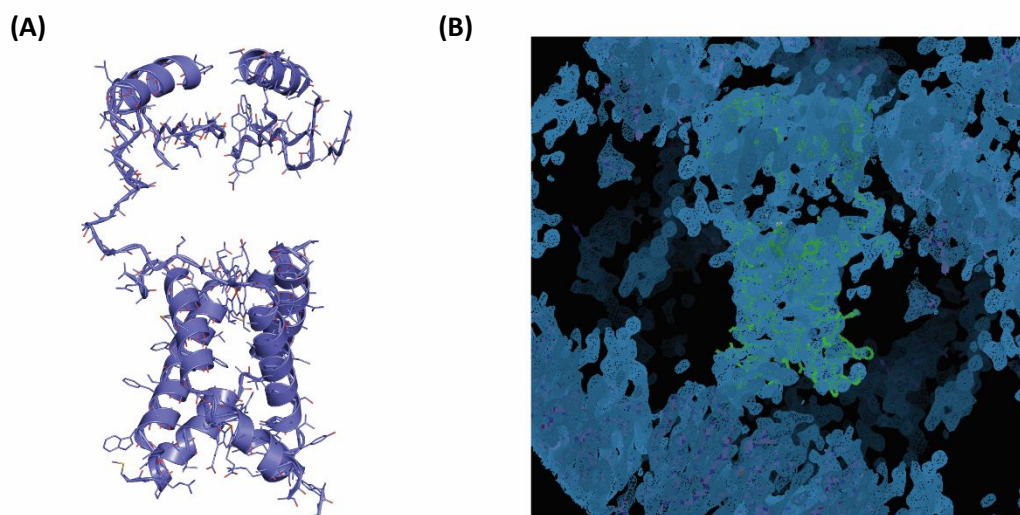
1
2
3 across the inner membrane (IM) is less straightforward and differs across many bacterial species
4 and siderophore chemistries. Ferric-siderophore complexes are either recognized by the
5
6 dedicated periplasmic-binding proteins for delivery to IM transporters for uptake into the
7
8 cytoplasm or the iron is released from the ferric-siderophore complexes by a reduction
9
10 mechanism. The Gram-negative bacterium *Pseudomonas aeruginosa* (an opportunistic human
11
12 pathogen) is able to take up Fe-siderophore complexes called ferrioxamines via a dedicated
13
14 TBDT FoxA in an act of siderophore piracy¹⁶. For several years we also worked towards the
15
16 structure determination of FoxB, another protein of unknown function located in the same
17
18 operon as FoxA. With the help of the AlphaFold2 model generated in the course of the CASP14
19
20 competition, we were able to determine the structure of FoxB. It possesses a novel fold with
21
22 the transmembrane domain harboring two heme molecules indicating a role as inner membrane
23
24 reductase involved in Fe-siderophore uptake and processing¹⁷.
25
26
27
28
29
30

31 **1.2. Workflow of how an AlphaFold2 model helped to solve the structure**

32
33
34 Native FoxB crystals obtained in decyl-maltopyranoside (DM) diffracted to approximately
35
36 5 Å resolution on average. Most of the crystals belonged to the $P2_12_12_1$ space group. All crystals
37
38 were obtained in 30% PEG 600, 0.1 M BICINE pH 9, 0.1 M ZnSO₄. Use of a lipid-like peptide
39
40 (LLP7) as additive allowed us to collect several datasets extending to 3.4-3.5 Å¹⁸.
41
42
43
44

45 All molecular replacement attempts using distant homologs and homology models thereof
46
47 failed. We acquired Se-Met anomalous data to 3.5 Å resolution, with anomalous signal to 4.5
48
49 Å as well as anomalous data at the Fe edge with anomalous signal to 4 Å, as we knew from
50
51 spectroscopic characterization that FoxB most likely contained at least one heme group.
52
53 Combining all anomalous data provided some experimental phases and allowed partial model
54
55 building. A single FoxB was present in the asymmetric unit. However, phasing power was only
56
57 sufficient to build approximately 60-70% of the backbone structure (Fig. 1). Although two
58
59 heme groups could be successfully placed, further tracing of the protein backbone and confident
60

1
2
3 sequence assignment was prevented by the low number of Met residues (5/382) and low
4
5 resolution of the datasets. Lysozyme fusion at the N-terminus also resulted in crystals
6
7 diffracting to approximately 4.5 Å.
8
9



33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Figure 1

Figure 1. (A) Partial FoxB model obtained by experimental phasing before the CASP14 model became available. At this point the model could not be further improved and the project was stuck for a year. (B) Experimental phases with partial FoxB model (map shown at 1.2 σ level).

At that point we were stuck with experimental phasing and submitted the FoxB sequence to CASP14. The model provided by AlphaFold2 (T1058TS427_3) resulted in a clear molecular replacement (MR) solution (TFZ: 18.9 / LLG: 324). We have also succeeded in finding a suitable crystal which diffracted to better than 3.5 Å, the crystal belonged to the P2₁2₁2 space group and contained 2 molecules of FoxB in the asymmetric unit. The final resolution of the diffraction dataset after starANISO processing was 3.1 Å. Subsequent MR-SAD and several rounds of building/refinement using COOT¹⁹ and REFMAC²⁰ further improved the model and

1
2
3 resulted in a good electron density map for the entire protein. Anomalous difference maps were
4
5 also used to validate the model (Fig. 2).
6
7

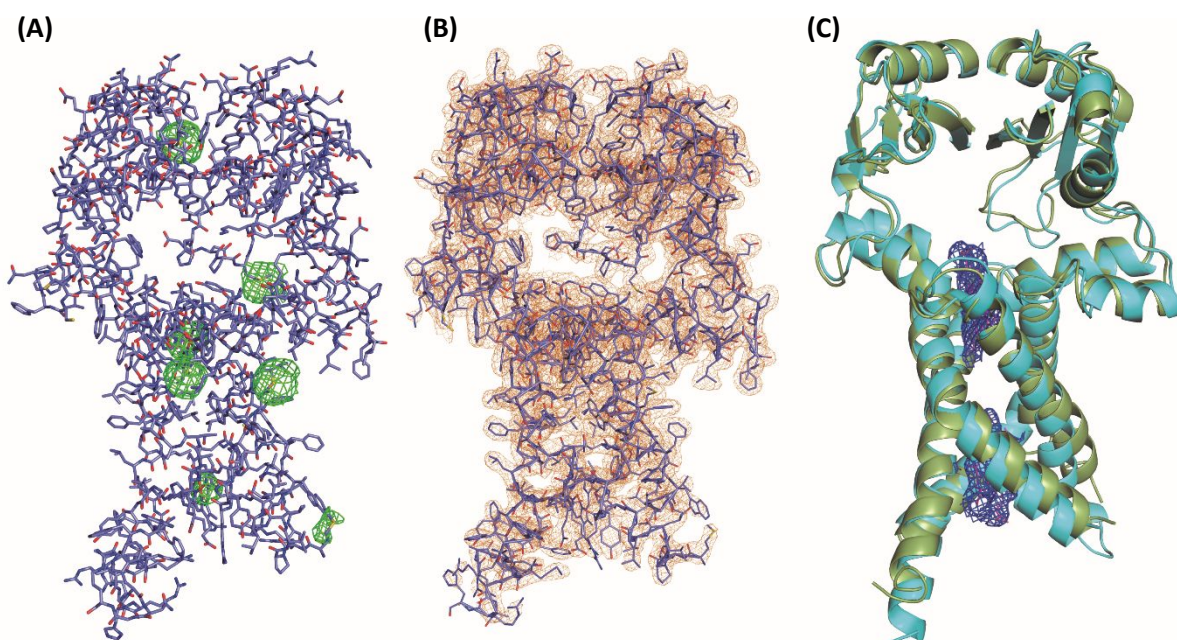


Figure 2

31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Figure 2. Workflow of FoxB structure determination. The structure was determined by MR-SAD using the AlphaFold2 model and experimental phases. (A) Anomalous difference map with Se and Fe sites at 2σ . (B) Overall map of FoxB after refinement (2σ). (C) Superposition of the final model (green) and AlphaFold2 model (cyan) shows excellent agreement. Density for heme groups (not present in AlphaFold2 model) is shown.

1.3. Model accuracy

The AlphaFold2 model that was used for the study (T1058TS427_3) shows a remarkable similarity to the final structure¹⁷. The overall RMSD is 1.17 Å for all atoms and 0.973 Å for C α atoms. Not only were all transmembrane helices built and registered correctly, but also the periplasmic domains containing several loops were modelled with high accuracy. There was no density for the cytoplasmic loop connecting TM helices 2 and 3 (residues 172-188), and it was

1
2
3 therefore omitted from the final model. Molecular replacement was only successful with the
4
5 AlphaFold2 model but not with server models from the CASP14 experiment (>30 models tried,
6
7 many of them with correct overall fold).
8
9

10
11 The success of the AlphaFold2 models seems to be due to their models “getting the details
12
13 right”, which was required for a clear MR solution. As one example for the accuracy of the
14
15 AlphaFold2 model, the His residues coordinating the two heme groups in FoxB were positioned
16
17 correctly, although this model did not contain heme groups (as we only provided the protein
18
19 sequence to CASP14). This fact however, also highlights a current limitation of the AlphaFold2
20
21 model: While it provides an astonishing good model for the apo protein, it is obviously still
22
23 lacking the functional groups (two heme groups in case of FoxB), which are responsible for the
24
25 biological function.
26
27
28
29
30
31

32 **2. The astounding accuracy of AlphaFold2 models of all subunits of phage AR9 non-** 33 34 **virion RNA polymerase (CASP: T1092-T1096) – by AF, MLS and PGL.** 35 36

37
38 From email to the CASP Prediction Center: *We are shocked... stunned... by the quality of the*
39
40 *model. You would not believe how much effort we have put into getting this structure. Years of*
41
42 *work... Both cryo-EM and crystallography... I mean, this is really shocking. Petr Leiman*
43

44 **2.1. Brief description of the target**

45
46 A group of large or “jumbo” bacteriophages, with genomes larger than 200 kbp, encode two
47
48 distinct DNA-dependent RNA polymerases (RNAPs), allowing these phages to assemble
49
50 independently from the host RNAP²¹⁻²⁴. One of these phage-encoded RNAPs is packaged into
51
52 the phage capsid and hence is called the virion RNAP (vRNAP). Following the attachment to
53
54 the host cell, the virus injects the vRNAP together with its DNA into the host cytoplasm. After
55
56 injection, the vRNAP transcribes early phage genes, including those of the second RNAP (the
57
58 non-virion RNAP, nvRNAP). The latter transcribes late genes, including those that encode for
59
60

1
2
3 the ν RNAP, which is then packaged into newly assembled phage particles. The exact
4
5 mechanism of this temporal and spatial activation/regulation of transcription is unclear but it is
6
7 known that ν - and ν RNAPs recognize different promoters ²³.
8
9

10
11 Both ν - and ν RNAPs are distantly related to multi-subunit RNAPs (msRNAPs) of bacteria,
12
13 eukaryotes, and archaea ²³. The universally conserved core of cellular msRNAPs contains six
14
15 subunits $\alpha_2\beta\beta'\omega$, and the catalytic cavity is formed by β and β' ²⁵. However, neither ν - or
16
17 ν RNAPs contain homologs of α or ω subunits, and their β and β' subunits are split into two or
18
19 three separate genes that are located in different regions of the phage genome. For sequence-
20
21 specific initiation of transcription, the phage AR9 ν RNAP core is required to form a complex
22
23 with a promoter specificity subunit gene product 226 (gp226) that shows no sequence similarity
24
25 to any known bacterial, eukaryotic, or archaeal transcription initiation factor. In fact, the amino
26
27 acid sequence of gp226 was a singleton in the GenBank database at the time of CASP14
28
29 experiment.
30
31
32
33

34
35 Besides employing a unique transcription factor, the AR9 ν RNAP possesses a number of
36
37 other distinct properties. Unlike any known msRNAP, the AR9 ν RNAP recognizes the
38
39 promoter in the template strand of double stranded DNA and can initiate promoter-specific
40
41 transcription on single stranded DNA ²⁶. Furthermore, as the genomic DNA of bacteriophage
42
43 AR9 contains deoxyuridine instead of thymidine ²¹, the AR9 ν RNAP is critically sensitive to
44
45 the presence of uracils in two key positions of its promoter sequence, and promoters with
46
47 thymines in these positions are not recognized ²⁶. To understand the novel and unusual
48
49 mechanism of promoter recognition by the AR9 ν RNAP, we decided to determine the
50
51 structure of this enzyme in various states: in complex with the specificity subunit and without
52
53 it, and in DNA template-bound and DNA-free forms. For the template, we used a short DNA
54
55 oligonucleotide that contained a promoter recognized by the AR9 ν RNAP *in vivo* and *in vitro*.
56
57
58
59
60

2.2. How AlphaFold2 models helped solve the structure

1
2
3 The most feature-full and continuous electron density map of the AR9 nvRNAP was
4 initially obtained by cryo-electron microscopy (cryo-EM) imaging of the nvRNAP holoenzyme
5 (i.e. containing the specificity subunit) in complex with the promoter-containing DNA
6 oligonucleotide. This complex contained five polypeptide chains – the specificity subunit
7 gp226, the N- and C-terminal parts of the β subunit gp105 and gp089 (respectively), and the N-
8 and C-terminal parts of the β' subunit gp270 and gp154 (respectively) – and the DNA
9 oligonucleotide, the structure of which will be described elsewhere. The cryo-EM
10 reconstruction was calculated using cryoSPARC²⁷ and had a resolution of 3.8 Å.
11
12
13
14
15
16
17
18
19
20
21

22 In parallel, several maps of the AR9 nvRNAP β - β' core (i.e. without the specificity subunit)
23 of varying quality and resolutions were obtained using X-ray crystallography. The dataset that
24 produced the best electron density also extended to 3.8 Å resolution, albeit this map was
25 significantly worse (poorer connectivity and quality of side chain features) than the cryo-EM
26 map. The phases for this dataset were obtained by eight-fold non-crystallographic
27 averaging^{28,29} of molecular replacement phases³⁰ calculated with the help of a partial model.
28 The latter was built using a single wavelength anomalous dispersion map of a dataset with a
29 smaller unit cell³¹⁻³³.
30
31
32
33
34
35
36
37
38
39
40

41 According to HHpred analysis at the time³⁴, the most similar RNAP with a known atomic
42 structure was that of *Mycobacterium tuberculosis* (PDB code 5ZX3³⁵). The AR9 nvRNAP
43 gp089, gp270, and gp154 proteins could all be aligned – with a 20-24% sequence identity and
44 100% probabilities – to continuous stretches of the *M. tuberculosis* RNAP β and β' subunits.
45 Gp105 was a more difficult target, with only its C-terminal half being predicted to be similar to
46 a fragment of the *M. tuberculosis* RNAP β subunit with an 80% probability and an E value of
47 2.3. The structure of gp226, as it was a unique sequence in the entire GenBank, could not be
48 reliably predicted by any tool.
49
50
51
52
53
54
55
56
57
58
59
60

Using both the best cryo-EM and X-ray maps of the AR9 nvRNAP and the structure of the *M. tuberculosis* RNAP as a chain-tracing guide in stretches of high sequence similarity, we manually built ~90% of the AR9 nvRNAP structure¹⁹. Some peripheral domains of gp105, gp154, and gp226 and regions for which no homology models existed were particularly challenging. Fortunately, while we were working on improving the cryo-EM map and X-ray phases to make the structure building process for these regions possible, the models of all five proteins produced by the AlphaFold2 team were made available to us by the CASP14 organizers. To our amazement, the AlphaFold2 models were of excellent quality and fit the cryo-EM and X-ray maps near perfectly almost everywhere including the no-homology regions (Fig. 3). This made the completion of the structure building process nearly trivial.

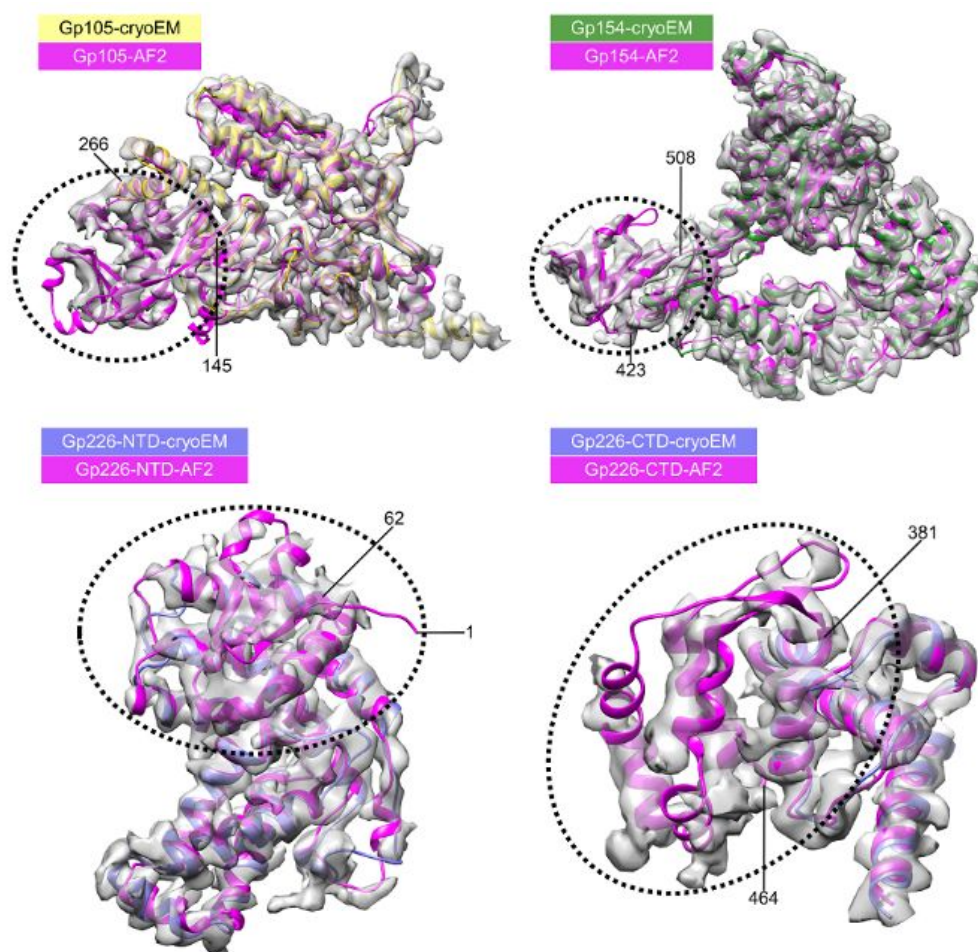


Figure 3

1
2
3 **Figure 3.** AlphaFold2 models of AR9 nvRNAP proteins fit the cryo-EM density nearly
4 perfectly. The cryo-EM-derived structures of gp105, gp154, and two gp226 domains are
5 colored according to the color code given in the upper left corner of each panel. All AlphaFold2
6 models are colored magenta. The electron density is contoured at 4.25 standard deviations
7 above the mean and colored semi-transparent grey. Regions where no cryo-EM-derived
8 structure existed prior to the availability of the AlphaFold2 models are indicated with a dashed
9 line and their boundary residues are labeled.
10
11
12
13
14
15
16
17
18
19
20
21
22

23. The accuracy of AlphaFold2 models

23
24
25
26 The AlphaFold2 models of individual domains were extremely similar to the cryo-EM-
27 derived structures. The only notable disagreement of AlphaFold2 models with experimental
28 data was in several regions of subunit contacts some of which are shown in Fig. 4. The
29 superposition of cryo-EM-derived structures and AlphaFold2 models resulted in the following
30 root mean square deviations between the equivalent C α atoms: 3.08 Å (465 out of 484 atoms)
31 for gp105, 2.00 Å (628 out of 649 atoms) for gp089, 2.50 Å (417 out of 426 atoms) for gp270,
32 2.42 Å (629 out of 631 atoms) for gp154, 1.54 Å (256 out of 261 atoms) for the N-terminal
33 domain of gp226 (gp226 NTD), and 2.76 Å (169 out of 169 atoms) for the C-terminal domain
34 of gp226 (gp226 CTD). Notably, the AlphaFold2 models were excellent at predicting the
35 structure of dynamic peripheral domains (e.g. the NTD and CTD of gp226). Additionally, the
36 AlphaFold2 helped to identify several *cis* prolines, which significantly improved the geometry
37 of the surrounding regions.
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

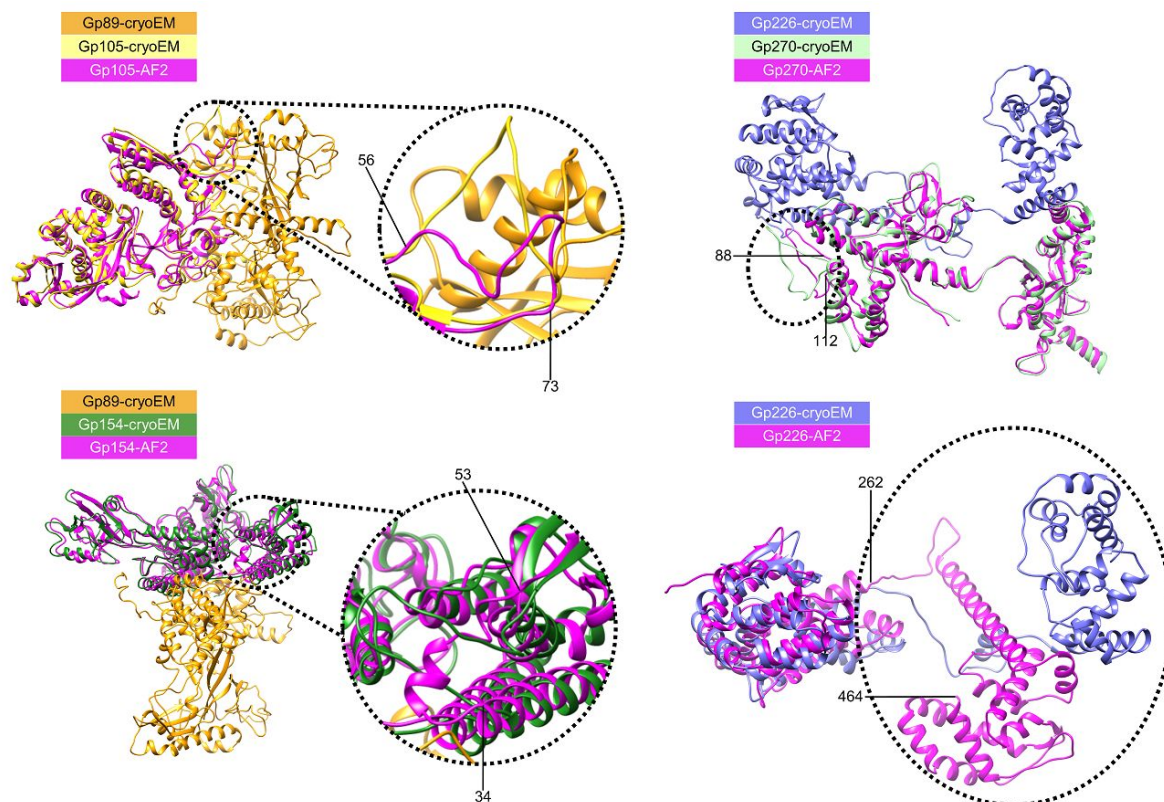


Figure 4

Figure 4. Inaccuracies in AlphaFold2 models. Cryo-EM-derived structures and AlphaFold2 models of several AR9 nvRNAP subunits are superimposed and regions where the conformation of the AlphaFold2 model deviates significantly from the cryo-EM-derived structure are indicated with a dashed line and their boundary residues are labeled. Note that the folds of both the N- and C-terminal domains of gp226 were predicted correctly, but the structure of the interdomain linker and the relative orientation of the two domains were incorrect.

The overall accuracy of AlphaFold2 models on multidomain targets was lower than that on individual domains, albeit still remarkably good (Fig. 4), and the structures of the four multidomain proteins that comprise the β - β' core of the AR9 nvRNAP were predicted correctly. The model of the gp226 interdomain linker and, as a consequence, the complete model of gp226

1
2
3 was incorrect, although this is hardly surprising considering the fact that the interdomain linker
4
5 does not have a well-defined secondary structure.
6
7

8
9 Besides collecting cryo-EM data on the AR9 nvRNAP holoenzyme in complex with the
10 promoter-containing DNA oligonucleotide, we crystallized it separately and collected X-ray
11 diffraction data to 3.4 Å resolution. This dataset had a solvent content of ~64% and contained
12
13 one molecule of the complete holoenzyme-DNA complex in the asymmetric unit of the *C*2
14 space group. As a final test of the accuracy of the AlphaFold2 models, we examined whether
15 they could serve as search models for solving the phase problem of this dataset by molecular
16 replacement. The models of gp105, gp089, gp270, and gp154 were used as is, without any
17 modification. The gp226 model consisted of two spatially separated globular domains (NTD)
18 and (CTD) connected by a long linker, so we treated the two domains as independent entities.
19 We then used Phaser³⁰ to perform an automatic molecular replacement procedure with these
20 six sets of coordinates as search models. The four proteins comprising the β-β' core of the
21 enzyme (gp105, gp089, gp270, and gp154) were placed correctly while the placement of both
22 gp226 domains was incorrect. Manual inspection of the map showed that an electron density
23 for both domains of gp226 was present although was weak, and that the density of a peripheral
24 domain of gp154 was slightly shifted compared to its location in the AlphaFold2 model. We
25 proceeded with fitting the AlphaFold2 models of both gp226 domains into the density and
26 adjusting the location of the peripheral gp154 domain – all as rigid bodies – using Coot¹⁹. A
27 subsequent 20-cycle restrained refinement run with Refmac5²⁰ brought the R-free factor to
28 39%, which resulted in a much better and cleaner electron density in which many of the minor
29 model inaccuracies (some of which are shown in Fig. 4) became obvious and could be easily
30 corrected using a long segment refine/morph procedure implemented in Coot. Further
31 corrections and refinement of the atomic model with Refmac5²⁰ and Phenix³⁶ improved the
32 density and revealed the presence of the DNA oligonucleotide. Subsequent rounds of
33 refinement and model building made the AlphaFold2-derived structure indistinguishable
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

(within the expected accuracy) from that obtained by an MR procedure that used the complete cryo-EM-derived holoenzyme complex structure as a search model.

In conclusion we note that the AlphaFold2 team has clearly developed a methodology to accurately predict the tertiary structure of individual domains not only for proteins for which deep sequence alignments could be built but even for unique proteins, such as AR9 gp226. Furthermore, the structures of multidomain proteins, such as those comprising individual subunits of the β - β' core of the AR9 nvRNAP enzyme, were also predicted with astounding accuracy. This places the AlphaFold2 team within reach of predicting the quaternary structure of larger complexes, and one can argue that they already demonstrated this by the accuracy of their prediction of individual subunits of the AR9 nvRNAP β - β' core that could be assembled into a complex that closely resembles the experimentally determined structure.

3. AlphaFold2 helped correct cis and trans proline assignments and the subsequent tracing of 20 amino acid residues in the crystal structure of the baseplate anchor and partner TSP assembly region of TSP4 from Bacteriophage CBA120 (CASP T1070) – by OH, KC, XS, JG, SBL and DCN

From email to the CASP Prediction Center: *Unbelievable. They predicted residues 16-75 correctly with an RMS of 1.26 Å. Also, the prediction includes a different assignment of a cis proline (P236) than my original assignment. It turned out that the predicted version is correct because it enables repositioning of a tyrosine residue (Y247) in the right place. The change, together with another adjustment ultimately results in a 2-residue shift of 20 residues (237-256). Osnat Herzberg*

3.1. Brief description of the target

As a member of the recently defined *Kuttervirus* genus, the *Escherichia coli* O157:H7 bacteriophage CBA120 infects multiple hosts using four tailspike proteins (TSP1-4). Each TSP

1
2
3 has a distinct endo-glycosidase activity specific to the lipopolysaccharides of different bacterial
4 hosts. The four phage CBA120 TSPs are so far the best characterized, thus they served as a
5 paradigm for understanding the infection mechanism and host range expansion characteristic
6 to the *Kuttervirus* genus. All TSPs assemble into trimers and employ the same overall fold of
7 their catalytic domains (trimers of β -helix subunits). Nevertheless, within this fold, the different
8 active site architectures confer different endo-glycosidase substrate specificities, which in turn
9 facilitates the host range expansion of the phage³⁷⁻⁴⁰. The four TSPs form a complex, seen on
10 negative-stained electron micrographs as a branched appendage emanating from the phage tail
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
41. The 335 N-terminal amino acids of TSP4 mediate this assembly and anchoring function. The
sequence of this region (herewith termed TSP4-N) comprise the target submitted for CASP14
structure prediction (target T1070). The crystal structure of TSP4-N was determined initially at
a resolution limit of 3.2 Å using Single-wavelength Anomalous Dispersion at the Se absorption
edge of crystals containing SeMet protein. This structure served as a Molecular Replacement
search model to determine the crystal structure of the wild-type TSP4-N using crystals that
diffracted to a resolution limit of 2.6 Å. Structure refinement of this crystal form yielded $R =$
0.206 and $R_{\text{free}} = 0.229$.

Consistent with the full-length TSP4, the TSP4-N also assembled into trimers. The
structure revealed four domains connected by flexible linkers. The 75 N-terminal amino acids
comprise the domain that anchors TSP4 to the phage tail baseplate (herewith termed AD). Of
these, approximately 50 amino acid residues fold into an intertwined triple β -helix, which then
disengage to form an antiparallel β -prism II from the ensuing 25 residues, with each subunit
contributing 3-stranded antiparallel β -sheet to the trimer prism (Fig. 5A). This was the most
challenging region for structure prediction because of its lack of sequence homology to
sequences of known protein structure. Following a short linker region, the polypeptide chain
folds into three domains (herewith termed XD1-3) that recruit the partner TSPs. While XD1
exhibits a low but clear sequence identity to a domain of gp9 from phage T4 baseplate (18%

over 95 of 100 shared amino acid residues), XD2 and XD3 exhibit only remote sequence homology to proteins of known crystal structure, which can be detected by Hidden Markov Model methods. Domain XD1 adopts a mixed β -sandwich fold, while both XD2 and XD3 adopt a jellyroll fold. In the crystal structures, whether the trimers employ a crystallographic or non-crystallographic 3-fold symmetry axis, all domains obey the same 3-fold symmetry axis. The XD1 and XD3 monomers form closely packed trimeric assemblies. However, XD2 subunits splay apart and do not interact with one another even though they remain related by the 3-fold symmetry axis. This spatial separation of XD2 subunits prevents binding of a trimeric partner TSP, and is probably a crystal packing artifact. Indeed, a crystal structure of a protein construct lacking the XD3 domain revealed closely packed XD2 subunits, as necessary for binding of a trimeric TSP partner.

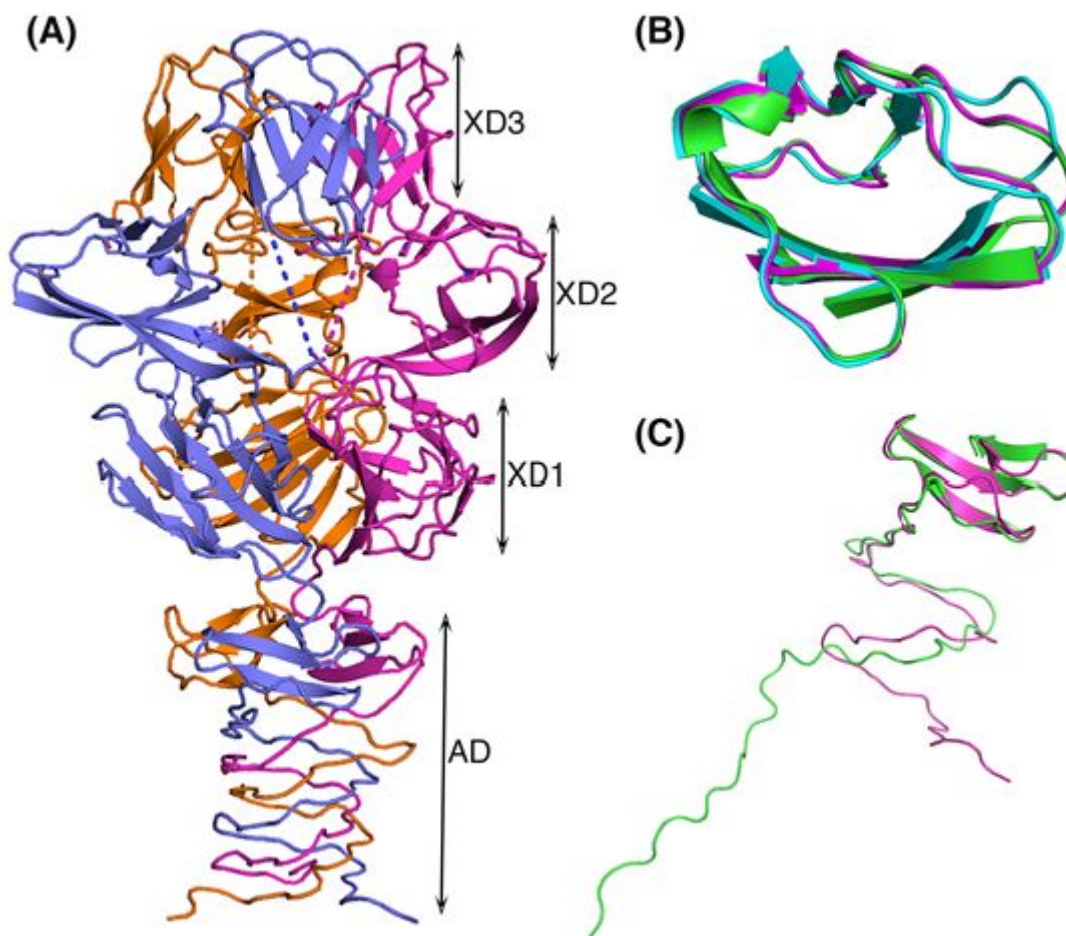


Figure 5.

Figure 5. (A) The structure of TSP4-N homo-trimer with each subunit in different color. The dash lines indicate structurally disordered linkers between XD2 and XD3. (B) Superposition of XD2 as seen in the crystal structure (magenta) and the structure predicted by group 427 (green) and group 226 (sky blue). (C) Superposition of AD crystal structure (magenta) and the structure predicted by group 427 (green).

3.2. How AlphaFold2 models helped improve the structure

The 2.6 Å resolution crystal structure of TSP4-N was determined by Molecular Replacement using a low-resolution structure (3.2 Å) that was initially built using the automated tracing program Autobuild⁴². Although the refinement progressed well, a strong residual difference electron density associated with Ile247 in the XD2 domain suggested that the experimental model required modification. A close examination of the AlphaFold2 model and the crystal structure revealed a polypeptide tracing error due to a wrong assignment of two neighboring proline residues (Pro236 and Pro239) carried over from the initial 3.2 Å structure. To better fit the electron density map, Pro236, located on a tight turn, was assigned a cis conformation, and Pro239 was assigned a trans conformation (Fig. 6A). Guided by the AlphaFold2 model, the two proline conformations were switched (now trans Pro236 and cis Pro239), and their positions shifted. The position of 20 amino acids were adjusted concomitantly so that Tyr249 was placed at the initial Ile247 position, which eliminated the residual difference electron density (Fig. 6B). This example demonstrates that in the future, the highly accurate models generated by AI methods will guide correct interpretations of low-resolution electron density maps generated by x-ray crystallography and cryo-EM, whenever difficulties exist in differentiating between cis and trans peptides.

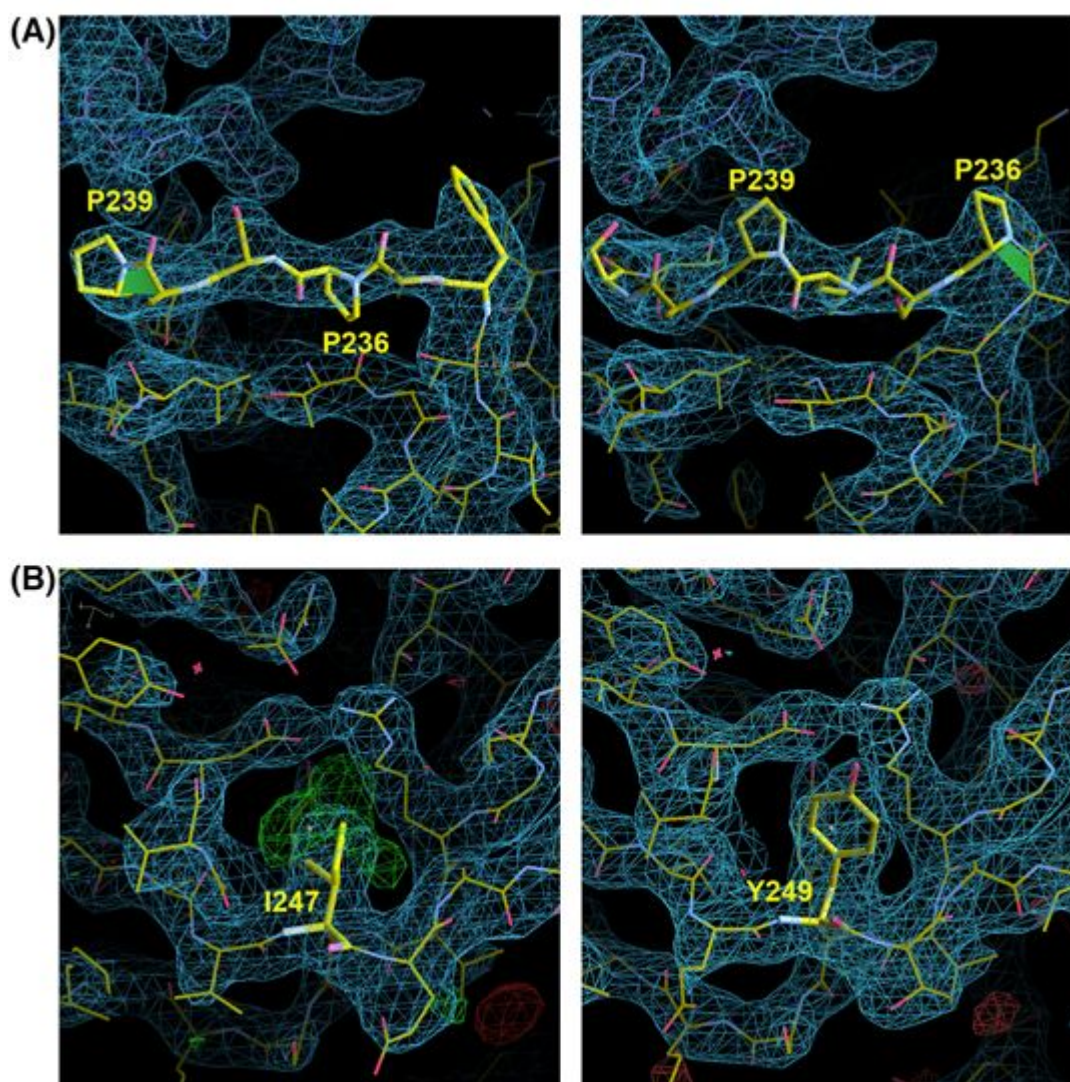


Figure 6.

Figure 6. Polypeptide chain tracing errors that were corrected by examination of the AlphaFold2 (group 427) structure. (A) The incorrect model in the vicinity of two neighboring proline residues (Pro236 and Pro239) together with the associated difference electron density map with the coefficient $2F_o - F_c$ colored blue (left) and the model corrected based on the AlphaFold2 predicted structure with the associated $2F_o - F_c$ difference electron density map (right). The cis bond conformations are highlighted in green (B) The incorrect placement of Ile247 with the associated $2F_o - F_c$ difference electron density map colored blue and the $F_o - F_c$ difference electron density map colored green (left). Correcting the positions of Pro236 and

1
2
3 Pro239 allowed placement of Tyr249 instead of Ile247 and eliminated the residual $F_o - F_c$
4
5 difference electron density (right).
6
7
8
9

10 11 **3.3. Model accuracy** 12 13

14
15 We assessed the CASP14 predictions of TSP4-N by individual domains because the
16
17 flexible interdomain linkers may adopt different conformations than those seen in the crystal
18
19 structures. Several groups predicted the correct folds of domains XD1-XD3, with different level
20
21 of accuracy. Overall, the AlphaFold2 predictions (DeepMind team, group 427) were the most
22
23 accurate with respect to the XD domains. Superposition of the crystal structure and the
24
25 AlphaFold2 model 1 using PyMOL ⁴³ yielded root mean squared deviation (RMSD) values for
26
27 aligned $C\alpha$ atoms of 0.38 Å for XD1 (90 of 100 superposed amino acid residues), 0.36 Å for
28
29 XD2 (52 of 60 superposed amino acid residues) and 0.63 Å for XD3 (63 of 65 superposed
30
31 amino acid residues). Several other CASP14 participants predicted the structures of these
32
33 domains successfully, typically with twice or more the RMSD values. Fig. 5B shows the
34
35 superposition of the XD2 domain, illustrating the remarkable similarity between the
36
37 experimental and Alphafold2 structures, and also an excellent structure similarity produced by
38
39 another group, even though not as good as the AlphaFold2 structure (ZhangTBM server, group
40
41 226, RMSD = 0.7-0.8 Å).
42
43
44
45
46
47

48
49 None of the structure predictions of the AD domain resembled the entire triple β -helix
50
51 region. Nevertheless, the AlphaFold2 model of the AD subunit contains a meandered
52
53 polypeptide chain covering residues 20-50 that resembles the β -helix trace seen in the crystal
54
55 structure, with RMSD value for 28 of the 31 aligned $C\alpha$ atoms of 1.7 Å. Moreover, the
56
57 predicted ensuing 3-stranded antiparallel β -sheet that forms the trimeric antiparallel β -prism II
58
59 (residues 51-75) is quite accurate, with RMSD value for all 25 aligned $C\alpha$ atoms of 0.65 Å. In
60

1
2
3 contrast, the AlphaFold2 residues 1-19 diverge from the experimental structure and the five
4 deposited models exhibit a wide range of extended polypeptide chain orientations. Fig. 5C
5
6 illustrates this by superposing only the closely related 3-stranded antiparallel β -sheet regions of
7
8 the X-ray structure and the AlphaFold2 model 1. Considering that the AD triple β -helix
9
10 polypeptide lacks significant amino acid sequence homology to those of known protein
11
12 structures, and that there are no intra subunit interactions in triple β -helices, it is surprising that
13
14 the fold calculated by the AI methods resembles at all the actual fold.
15
16
17
18
19
20
21
22

23 **4. AlphaFold2 models enable solving crystal structure of Af1503 transmembrane** 24 **receptor (CASP: T1100) that withstood experimental approaches for years – by** 25 **MDH, RA and ANL.** 26 27 28 29

30 From email to the CASP Prediction Center: *I cannot overstate my excitement at the fact that*
31 *Marcus Hartmann solved the structure of Af1503 by molecular replacement with the models*
32 *of group g427. Andrei Lupas*
33
34
35
36

37 **4.1. Brief description of the target** 38

39 Our department has a long-standing interest in coiled coils and their role in
40 transmembrane signal transduction. Coiled coils are bundles of α -helices with a specific regular
41 and repetitive packing ⁴⁴; they are found in innumerable structural contexts in essentially all
42 aspects of cell biology ⁴⁵. While their structural and functional roles are well understood in
43
44 many contexts, their role in transmembrane signal transduction is still debated. Many
45
46 transmembrane receptors are homo-dimeric proteins in which a membrane-spanning coiled coil
47
48 connects extracellular sensor domains to intracellular effector domains, such that signals have
49
50 to be passed along the coiled-coil segment. To study this process, we have been working on the
51
52 minimalistic putative receptor Af1503 from *Archaeoglobus fulgidus* - fortuitously, we had
53
54 already entered its genomic neighbor, Af1502, into the CASP 11 experiment ^{46,47}. Sequence
55
56
57
58
59
60

1
2
3 analysis suggested that Af1503 forms a homo-dimer merely consisting of an extracellular PAS
4 domain connected to an intracellular HAMP domain via an antiparallel tetrameric coiled coil.
5
6 While we conducted several structural studies on the isolated HAMP domain ^{48,49} and on
7
8 chimeric fusion proteins in which we fused the Af1503 HAMP domain to other coiled coil-
9
10 based signaling domains ^{50,51}, we were so far unable to determine the structure of the full
11
12 receptor ⁵².
13
14
15

16 17 18 **4.2. How AlphaFold2 models helped to solve the structure**

19
20 Our problems in obtaining the structure of the full receptor did not lie in the behavior of
21
22 the protein. The protein was very well behaved, stable, and readily crystallized in a range of
23
24 conditions. However, crystal quality was very erratic, could not be improved systematically,
25
26 and diffraction was generally strongly anisotropic and not to high resolution. This led to the
27
28 failure of experimental phasing approaches, despite several different strategies employed. On
29
30 the other hand, molecular replacement (MR) was not successful, as we only had the structure
31
32 of the HAMP domain as an available search model, and as the approach was further complicated
33
34 by the presence of translational non-crystallographic symmetry. To aid MR, we decided to
35
36 tackle a truncated construct covering the extracellular PAS domain, but this construct failed to
37
38 crystallize. In contrast, we succeeded with an NMR analysis of this construct, revealing the fold
39
40 of the PAS domain, but the structural models derived from the NMR data were too far from the
41
42 actual crystal structure to succeed in MR attempts.
43
44
45
46
47

48
49 Finally, years later, we easily managed to solve the crystal structure using the
50
51 AlphaFold2 models. As the predictions were modeled as monomers, without constraints for the
52
53 homo-dimeric state, they were not fully compatible with the dimeric state along the whole
54
55 chain, and a very first, naive MR attempt employing a full model did not succeed. However, in
56
57 the second attempt, only employing a single PAS domain with a short coiled-coil segment as a
58
59 search model, the structure was essentially solved using MOLREP with standard parameters ⁵³;
60

after the correct placement and initial refinement of the PAS domains, the electron density for the rest of the protein was clearly traceable.

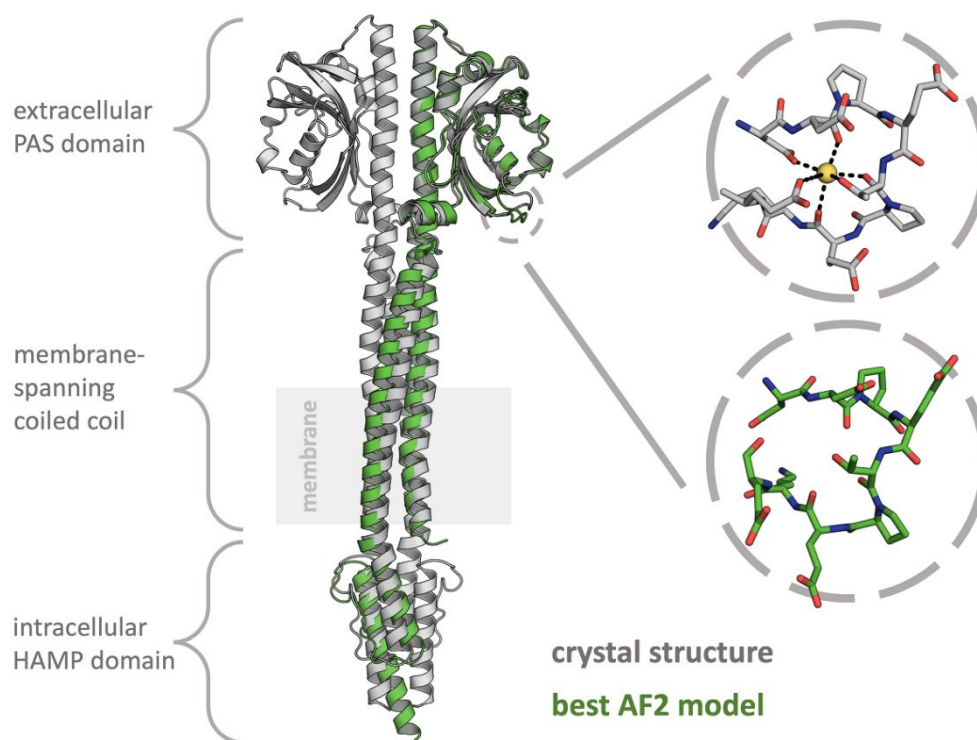


Figure 7

Figure 7. The crystal structure of dimeric Af1503 (grey) is shown in a superposition with the best AlphaFold2 model (green, monomer). The only noteworthy difference between the prediction and the crystal structure is found in a loop in the PAS domain, which was found to coordinate an ion in the crystal structure.

4.3. Model accuracy

Although the AlphaFold2 predictions were modeled as monomers that are not fully compatible with the dimeric state of Af1503, the predicted models superimpose closely on the final crystal structure (Fig. 7). Of the five AlphaFold2 models, four are in a conformation that closely matches the dimeric state, and all of them superimpose with an RMSD below 2.5 Å over their full length on all chains of the crystal structure. Consequently, more focused, local

1
2
3 superimpositions yield RMSD values far below 2 Å. In short, the model accuracy is fairly close
4
5 to what one would expect for another crystal structure of the same protein. There is just one
6
7 region that deviates from the crystal structure: The electron density revealed that an elongated
8
9 loop within the PAS domain is actually coordinating a metal ion, which has a pronounced
10
11 impact on its structure. Needless to say, AlphaFold2 did not predict the presence and
12
13 coordination of that ion, but nevertheless, it predicted this loop in a conformation that is at least
14
15 close to the ion-bound state.
16
17
18
19

20 **5. AlphaFold2 models aid in crystal structure determination of the bacterial exo-** 21 22 **sialidase Sia24 (CASP: T1089) by molecular replacement – by SDR and GAC.** 23 24

25 From email to the CASP Prediction Center: *Models 1, 2, 3, and 5 worked quite well as an*
26
27 *ensemble for molecular replacement, and quite well on their own. We eventually achieved*
28
29 *similar results with an ensemble of current PDB models, but this one scored much higher in*
30
31 *MR from the beginning. Steven Rees*
32
33

34 **5.1. Brief description of the target** 35

36 Sialidase enzymes (or neuraminidases) cleave sialic acid (SA) moieties found on mucin
37
38 glycoproteins of the gastrointestinal (GI) tract, and are utilized by microbial communities for
39
40 the sequestration of SAs as metabolic substrates, or (in the case of some pathogenic species) a
41
42 means of biofilm formation, surface adhesion, and revealing toxin-binding sites^{54,55}. Exo-
43
44 sialidases, which cleave terminal SAs, are typically classified in the carbohydrate-active
45
46 enzymes database (CAZy) as GH family 33 (GH33) and are the most common sialidases
47
48 identified^{54,56,57}, typically utilizing a two-step catalytic mechanism where a conserved Glu
49
50 activates a spatially proximal Tyr for nucleophilic attack of C5 of the SA, prompting acid-base
51
52 catalysis at C5 by an Asp residue^{58,59}. While most sialidases characterized to date are
53
54 ambivalent towards the mammalian SAs Neu5Ac and Neu5Gc (differing only by a hydroxyl
55
56 group at the acetamido C5 on the latter), we and others characterized a series of Neu5Gc-
57
58
59
60

1
2
3 favoring sialidases in both the microbial communities of mice fed Neu5Gc-enriched diets and
4
5 a human population during Neu5Gc-enriched dietary seasons ⁵⁹. This study identified an
6
7 upregulation of Sia24, a Neu5Gc-favoring sialidase likely from *Bacteroides acidifaciens* with
8
9 low sequence homology to published sialidase structures.
10
11
12

13 **5.2. Methodology**

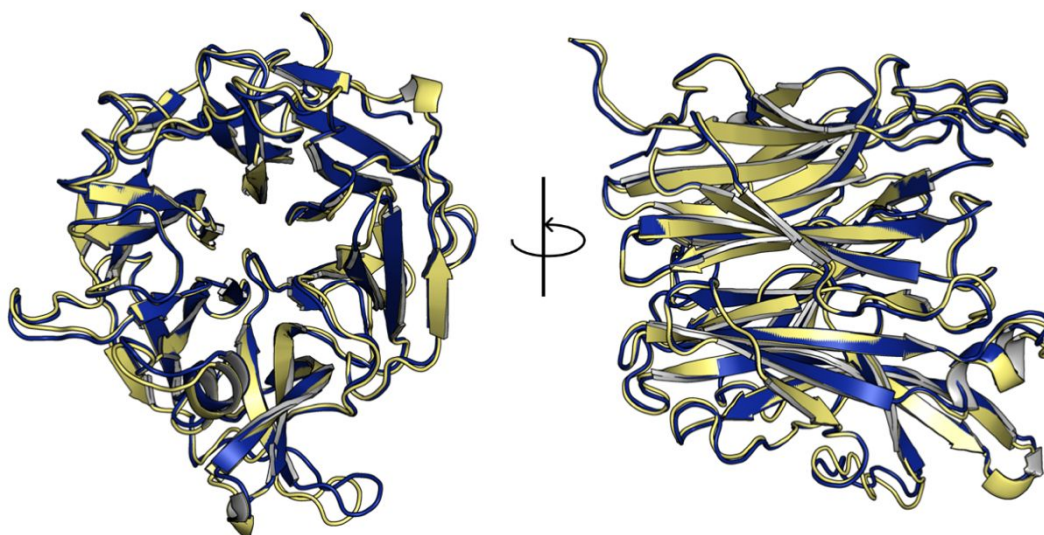
14
15 Sia24 was purified and concentrated to 10-12 mg/mL ⁵⁹, and crystallized in 100 mM Bis-
16
17 Tris pH 6.5 and 20% polyethylene glycol monomethyl ether 5,000. Crystals in the P4₁ space
18
19 group typically diffracted to 2.2-2.6 Å, with a single high-resolution dataset collected at 2.0 Å.
20
21 A more detailed description of the protein production and crystallization are provided in the
22
23 Supplementary Material, and will be presented in a future study.
24
25

26
27 Our initial molecular replacement attempts used cross-species homolog structures
28
29 identified by sequence-based searches in the PDB. These searches focused on using the catalytic
30
31 domain of exo-sialidase models derived from the GH33 family, as Sia24 lacks the carbohydrate-
32
33 binding motif found in some members. Various identified catalytic domain search models (from
34
35 PDB accession codes 1DIL, 1EUR, 1WCQ, 2VK5, 4FJ6, 4J9T, 4BBW, 4Q6K, and 5TSP)
36
37 initially failed to find a reasonable phasing solution by molecular replacement regardless of
38
39 model modification (e.g., poly-alanine, CCP4's Chainsaw-mediated side-chain pruning and
40
41 mutagenesis, and removal of flexible loop regions outside of canonical beta-propeller domain
42
43 secondary structure). *Ab initio* models generated by Robetta (<https://robetta.bakerlab.org/>) and
44
45 I-TASSER ⁶⁰⁻⁶² did not yield a solution by molecular replacement. Phyre2 ⁶³ offered reasonable
46
47 solutions (TFZ=14.1, LLG=194), as did using PHENIX.ENSEMBLER to generate an ensemble
48
49 of the nine models mentioned above (TFZ=16.9, LLG=256). Both of these approaches
50
51 struggled during subsequent refinement and manual building steps, and the latter ensemble
52
53 models lacked much of the Sia24 sequence because of low homology. Concurrently, we tried
54
55 models of Sia24 generated by the AlphaFold2 team and provided by the CASP14 organizers.
56
57
58
59
60 Four of their five coordinate models were quickly successful in initial phase estimation by

1
2
3 molecular replacement after removal of flexible N- and C-terminal regions, with model 2
4
5 (T1089TS427_2) showing the highest performance (TFZ=62.5, LLG=3791).
6
7

8 **5.3. Model accuracy**

9
10 AlphaFold2's model had high coordinate similarity (RMSD=1.08 Å on all atoms, 0.55
11 Å after 5 cycles of outlier rejection) to the crystallographic structure (PDB code 7MHU), and
12
13 displays the beta-propeller structure of the canonical exo-sialidase catalytic domain (Fig. 8).
14
15 Most side-chains are also reasonably oriented. The largest deviations in the models were
16
17 localized to the N- and C-termini and regions between anti-parallel beta strand propeller motifs.
18
19 The low-homology model ensemble described above has a similar consistency, albeit lacking
20
21 most information on side-chain and flexible loop placement. Similarly, models from other *ab*
22
23 *initio* methods display reasonable overlap, but were not successful in initial molecular
24
25 replacement attempts.
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52



53 **Figure 8**

54
55 **Figure 8.** Superposition of Sia24 predictive and crystallographic models. The structure of Sia24
56
57 (dark blue) was solved with initial phase determination by molecular replacement using a model
58
59 generated by AlphaFold2 (yellow).
60

1
2
3
4
5
6 Given the initial results for both Phyre2 and ensembled low-homology models in
7
8 molecular replacement, a solution for Sia24 without AlphaFold2 would have likely been
9
10 eventually determined. However, the greatly improved accuracy of the AlphaFold2 models,
11
12 with all side chains and flexible loops in place, is undeniable. Sia24 exhibits roughly 20%
13
14 sequence homology with previously known structures, which in most cases is at the threshold
15
16 of likely success with molecular replacement for most targets. AlphaFold2 was able to use
17
18 information from known structures in a novel way compared to previous algorithms, and
19
20 provide an effective solution where the alternatives struggled.
21
22
23
24
25
26
27

28 **Discussion**

29
30 This paper describes the solution of four experimental structures using molecular
31
32 replacement with models submitted to CASP14. We also report improvement of an already
33
34 solved target using models.
35
36
37

38
39 Molecular replacement is a very well established technique, but high accuracy models
40
41 are needed, and until now that has almost always required the availability of templates based
42
43 on high levels of sequence identity ⁶⁴. The three most recent CASPs have seen dramatic
44
45 improvements in the accuracy of non-homologous models, first from the successful application
46
47 of 3D contact prediction methods using statistical approaches ⁶⁵ and then from the use of deep
48
49 learning methods ^{13,66}. In CASP14, the Alphafold2 group submitted models for many targets
50
51 that rival the corresponding experimental structures in accuracy ¹⁰⁻¹³. The difficulties in
52
53 obtaining experimental structures for seven of the CASP14 targets provided an opportunity to
54
55 objectively test the ability of new methods in this respect. As the accounts in this paper show,
56
57 the models are indeed powerful.
58
59
60

1
2
3 A post-CASP analysis by Randy Read and colleagues ¹⁵ found that all CASP14 targets
4 with available diffraction data could be solved or at least partially solved using molecular
5 replacement with Alphafold2 models. The analysis implies that in future, structure modeling
6
7
8
9
10 should be a viable means of solving all but the most challenging crystal structures with
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Acknowledgements

Initials of the authors contributing to specific sections are provided in the sections' titles.

Idea, concept, abstract, introduction, discussion, editing and coordination - by AK and JM.

CASP experiment is supported by the US National Institute of General Medical Sciences (NIGMS/NIH), grant number GM100482.

OH and DNC studies were supported by US National Institute of General Medical Sciences (NIGMS/NIH), grant number RO1GM110202.

FIGURE CAPTIONS

Figure 1. (A) Partial FoxB model obtained by experimental phasing before the CASP14 model became available. At this point the model could not be further improved and the project was stuck for a year. (B) Experimental phases with partial FoxB model (map shown at 1.2σ level).

Figure 2. Workflow of FoxB structure determination. The structure was determined by MR-SAD using the AlphaFold2 model and experimental phases. (A) Anomalous difference map with Se and Fe sites at 2σ . (B) Overall map of FoxB after refinement (2σ). (C) Superposition

1
2
3 of the final model (green) and AlphaFold2 model (cyan) shows excellent agreement. Density
4
5 for heme groups (not present in AlphaFold2 model) is shown.
6
7
8
9

10
11 **Figure 3.** AlphaFold2 models of AR9 nvRNAP proteins fit the cryo-EM density nearly
12 perfectly. The cryo-EM-derived structures of gp105, gp154, and two gp226 domains are
13 colored according to the color code given in the upper left corner of each panel. All AlphaFold2
14 models are colored magenta. The electron density is contoured at 4.25 standard deviations
15 above the mean and colored semi-transparent grey. Regions where no cryo-EM-derived
16 structure existed prior to the availability of the AlphaFold2 models are indicated with a dashed
17 line and their boundary residues are labeled.
18
19
20
21
22
23
24
25
26
27
28
29
30

31 **Figure 4.** Inaccuracies in AlphaFold2 models. Cryo-EM-derived structures and AlphaFold2
32 models of several AR9 nvRNAP subunits are superimposed and regions where the
33 conformation of the AlphaFold2 model deviates significantly from the cryo-EM-derived
34 structure are indicated with a dashed line and their boundary residues are labeled. Note that the
35 folds of both the N- and C-terminal domains of gp226 were predicted correctly, but the structure
36 of the interdomain linker and the relative orientation of the two domains were incorrect.
37
38
39
40
41
42
43
44
45
46
47
48

49 **Figure 5.** (A) The structure of TSP4-N homo-trimer with each subunit in different color. The
50 dash lines indicate structurally disordered linkers between XD2 and XD3. (B) Superposition of
51 XD2 as seen in the crystal structure (magenta) and the structure predicted by group 427 (green)
52 and group 226 (sky blue). (C) Superposition of AD crystal structure (magenta) and the structure
53 predicted by group 427 (green).
54
55
56
57
58
59
60

1
2
3
4
5
6 **Figure 6.** Polypeptide chain tracing errors that were corrected by examination of the
7 AlphaFold2 (group 427) structure. (A) The incorrect model in the vicinity of two neighboring
8 proline residues (Pro236 and Pro239) together with the associated difference electron density
9 map with the coefficient $2F_o-F_c$ colored blue (left) and the model corrected based on the
10 AlphaFold2 predicted structure with the associated $2F_o-F_c$ difference electron density map
11 (right). The cis bond conformations are highlighted in green (B) The incorrect placement of
12 Ile247 with the associated $2F_o-F_c$ difference electron density map colored blue and the F_o-F_c
13 difference electron density map colored green (left). Correcting the positions of Pro236 and
14 Pro239 allowed placement of Tyr249 instead of Ile247 and eliminated the residual F_o-F_c
15 difference electron density (right).
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32

33 **Figure 7.** The crystal structure of dimeric Af1503 (grey) is shown in a superposition with the
34 best AlphaFold2 model (green, monomer). The only noteworthy difference between the
35 prediction and the crystal structure is found in a loop in the PAS domain, which was found to
36 coordinate an ion in the crystal structure.
37
38
39
40
41
42
43
44
45

46 **Figure 8.** Superposition of Sia24 predictive and crystallographic models. The structure of Sia24
47 (dark blue) was solved with initial phase determination by molecular replacement using a model
48 generated by AlphaFold2 (yellow).
49
50
51
52
53
54
55

56 REFERENCES

- 57
58
59
60 1. Schwede T, Sali A, Honig B, et al. Outcome of a workshop on applications of protein models in biomedical research. *Structure*. 2009;17(2):151-159.

- 1
 - 2
 - 3
 - 4
 - 5
 - 6
 - 7
 - 8
 - 9
 - 10
 - 11
 - 12
 - 13
 - 14
 - 15
 - 16
 - 17
 - 18
 - 19
 - 20
 - 21
 - 22
 - 23
 - 24
 - 25
 - 26
 - 27
 - 28
 - 29
 - 30
 - 31
 - 32
 - 33
 - 34
 - 35
 - 36
 - 37
 - 38
 - 39
 - 40
 - 41
 - 42
 - 43
 - 44
 - 45
 - 46
 - 47
 - 48
 - 49
 - 50
 - 51
 - 52
 - 53
 - 54
 - 55
 - 56
 - 57
 - 58
 - 59
 - 60
2. Tramontano A. The role of molecular modelling in biomedical research. *FEBS Lett.* 2006;580(12):2928-2934.
 3. Kuhlman B, Bradley P. Advances in protein structure prediction and design. *Nat Rev Mol Cell Biol.* 2019;20(11):681-697.
 4. Simpkin AJ, Thomas JMH, Simkovic F, Keegan RM, Rigden DJ. Molecular replacement using structure predictions from databases. *Acta Crystallogr D Struct Biol.* 2019;75(Pt 12):1051-1062.
 5. Case DA. Using quantum chemistry to estimate chemical shifts in biomolecules. *Biophys Chem.* 2020;267:106476.
 6. Thompson JM, Sgourakis NG, Liu G, et al. Accurate protein structure modeling using sparse NMR data and homologous structure information. *Proc Natl Acad Sci U S A.* 2012;109(25):9875-9880.
 7. Taylor NM, Prokhorov NS, Guerrero-Ferreira RC, et al. Structure of the T4 baseplate and its function in triggering sheath contraction. *Nature.* 2016;533(7603):346-352.
 8. Malhotra S, Trager S, Dal Peraro M, Topf M. Modelling structures in cryo-EM maps. *Curr Opin Struct Biol.* 2019;58:105-114.
 9. Kryshtafovych A, Malhotra S, Monastyrskyy B, et al. Cryo-electron microscopy targets in CASP13: Overview and evaluation of results. *Proteins.* 2019;87(12):1128-1140.
 10. Jumper J, Hassabis D, et al. AlphaFold2 in CASP14. *Proteins.* 2021(This issue).
 11. Hartmann M, Lupas A, et al. Evaluation of high accuracy modeling in CASP14. *Proteins.* 2021(This issue).
 12. Kinch L, Grishin N, et al. CASP14 Topology Evaluation of Difficult Targets. *Proteins.* 2021(This issue).
 13. Kryshtafovych A, Moulton J, et al. Critical Assessment of protein Structure Prediction, round 14 - tentative name. *Proteins.* 2021(This issue).
 14. Garcia-Alai MM, Heidemann J, Skruzny M, et al. Epsin and Sla2 form assemblies through phospholipid interfaces. *Nat Commun.* 2018;9(1):328.
 15. Read RJ, Rigden DJ, Lupas A, Hartmann M, et al. Effectiveness of molecular replacement on CASP14 targets. *Proteins.* 2021;This issue.
 16. Josts I, Veith K, Tidow H. Ternary structure of the outer membrane transporter FoxA with resolved signalling domain provides insights into TonB-mediated siderophore uptake. *Elife.* 2019;8.
 17. Josts I, Veith K, Normant V, Schalk I, H. T. Structural insights into a novel family of integral membrane siderophore reductases. *BioRxiv.* 2021;2021.01.28.428567.
 18. Veith K, Martinez Molledo M, Almeida Hernandez Y, et al. Lipid-like Peptides can Stabilize Integral Membrane Proteins for Biophysical and Structural Studies. *ChemBiochem.* 2017;18(17):1735-1742.
 19. Emsley P, Lohkamp B, Scott WG, Cowtan K. Features and development of Coot. *Acta Crystallogr D Biol Crystallogr.* 2010;66(Pt 4):486-501.
 20. Murshudov GN, Skubak P, Lebedev AA, et al. REFMAC5 for the refinement of macromolecular crystal structures. *Acta Crystallogr D Biol Crystallogr.* 2011;67(Pt 4):355-367.
 21. Lavysh D, Sokolova M, Minakhin L, et al. The genome of AR9, a giant transducing Bacillus phage encoding two multisubunit RNA polymerases. *Virology.* 2016;495:185-196.
 22. Lavysh D, Sokolova M, Slashcheva M, Forstner KU, Severinov K. Transcription Profiling of Bacillus subtilis Cells Infected with AR9, a Giant Phage Encoding Two Multisubunit RNA Polymerases. *mBio.* 2017;8(1).
 23. Sokolova ML, Misovetec I, K VS. Multisubunit RNA Polymerases of Jumbo Bacteriophages. *Viruses.* 2020;12(10).
 24. Ceyssens PJ, Minakhin L, Van den Bossche A, et al. Development of giant bacteriophage varphiKZ is independent of the host transcription apparatus. *J Virol.* 2014;88(18):10501-10510.
 25. Lee J, Borukhov S. Bacterial RNA Polymerase-DNA Interaction-The Driving Force of Gene Expression and the Target for Drug Action. *Front Mol Biosci.* 2016;3:73.

- 1
 - 2
 - 3
 - 4
 - 5
 - 6
 - 7
 - 8
 - 9
 - 10
 - 11
 - 12
 - 13
 - 14
 - 15
 - 16
 - 17
 - 18
 - 19
 - 20
 - 21
 - 22
 - 23
 - 24
 - 25
 - 26
 - 27
 - 28
 - 29
 - 30
 - 31
 - 32
 - 33
 - 34
 - 35
 - 36
 - 37
 - 38
 - 39
 - 40
 - 41
 - 42
 - 43
 - 44
 - 45
 - 46
 - 47
 - 48
 - 49
 - 50
 - 51
 - 52
 - 53
 - 54
 - 55
 - 56
 - 57
 - 58
 - 59
 - 60
26. Sokolova M, Borukhov S, Lavysh D, Artamonova T, Khodorkovskii M, Severinov K. A non-canonical multisubunit RNA polymerase encoded by the AR9 phage recognizes the template strand of its uracil-containing promoters. *Nucleic Acids Res.* 2017;45(10):5958-5967.
27. Punjani A, Rubinstein JL, Fleet DJ, Brubaker MA. cryoSPARC: algorithms for rapid unsupervised cryo-EM structure determination. *Nat Methods.* 2017;14(3):290-296.
28. Blow DM, Rossmann MG, Jeffery BA. The Arrangement of Alpha-Chymotrypsin Molecules in the Monoclinic Crystal Form. *J Mol Biol.* 1964;8:65-78.
29. Cowtan K. Recent developments in classical density modification. *Acta crystallographica Section D, Biological crystallography.* 2010;66(Pt 4):470-478.
30. McCoy AJ, Grosse-Kunstleve RW, Adams PD, Winn MD, Storoni LC, Read RJ. Phaser crystallographic software. *J Appl Crystallogr.* 2007;40(Pt 4):658-674.
31. Wang BC. Resolution of phase ambiguity in macromolecular crystallography. *Methods Enzymol.* 1985;115:90-112.
32. Hendrickson WA. Determination of macromolecular structures from anomalous diffraction of synchrotron radiation. *Science.* 1991;254(5028):51-58.
33. Read RJ, McCoy AJ. Using SAD data in Phaser. *Acta crystallographica Section D, Biological crystallography.* 2011;67(Pt 4):338-344.
34. Zimmermann L, Stephens A, Nam SZ, et al. A Completely Reimplemented MPI Bioinformatics Toolkit with a New HHpred Server at its Core. *J Mol Biol.* 2018;430(15):2237-2243.
35. Li L, Fang C, Zhuang N, Wang T, Zhang Y. Structural basis for transcription initiation by bacterial ECF sigma factors. *Nature communications.* 2019;10(1):1153.
36. Adams PD, Afonine PV, Bunkoczi G, et al. PHENIX: a comprehensive Python-based system for macromolecular structure solution. *Acta Crystallogr D Biol Crystallogr.* 2010;66(Pt 2):213-221.
37. Chen C, Bales P, Greenfield J, Heselpoth RD, Nelson DC, Herzberg O. Crystal structure of ORF210 from E. coli O157:H1 phage CBA120 (TSP1), a putative tailspike protein. *PLoS One.* 2014;9(3):e93156.
38. Greenfield J, Shang X, Luo H, et al. Structure and tailspike glycosidase machinery of ORF212 from E. coli O157:H7 phage CBA120 (TSP3). *Sci Rep.* 2019;9(1):7349.
39. Greenfield J, Shang X, Luo H, et al. Structure and function of bacteriophage CBA120 ORF211 (TSP2), the determinant of phage specificity towards E. coli O157:H7. *Sci Rep.* 2020;10(1):15402.
40. Plattner M, Shneider MM, Arbatsky NP, et al. Structure and Function of the Branched Receptor-Binding Complex of Bacteriophage CBA120. *Journal of molecular biology.* 2019;431(19):3718-3739.
41. Adriaenssens EM, Ackermann H-W, Anany H, et al. A suggested new bacteriophage genus: "Viunaliikevirus". *Archives of Virology.* 2012;157(10):2035-2046.
42. Terwilliger TC, Grosse-Kunstleve RW, Afonine PV, et al. Iterative model building, structure refinement and density modification with the PHENIX AutoBuild wizard. *Acta Crystallogr D Biol Crystallogr.* 2008;64(Pt 1):61-69.
43. *The PyMOL Molecular Graphics System* [computer program]. Palo Alto, CA, USA: DeLano Scientific; 2002.
44. Lupas AN, Bassler J, Dunin-Horkawicz S. The Structure and Topology of alpha-Helical Coiled Coils. *Subcell Biochem.* 2017;82:95-129.
45. Hartmann MD. Functional and Structural Roles of Coiled Coils. *Subcell Biochem.* 2017;82:63-93.
46. Korycinski M, Albrecht R, Ursinus A, et al. STAC--A New Domain Associated with Transmembrane Solute Transport and Two-Component Signal Transduction Systems. *J Mol Biol.* 2015;427(20):3327-3339.
47. Kryshatfovych A, Moulton J, Basle A, et al. Some of the most interesting CASP11 targets through the eyes of their authors. *Proteins.* 2016;84 Suppl 1:34-50.
48. Ferris HU, Dunin-Horkawicz S, Mondejar LG, et al. The mechanisms of HAMP-mediated signaling in transmembrane receptors. *Structure.* 2011;19(3):378-385.

- 1
 - 2
 - 3
 - 4
 - 5
 - 6
 - 7
 - 8
 - 9
 - 10
 - 11
 - 12
 - 13
 - 14
 - 15
 - 16
 - 17
 - 18
 - 19
 - 20
 - 21
 - 22
 - 23
 - 24
 - 25
 - 26
 - 27
 - 28
 - 29
 - 30
 - 31
 - 32
 - 33
 - 34
 - 35
 - 36
 - 37
 - 38
 - 39
 - 40
 - 41
 - 42
 - 43
 - 44
 - 45
 - 46
 - 47
 - 48
 - 49
 - 50
 - 51
 - 52
 - 53
 - 54
 - 55
 - 56
 - 57
 - 58
 - 59
 - 60
49. Hulko M, Berndt F, Gruber M, et al. The HAMP domain structure implies helix rotation in transmembrane signaling. *Cell*. 2006;126(5):929-940.
50. Ferris HU, Coles M, Lupas AN, Hartmann MD. Crystallographic snapshot of the Escherichia coli EnvZ histidine kinase in an active conformation. *J Struct Biol*. 2014;186(3):376-379.
51. Ferris HU, Dunin-Horkawicz S, Hornig N, et al. Mechanism of regulation of receptor histidine kinases. *Structure*. 2012;20(1):56-66.
52. Hartmann MD, Dunin-Horkawicz S, Hulko M, Martin J, Coles M, Lupas AN. A soluble mutant of the transmembrane receptor Af1503 features strong changes in coiled-coil periodicity. *J Struct Biol*. 2014;186(3):357-366.
53. Vagin A, Teplyakov A. Molecular replacement with MOLREP. *Acta Crystallogr D Biol Crystallogr*. 2010;66(Pt 1):22-25.
54. Juge N, Tailford L, Owen CD. Sialidases from gut bacteria: a mini-review. *Biochem Soc Trans*. 2016;44(1):166-175.
55. Lewis AL, Lewis WG. Host sialoglycans and bacterial sialidases: a mucosal perspective. *Cell Microbiol*. 2012;14(8):1174-1182.
56. Cantarel BL, Coutinho PM, Rancurel C, Bernard T, Lombard V, Henrissat B. The Carbohydrate-Active EnZymes database (CAZy): an expert resource for Glycogenomics. *Nucleic Acids Res*. 2009;37(Database issue):D233-238.
57. Lombard V, Golaconda Ramulu H, Drula E, Coutinho PM, Henrissat B. The carbohydrate-active enzymes database (CAZy) in 2013. *Nucleic Acids Res*. 2014;42(Database issue):D490-495.
58. Amaya MF, Watts AG, Damager I, et al. Structural insights into the catalytic mechanism of Trypanosoma cruzi trans-sialidase. *Structure*. 2004;12(5):775-784.
59. Zaramela LS, Martino C, Alisson-Silva F, et al. Gut bacteria responding to dietary change encode sialidases that exhibit preference for red meat-associated carbohydrates. *Nat Microbiol*. 2019;4(12):2082-2089.
60. Roy A, Kucukural A, Zhang Y. I-TASSER: a unified platform for automated protein structure and function prediction. *Nature Protocols*. 2010;5(4):725-738.
61. Yang J, Yan R, Roy A, Xu D, Poisson J, Zhang Y. The I-TASSER Suite: protein structure and function prediction. *Nature Methods*. 2015;12(1):7-8.
62. Zhang Y. I-TASSER server for protein 3D structure prediction. *BMC Bioinformatics*. 2008;9:40.
63. Kelley LA, Mezulis S, Yates CM, Wass MN, Sternberg MJE. The Phyre2 web portal for protein modeling, prediction and analysis. *Nature Protocols*. 2015;10(6):845-858.
64. Uson I, Ballard CC, Keegan RM, Read RJ. Integrated, rational molecular replacement. *Acta Crystallogr D Struct Biol*. 2021;77(Pt 2):129-130.
65. Moulton J, Fidelis K, Kryshtafovych A, Schwede T, Tramontano A. Critical assessment of methods of protein structure prediction (CASP)-Round XII. *Proteins*. 2018;86 Suppl 1:7-15.
66. Kryshtafovych A, Schwede T, Topf M, Fidelis K, Moulton J. Critical assessment of methods of protein structure prediction (CASP)-Round XIII. *Proteins*. 2019;87(12):1011-1020.

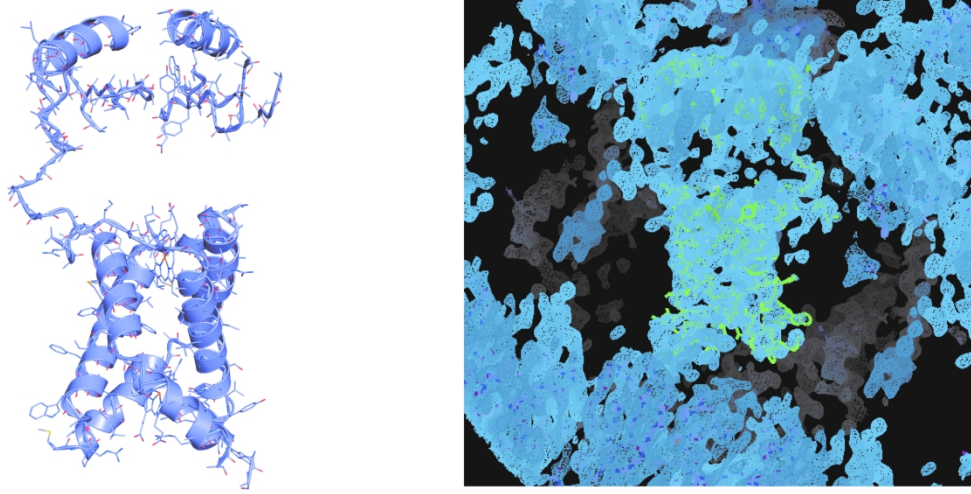


Figure 1. (A) Partial FoxB model obtained by experimental phasing before the CASP14 model became available. At this point the model could not be further improved and the project was stuck for a year. (B) Experimental phases with partial FoxB model (map shown at 1.2σ level).

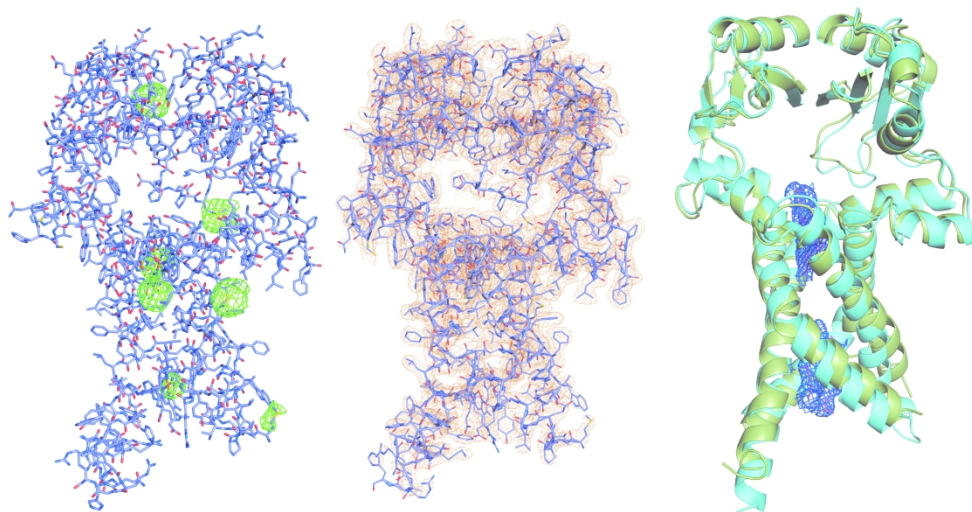


Figure 2. Workflow of FoxB structure determination. The structure was determined by MR-SAD using the AlphaFold2 model and experimental phases. (A) Anomalous difference map with Se and Fe sites at 2σ . (B) Overall map of FoxB after refinement (2σ). (C) Superposition of the final model (green) and AlphaFold2 model (cyan) shows excellent agreement. Density for heme groups (not present in AlphaFold2 model) is shown.

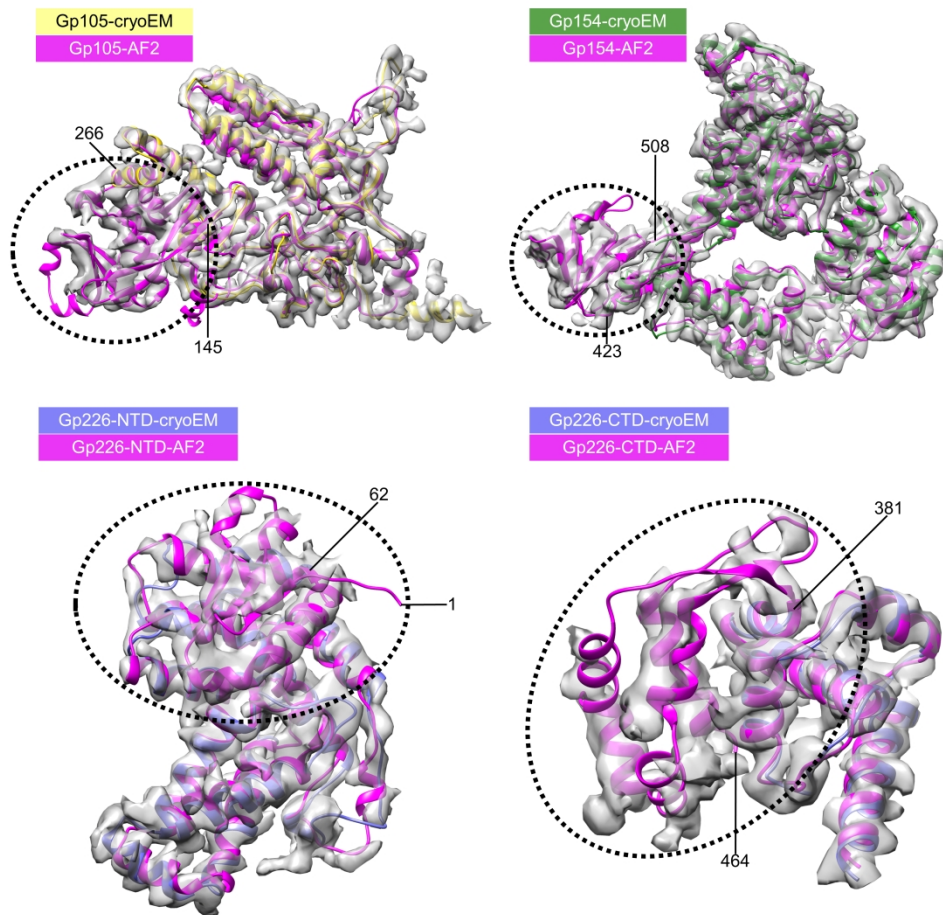


Figure 3. AlphaFold2 models of AR9 nvRNAP proteins fit the cryo-EM density nearly perfectly. The cryo-EM-derived structures of gp105, gp154, and two gp226 domains are colored according to the color code given in the upper left corner of each panel. All AlphaFold2 models are colored magenta. The electron density is contoured at 4.25 standard deviations above the mean and colored semi-transparent grey. Regions where no cryo-EM-derived structure existed prior to the availability of the AlphaFold2 models are indicated with a dashed line and their boundary residues are labeled.

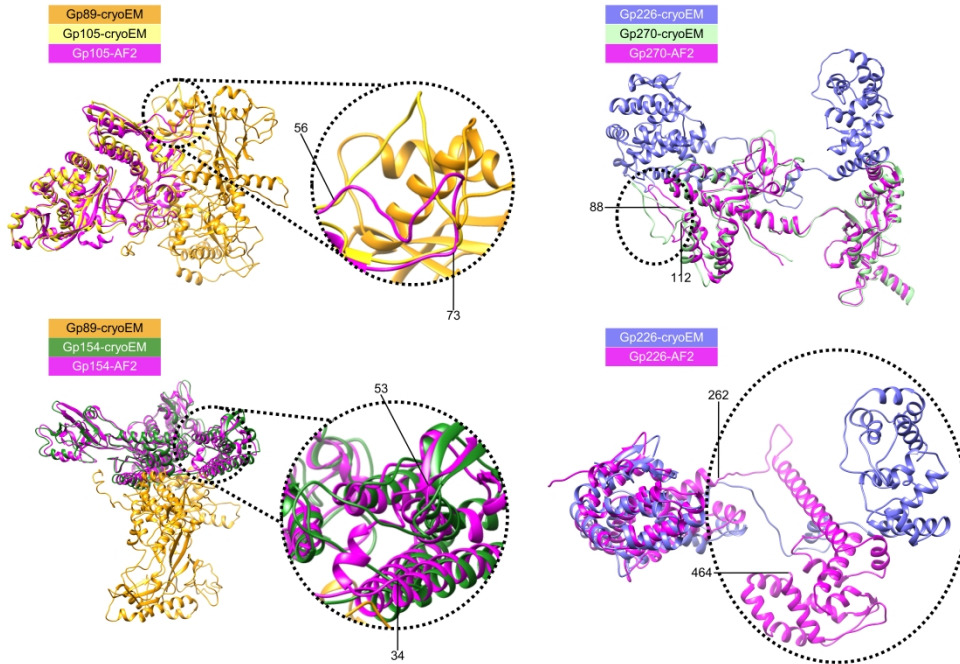


Figure 4. Inaccuracies in AlphaFold2 models. Cryo-EM-derived structures and AlphaFold2 models of several AR9 nvRNAP subunits are superimposed and regions where the conformation of the AlphaFold2 model deviates significantly from the cryo-EM-derived structure are indicated with a dashed line and their boundary residues are labeled. Note that the folds of both the N- and C-terminal domains of gp226 were predicted correctly, but the structure of the interdomain linker and the relative orientation of the two domains were incorrect.

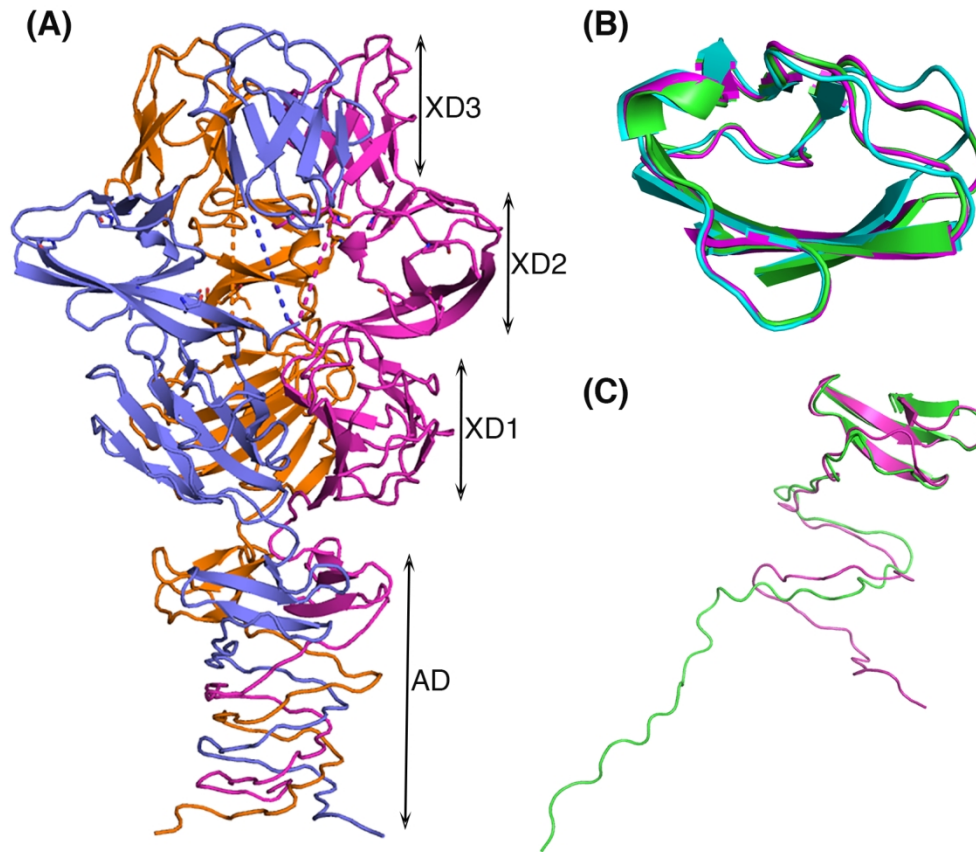


Figure 5. (A) The structure of TSP4-N homo-trimer with each subunit in different color. The dash lines indicate structurally disordered linkers between XD2 and XD3. (B) Superposition of XD2 as seen in the crystal structure (magenta) and the structure predicted by group 427 (green) and group 226 (sky blue). (C) Superposition of AD crystal structure (magenta) and the structure predicted by group 427 (green).

145x128mm (300 x 300 DPI)

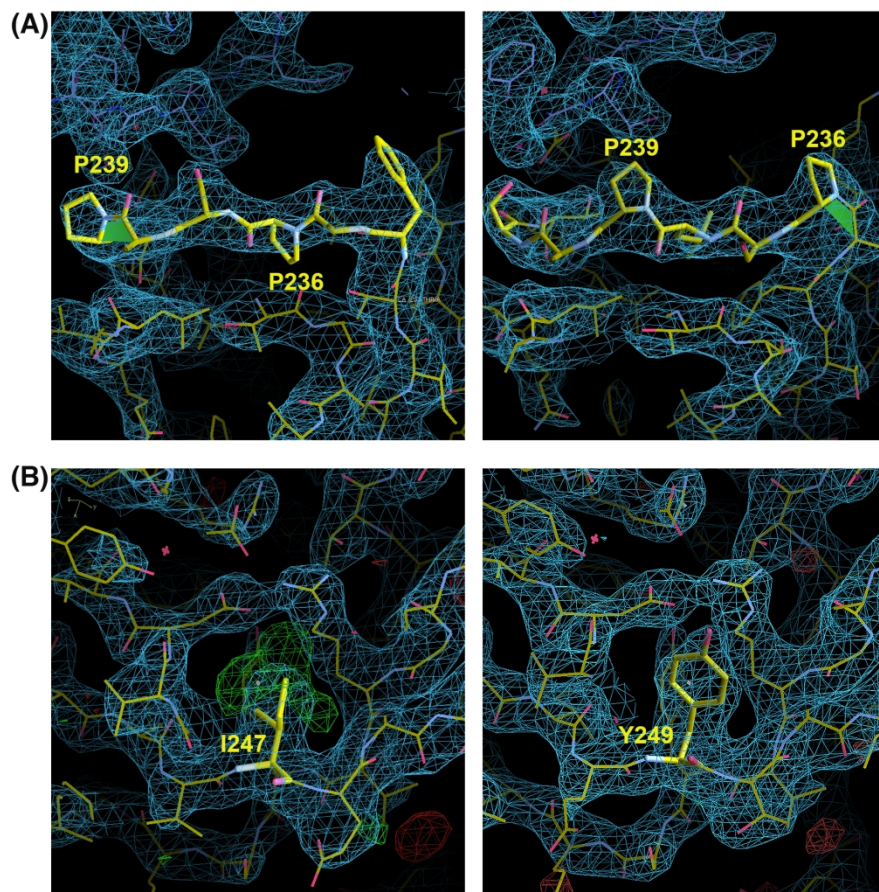
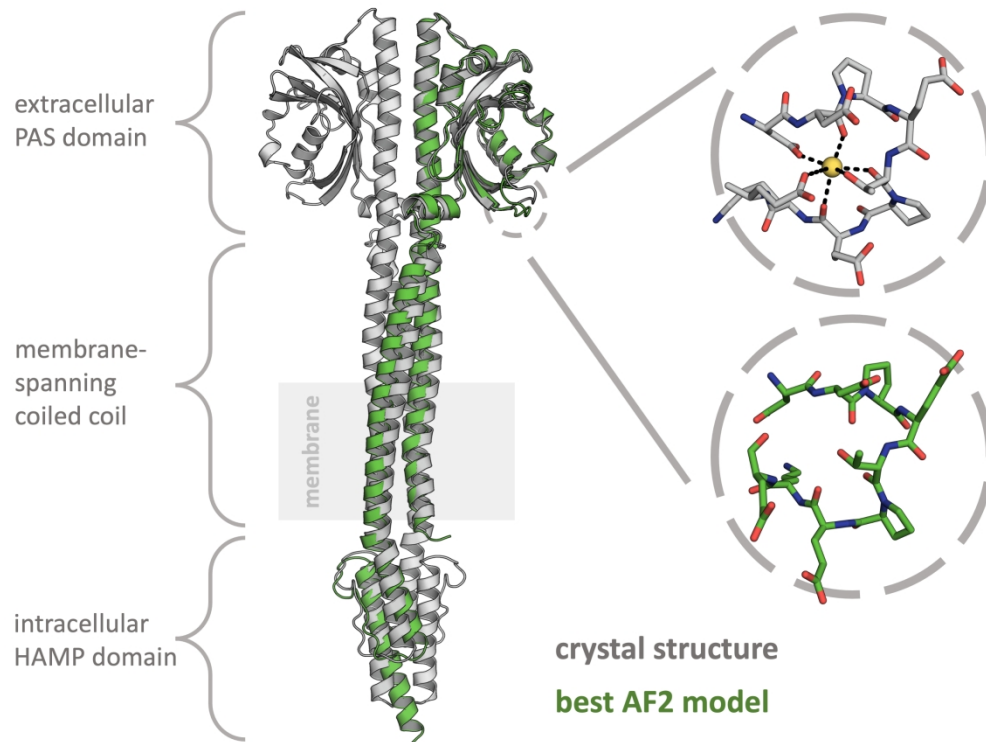


Figure 6. Polypeptide chain tracing errors that were corrected by examination of the AlphaFold2 (group 427) structure. (A) The incorrect model in the vicinity of two neighboring proline residues (Pro236 and Pro239) together with the associated difference electron density map with the coefficient $2F_o - F_c$ colored blue (left) and the model corrected based on the AlphaFold2 predicted structure with the associated $2F_o - F_c$ difference electron density map (right). The cis bond conformations are highlighted in green (B) The incorrect placement of Ile247 with the associated $2F_o - F_c$ difference electron density map colored blue and the $F_o - F_c$ difference electron density map colored green (left). Correcting the positions of Pro236 and Pro239 allowed placement of Tyr249 instead of Ile247 and eliminated the residual $F_o - F_c$ difference electron density (right).

177x177mm (300 x 300 DPI)



31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Figure 7. The crystal structure of dimeric Af1503 (grey) is shown in a superposition with the best AlphaFold2 model (green, monomer). The only noteworthy difference between the prediction and the crystal structure is found in a loop in the PAS domain, which was found to coordinate an ion in the crystal structure.

253x190mm (300 x 300 DPI)

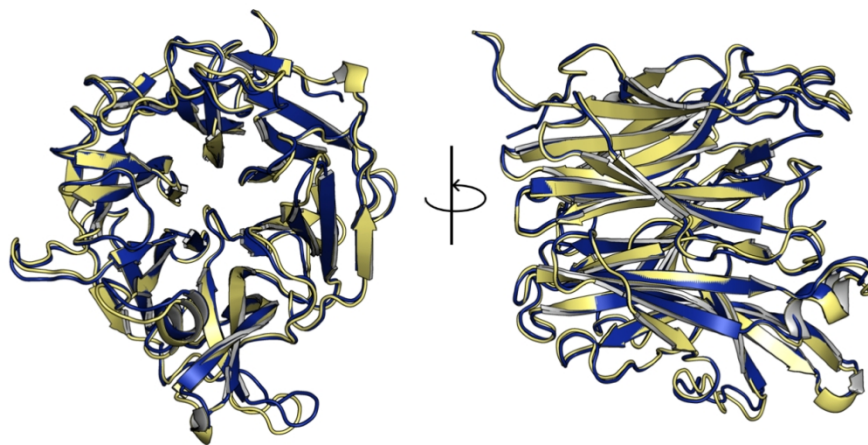


Figure 8. Superposition of Sia24 predictive and crystallographic models. The structure of Sia24 (dark blue) was solved with initial phase determination by molecular replacement using a model generated by AlphaFold2 (yellow).

216x101mm (169 x 169 DPI)