

AlphaFold Protein Structure Database for Sequence-Independent Molecular Replacement

Lawrence Chai^{a,‡}, Ping Zhu^{a,‡}, Jin Chai^a, Changxu Pang^a, Babak Andi^b, Sean McSweeney^b, John Shanklin^a, Qun Liu^{a,b*}

^a Biology Department, Brookhaven National Laboratory, Upton, NY 11973, USA

^b Photon Science, NSLS-II, Brookhaven National Laboratory, Upton, NY 11973, USA

[‡]These authors contributed equally to this study

*Corresponding author Email: qunliu@bnl.gov

Keywords: AlphaFold, molecular replacement, crystallization contaminants, structure determination, YncE, YadF.

Abstract

Crystallographic phasing recovers the phase information that is lost during a diffraction experiment. Molecular replacement is a dominant phasing method for the crystal structures in the protein data bank. In one form it uses a protein sequence to search a structure database for finding suitable templates for phasing. However, such sequence information is not always available such as when proteins are crystallized with unknown binding partner proteins or when the crystal is that of a contaminant. The recent development of AlphaFold has resulted in the availability of predicted protein structures for all proteins from twenty species. In this work, we tested whether AlphaFold-predicted *E. coli* protein structures were accurate enough for sequence-independent phasing of diffraction data from two crystallization contaminants for which we had not identified the protein. Using each of more than 4000 predicted structures as a search model, robust molecular replacement solutions were obtained which allowed the identification and structure determination of both structures, YncE and YadF. Our results advocate a general utility of AlphaFold-predicted structure database with respect to crystallographic phasing.

1. Introduction

Crystallographic phasing requires the retrieval of phase information that is lost during diffraction experiments. When there are no homology models, such phase information is recovered experimentally using isomorphous replacement preferably with their anomalous signals (Hendrickson, 2014, Liu & Hendrickson, 2015). With the accumulation of experimentally determined structures, molecular replacement (Rossmann, 1990, Evans & McCoy, 2008) is becoming a routine method for crystallographic phasing. For example, 71% of deposited crystal structures in the PDB database (www.pdb.org), were determined using molecular replacement.

Molecular replacement exploits the similarity between known structures and the structure to be determined. Programs such as MOLREP (Vagin & Teplyakov, 1997), PHASER (McCoy *et al.*, 2007), and AMoRe (Navaza, 2001) have been developed. When protein sequence information is known, molecular replacement pipelines may be used to automate the process as implemented in MrBUMP (Keegan *et al.*, 2018), BALBES (Long *et al.*, 2008), and MRage (Bunkoczi *et al.*, 2013). Using an *ab initio* modelling software such as ROSETTA (Das & Baker, 2008), predicted structures may also be used for molecular replacement as implemented in AMPLE (Bibby *et al.*, 2012). However, there are cases in which the sequence information is unknown. Examples include the crystallization of contaminant proteins or unknown protein-binding partner proteins (Hungler *et al.*, 2016). Under such scenarios, MoRDa (a non-redundant, annotated PDB database) (Vagin & Lebedev, 2015), ContaMiner/ContaBase (a collection of previously reported contaminant protein structures) (Hungler *et al.*, 2016), and a SIMBAD pipeline (Sequence-Independent Molecular Replacement Based on Available Databases) (Simpkin *et al.*, 2018) may be used for sequence-independent molecular replacement using database searching approaches. Among these tools, SIMBAD searches contaminant and MoRDa databases for a protein sequence-independent molecular-replacement (Simpkin *et al.*, 2018).

Machine learning has been extensively used for protein structure predictions with the recent development of the revolutionary attention-based AlphaFold (Bouatta *et al.*, 2021, Jumper *et al.*, 2021) and RoseTTAFold algorithms (Baek *et al.*, 2021). Both methods have enabled accurate prediction of protein structures approaching the fidelity of their crystal structures. In collaboration with an European Molecular Biology Laboratory (EMBL) team, AlphaFold released more than 350,000 predicted protein structures for twenty species including humans, and the predominant model systems including yeast, *Arabidopsis*, and *E. coli* (<https://alphafold.ebi.ac.uk>) (Tunyasuvunakool *et al.*, 2021). The predicted structures cover all the coded proteins within each species. The AlphaFold-predicted structures may serve as a valuable new resource to support crystallographic phasing. It is therefore possible

to use the species-wide structural databases for a protein sequence-independent molecular replacement for phasing diffraction data. This database approach may be of particular use for phasing proteins crystallized unexpectedly, proteolysis products, and structures of significant conformational changes. When crystallization of a protein with an unexpected binding partner protein, the AlphaFold database could be also used to identify potential binding proteins without the need for using mass spectrometry or protein sequencing.

During structural studies using X-ray crystallography, many proteins are expressed in *E. coli* and purified using various affinity columns. Often, in addition to protein of interest, *E. coli* contaminant proteins may bind either to the affinity resin or the protein of interest, and may co-purified and used for crystallization. Although crystallization of a contaminant protein is relatively rare, more contaminant structures have been identified as reported in the ContaBase database (Hungler *et al.*, 2016). For new contaminant proteins it may take some effort to identify it through experimental phasing, mass spectrometry, protein sequencing, or using database searches. AlphaFold has generated a complete database of predicted structures for all folded protein sequences in *E. coli*. Here we sought to test whether this database supports crystallographic phasing in the absence of protein sequence information.

In our crystallization work on two plant proteins that were over-expressed in *E. coli*, we unexpectedly crystallized two contaminants and collected diffraction data to about 2.3-2.5 Å resolution. For one of them, we could not solve its structure using existing methods. In this work, we used the two contaminant data sets for sequence-independent molecular replacement. Using a relatively straightforward workflow, we show that AlphaFold-predicted structures can be used to phase both structures without any protein sequence information. Our work highlights the broad utility of the AlphaFold-predicted structure database for crystallographic analysis.

2. Methods

2.1 Sample preparation for YncE/P76116

E. coli contaminant protein YncE was co-purified while we worked on the expression of a plant $\Delta 6$ desaturase in BL21-Gold (DE3) cells (Novagen). The desaturase protein was over-expressed at 30 °C for 4 hours by addition of 0.2 mM IPTG to the cell culture with an A_{600} of 0.6. Harvested cells were re-suspended in resuspension buffer (30 mM MES, 33 mM HEPES, 33mM NaOAc pH 7.5) supplemented with 2 mM $MgCl_2$ and 0.1 mg/ml DNase. The cells were lysed using a French Press and cell debris were removed by centrifugation at 25,000 x g for 30 min at 4°C. The clarified extract was loaded onto a Poros 20 HS column (Perceptive Biosystems, Framingham MA), washed with five column volumes of resuspension buffer, and eluted with a linear gradient of 0-1.2 M NaCl in the resuspension

buffer. Desaturase fractions were pooled and concentrated, subjected to a size-exclusion HPLC column (TSKgel G3000SW column, Tosoh Bioscience), and eluted with 20 mM HEPES, pH 7.0, and 100 mM NaCl. The desaturase fractions were pooled and concentrated to 15 mg/ml.

Crystals were grown using the hanging drop vapor diffusion method consisting of 0.6 μ l protein mixed with an equal volume of reservoir solution containing 0.2 M Li_2SO_4 , 0.1 M MES, pH 6.0, and 20% PEG 4000. Plate-shaped crystals were flash-frozen with liquid nitrogen. Cryo-protectant was not added prior to freezing.

2.2 Sample preparation for YadF/P61517

E. coli. contaminant protein YadF was co-purified with the production of Arabidopsis Metacaspase 4 (AtMC4) in BL21 pLysS (DE3) cells (Novagen). Cells were lysed using a homogenizer and the soluble fraction of AtMC4 was collected for a three-step purification by nickel-nitrilotriacetic acid (Ni-NTA) affinity chromatography (HisTrap FF column, GE Healthcare, Inc.), ion exchange chromatography (HiTrap Q HP column, GE Healthcare, Inc.), and gel filtration (Superdex 200 10/300 GL column, GE Healthcare, Inc.). Purified AtMC4 was then mixed and incubated with the excess molar amount of the inhibitor PPACK (Santa Cruz Biotechnology, Inc.). This mixture was further purified by gel filtration and the inhibitor-bound complex was concentrated to 8-10 mg/ml for crystallization.

Crystals were grown using the hanging drop vapor diffusion method. One μ l of inhibitor-bound AtMC4 was mixed with an equal volume of precipitant that contains 100 mM sodium cacodylate, pH 6.8, and 1.8 M ammonium sulfate. For cryo-crystallography, crystals were transferred into the precipitant supplemented with 10% glycerol and were flash-cooled into liquid nitrogen for cryogenic data collection.

2.3 Diffraction data collection and reduction

Diffraction data were collected at the NSLS-II beamline FMX (17ID-2) at 100 K (Fuchs *et al.*, 2016, Schneider *et al.*, 2021). The beamline is equipped with an Eiger 16M detector. For YncE, we collected data at an X-ray wavelength of 0.979 Å. A total of 1800 frames were collected from a single YncE crystal with a rotation angle of 0.2°. For YadF, we collected data at an X-ray wavelength of 1.891 Å. A total of ~1500 frames were collected from four YadF crystals with a rotation angle of 0.3°.

Single-crystal data sets were indexed and integrated independently using DIALS (Waterman *et al.*, 2016) and then scaled and merged using CCP4 programs POINTLESS and AIMLESS (Evans *et al.*, 2011, Evans & Murshudov, 2013) with the outlier rejection as implemented in PyMDA (Guo *et al.*, 2018, Takemaru *et al.*, 2020). For the YncE data, we rejected 700 radiation-damaged frames. For the YadF data, we rejected 948 radiation-

damaged frames using a decay value of 1.0 as defined by $\text{frame_cutoff} = [\text{Min}(\text{SmRmerge}) \times (1 + \text{decay})]$, where $\text{Min}(\text{SmRmerge})$ is the lowest SmRmerge (reported in AIMLESS log file) within a single-crystal data set; and decay is a rejection ratio (Takemaru *et al.*, 2020). The data collection and data processing statistics for the two data sets are shown in **Table 1**.

2.4 AlphaFold structures for database-driven molecular replacement

Fig. 1 shows the workflow of using AlphaFold-predicted *E. coli* structure database for sequence-independent molecular replacement. From these twenty AlphaFold-predicted structure databases (<https://alphafold.ebi.ac.uk/download>), we downloaded all 4363 *E. coli* protein structures. Among these structures, we removed those with less than 50 residues from further use. Then, we set up a molecular replacement search using the remaining 4175 structures. For each structure, we performed molecular replacement in MOLREP (Vagin & Teplyakov, 1997) with both rotation and translation searches with a high-resolution data cut-off at 3.0 Å resolution. The structures displaying the highest rotation and translation peaks were used to narrow the molecular replacement search. For YncE, we removed 34 disordered residues from N-terminal and used MOLREP for multi-copy molecular replacement (Vagin & Teplyakov, 2000). For YadF, we tried molecular replacement in different space groups to find the one with the highest translation peak height.

2.5 Model building and structure refinement

Iterative model building and refinement were performed in COOT (Emsley *et al.*, 2010) and PHENIX.REFINE (Afonine *et al.*, 2012, Echols *et al.*, 2014), respectively. For the YncE data, Bijvoet pairs were averaged for structure refinement. For the YadF data, Bijvoet pairs were treated as two different reflections in structure refinement, and the resultant Fourier coefficients were used for calculation of Bijvoet-difference Fourier maps. We also used anomalous signals for a f'' refinement (Liu *et al.*, 2013) to find anomalous scattering elements in the YadF structure. For the f'' refinement, the occupancies for the potassium and zinc ions were first estimated so that their refined individual B factors are close to the average B factors from their interacting protein and water molecules. We then refined f'' in PHENIX.REFINE starting with f'' values of zero for sulfur, potassium and zinc. The stereochemistry of the refined structures was validated with PROCHECK (Laskowski *et al.*, 1993) and MolProbity (Chen *et al.*, 2010) for quality assurance. The refinement statistics for the two data sets are shown in **Table 1**.

3. Results

3.1. AlphaFold structures for phasing YncE

During our work on the purification and crystallization of a plant desaturase, we co-purified YncE under crystallization conditions of 0.2 M Li₂SO₄, 0.1 M MES, pH 6.0, and 20% PEG 4000. We collected diffraction data and processed the data to d_{\min} 2.5 Å in space group P2₁ with unit dimensions a=53.2 Å, b=147.3 Å, c=96.9 Å, and β =104.4 °. Although the expected sequence identity for the desaturase to its homologous structures in PDB is beyond 80%, we were unable to solve the structure using the PDBs as search models, suggesting that this crystallized protein could be a contaminant. We used CCP4 online servers to search for contaminants but did not get a clear solution. To identify the contaminant, we also tried to repeat the crystallization and used mass spectrometry to identify the contaminant. Unfortunately, we were unable to reproduce the exact crystals.

With only the diffraction data available, we hypothesized that the contaminant protein must originate from *E. coli*. With the release of the AlphaFold-predicted *E. coli* structures, we reasoned that the crystallized contaminant should be represented in the AlphaFold structure database. We proceeded with the procedure described in **Fig. 1** to search for a monomer. All AlphaFold structures give their highest rotation and translation peaks beyond zero with a single structure, YncE (UNIPROT entry P76116), showed the highest RF/sig and TF/sig of 12.43 and 13.08, respectively (**Fig. 2a**).

Unit-cell content and self-rotation function analyses suggested the presence of multiple copies of YncE in the asymmetric unit (a.u.). We therefore performed focused molecular replacement searches for multiple copies using MOLREP. Visualization of the AlphaFold-predicted YncE structure indicated that it has a long N-terminal extension consisting of 34 poorly predicted/disordered residues. To assure that such a long extension would not affect the packing analysis in MOLREP, we removed the N-terminal 34 residues and used the truncated model for search of 2-5 monomers. We obtained the best results while searching for four monomers in a.u. and observed that both TF/sig and wRfac improved with an increasing number of monomers (**Fig. 2b**). With the four-monomer search, the final TF/sig and wRfac are 25.35 and 0.437, respectively, strongly indicating a correct solution for protein identification and structure determination.

The refined YncE structure has four molecules, each containing residues from 32 to 342 and forming a seven-bladed β -propeller structure (**Fig. 2c**). Except the N-terminal extension, the structure is very similar to the AlphaFold-predicted structure with an RMSD of 0.39 Å for 321 aligned C α atoms (**Fig. 2d**). However, we found that many side chains have different conformations, perhaps due to crystal contacts or disordered conformations.

In the UNIPROT entry for P76116/YncE, two PDBs (3VGZ and 3VH0) are reported, one crystallized in C222₁ lattice and the other crystallized in I4₁ lattice (Kagawa *et al.*, 2011). Our P2₁-form structure is a new contaminant structure. The P2₁-form structure has an RMSD of 0.44 Å with the C222₁-form structure and 0.37 Å with the I4₁-form structure, indicating that

all three structures are very similar although being crystallized in different space groups. **Table 2** summarizes the detailed crystallographic comparison of the YncE structure determined in three different lattices.

3.2 AlphaFold structures for phasing YadF

E. coli YadF is another contaminant protein that was co-purified with an Arabidopsis metacaspase 4 (AtMC4). AtMC4 is a cysteine protease and we have previously determined its structure in an apo form (Zhu *et al.*, 2020). To get a complex structure of AtMC4 with a protease inhibitor PPACK, we attempted to crystallize the complex for structural analysis. Crystals with dimensions of about 20-30 μm were obtained under the crystallization conditions of 0.1 M sodium cacodylate, pH 6.8 and 1.8 M $(\text{NH}_4)_2\text{SO}_4$. We collected diffraction data from four crystals at a relatively longer X-ray wavelength of 1.891 Å. The processed data at d_{min} 2.3 Å has a tetragonal lattice with unit-cell dimensions of $a=67.5$ Å and $c=85.3$ Å. However, we couldn't solve its structure using the AtMC4 structures of either the full-length or its truncations. Therefore, we suspected that this could be another *E. coli* contaminant and may be suitable for structure determination through searching the AlphaFold-predicted structure database.

Using the same workflow described above for YncE, we performed molecular replacement searches using MOLREP for each of the 4175 structures. **Fig. 3a** shows the histogram plot for RF/sig and TF/sig. Although there are four targets with highest translation peaks beyond 10 (UNIPROT entries P0CF69, P75971, P0CF68, and P61517), P61517/YadF is the only target with the highest rotation peak at 9.04, suggesting it is a possible solution for downstream model building and refinement. YadF has 220 residues, the unit-cell content analysis suggested a single molecule in a.u. with an estimated solvent content of 43%. The initial refinement in PHENIX.REFINE yield an R/free R of 0.30/0.39, suggesting larger structural differences relative to the AlphaFold-predicted structure. Therefore, we rebuilt the model using ARP/WARP (Langer *et al.*, 2008). ARP/WARP produced a nearly complete model of 208 residues in one chain with an R/free R of 0.194/0.252, indicating a correct identification and structure determination using the AlphaFold structure database.

The refined structure has 211 residues, and its structure is shown in **Fig. 3b**. The structure has an N-terminal α -helix domain and a C-terminal mixed $\alpha\beta$ domain. Compared with the AlphaFold-predicted structure, the RMSD is 1.18 Å for 206 aligned C_α atoms. Most structural differences are on the N-terminal helix and the loop connecting it to the $\alpha\beta$ domain (**Fig. 3c**).

YadF is a carbonic anhydrase whose activity is zinc dependent (Cronk *et al.*, 2001). We had collected data at an X-ray wavelength of 1.891 Å at which the theoretical anomalous signal f'' is 0.98 e. Therefore, we used an f'' refinement to characterize zinc anomalous signals (Liu *et al.*, 2013). With an estimated occupancy of 1.0 for the zinc site, the refined f'' is 0.94 e, clearly validating the specialization of the zinc site. Zinc is coordinated with two cysteine residues (Cys42 and Cys101), His98, and Asp44. **Fig. 3d** shows the Bijvoet difference Fourier densities for the active site. The Bijvoet densities cover zinc as well as two sulfur atoms. Surprisingly, next to the zinc/sulfur densities, we observed an extra electron density next to His98. To identify the type of anomalous scatterers associated with this density, we performed the f'' refinement with a candidate ion of Zn^{2+} , Ca^{2+} , K^+ , or Na^+ . Through the f'' refinements, the only reasonable fit for this anomalous scatterer is K^+ with an occupancy of 0.6 and a B-factor of 33.5 Å². However, we did not include K^+ either in protein purification or crystallization. Its exact origin and potential functional role will therefore be the subject of further investigation.

The AlphaFold-predicted structure does not contain any ions, neither Zn^{2+} nor K^+ . Structural superimposition of the AlphaFold structure with the ion-bound YadF structure indicates conformational changes of Asp44 (**Fig. 3d**). Interestingly, the same residue has been proposed to undergo conformational change so that substrate CO_2 can approach Zn^{2+} to form a CO-Zn^{2+} species (Cronk *et al.*, 2001). Thus, it is possible that the AlphaFold-predicted structure might resemble an intermediate state of YadF, at least for the active site structure.

Under UNIPROT entry P61517, there are four reported PDBs (1I6O, 1I6P, 1T75, and 2ESF) (Cronk *et al.*, 2001, Cronk *et al.*, 2006), all determined in tetragonal lattices but with different crystallization conditions. Our structure has a RMSD between 0.35 and 0.79 Å compared to these structures. **Table 3** summarizes detailed crystallographic comparison of YadF under different crystallization conditions.

4. Discussion

4.1 AlphaFold-predicted structure database

Crystallizing protein contaminants is a relatively common problem. In this work we demonstrate that AlphaFold-predicted *E. coli* structures can be useful for molecular replacement to identify unknown crystallized contaminant proteins and determine their structures. In our tests, we did not modify the predicted structures for the initial molecular replacement searches even though these predicted structures may contain unstructured extensions and poorly predicted regions such as we found with the N-terminal long extension in YncE.

For the two contaminant structures that we determined using AlphaFold-predicted structure database, YncE is a new contaminant. Although there are two crystal structures (PDB entries 3VGZ and 3VH0), we did not get a clear solution while trying database search approaches using the CCP4 online server. As a comparison, for YadF, in addition to using AlphaFold structure database, we can find a solution using its unit cell dimensions to search the PDB database, and PDBs 1I6P and 4ZNZ were identified. It turned out that PDB 4ZNZ was already reported as a crystallization contaminant (Niedzialkowska *et al.*, 2016) but crystallized in a different condition (**Table 3**). It is noted that the YadF structure in this work has a larger RMSD with the AlphaFold-predicted structure (1.2 Å) than with other crystal structures (0.35 – 0.79 Å). As shown in **Fig 3c**, largest structural differences are located at the N-terminal helix. In YadF crystal structures, this helix is stabilized by forming a dimer with its symmetry mate (Cronk *et al.*, 2001). As a contrast, AlphaFold-predicted structure is a monomer, and the N-terminal helix can thus be more flexible.

Phasing with an *E. coli* structure database has multiple advantages over the PDB database. First, the predicted structures contain only single-chain structures, which may be used directly for rotation search with no need for further processing, such as removing non-protein components or splitting a protein complex into individual components. Secondly, the predicted structure is based on the entire encoded protein sequences. Consequently, using such a database gives a higher chance to find a promising and suitable structure template for phasing. Although in this work we used *E. coli* structures for identification and determination of contaminant structures, AlphaFold has predicted 350,000 structures for proteins from 20 species (Tunyasuvunakool *et al.*, 2021); and those structure databases may be well suited for phasing contaminant structures from other expression hosts such as mammalian cells, yeast, Arabidopsis and etc. Thirdly, AlphaFold structures may be used to identify and phase unexpected proteolytic fragments or unexpected binding partner proteins.

Using a domain-structure database and modelled structure for phasing has been previously implemented in MoRDa and AMPLE, respectively (Bibby *et al.*, 2014, Vagin & Lebedev, 2015). However, due to the limited number of structural domains and the inaccuracy associated with the modelling, database-based phasing is not routine, and is normally used as a method of last resort after exhausting other phasing strategy options. As AlphaFold-predicted structures have an improved accuracy relative to experimental structures, molecular replacement using AlphaFold structures could have more routine applications even for *de novo* phasing for which there is no homologous structure. The AlphaFold algorithm uses an artificial intelligence model that was extensively trained with available PDBs and sequence databases (Jumper *et al.*, 2021). Hence the AlphaFold-predicted structures could be biased toward known structures. Accordingly, novel protein structures with novel folds are needed to improve the prediction accuracy of AlphaFold. Based on our findings, we speculate that an increasing number of crystal structures will be phased using AlphaFold-predicted structure workflows.

4.2 Combining AlphaFold phasing with anomalous signals

Perhaps due to the existence of prior crystal structures for both YncE and YadF, AlphaFold-predicted structures are quite accurate, with a RMSD of 0.30 Å and 1.18 Å, to their refined structures (**Fig. 2d, 3c**). When there are only remote or no homologous structures, AlphaFold-predicted structures may be not sufficient for phasing solely through molecular replacement. We propose that molecular replacement with anomalous signals, e.g. MR-SAD (Panjikar *et al.*, 2009), might be a highly productive option.

For YadF, we collected long-wavelength data at 1.891 Å which allowed the characterization of anomalous scatters zinc, potassium, and sulfur atoms within the structure. To see whether anomalous signals will enhance AlphaFold-based crystallographic phasing, we tested MR-SAD (Panjikar *et al.*, 2009) using the PHASER_EP pipeline (Mccoy *et al.*, 2007). With the initial phases from the AlphaFold structure, PHASER_EP identified seven anomalous scatterers with a figure-of-merit of 0.467. The MR-SAD map is of high quality; the pipeline can build 201 residues in eight fragments, with the longest fragment representing 71 residues. Subsequently, ARP/wARP can build the exact model as starting from the AlphaFold structure without using anomalous signals. For YadF, anomalous signals did not help much because ARP/wARP overcame the model errors (for example the N-terminal helix, **Fig. 3c**) through automated model building. For cases where the model is not accurate enough or the diffraction data are not of high enough resolution, MR-SAD may help to solve structures that are otherwise very challenging or even currently considered unsolvable. Most proteins contain intrinsic sulfur atoms that are native anomalous scatters of long-wavelength X-rays. So, to optimize the use of AlphaFold-predicted structures for

phasing a *de novo* structure, it might be advantageous to collect long-wavelength native-SAD data preferably using a helium flight path if available. Then these anomalous signals from sulfur atoms can be used for AlphaFold-based phasing using MR-SAD.

5. Concluding remarks

Using the AlphaFold-predicted *E. coli* structure database, we identified and determined structures for two crystallization contaminants without protein sequence information. The molecular replacement solutions and the structural comparison of refined structures with those AlphaFold-predicted structures suggest that the predicted structures are of high accuracy for crystallographic phasing and will likely be integrated into other structure determination pipelines.

Acknowledgements

This work was supported in part by Brookhaven National Laboratory LDRD 22-008 and NIH grant GM107462. Q. L. was supported by the U.S. Department of Energy, Office of Science, Office of Biological and Environmental Research, as part of the Quantitative Plant Science Initiative at BNL. J.C. and J.S. were supported by Division of Chemical Sciences, Geosciences, and Biosciences, Office of Basic Energy Sciences, United States Department of Energy Grant DOE KC0304000. The FMX beamline is part of the Center for BioMolecular Structure (CBMS) which is primarily supported by the National Institutes of Health, National Institute of General Medical Sciences (NIGMS) through a Center Core P30 Grant (P30GM133893), and by the DOE Office of Biological and Environmental Research (KP1607011). NSLS-II is supported in part by the U.S. Department of Energy, Office of Science, Office of Basic Energy Sciences Program under contract number DE-SC0012704 (KC0401040).

Data availability

Atomic coordinates and structure factor files have been deposited in the RCSB Protein Data Bank (PDB) under the accession codes XXXX and XXXX.

References

- Afonine, P. V., Grosse-Kunstleve, R. W., Echols, N., Headd, J. J., Moriarty, N. W., Mustyakimov, M., Terwilliger, T. C., Urzhumtsev, A., Zwart, P. H. & Adams, P. D. (2012). *Acta Cryst. D* **68**, 352-367.
- Baek, M., DiMaio, F., Anishchenko, I., Dauparas, J., Ovchinnikov, S., Lee, G. R., Wang, J., Cong, Q., Kinch, L. N. & Schaeffer, R. D. (2021). *Science* **871**-876.
- Bibby, J., Keegan, R. M., Mayans, O., Winn, M. D. & Rigden, D. J. (2012). *Acta Crystallogr D* **68**, 1622-1631.
- Bibby, J., Keegan, R. M., Mayans, O., Winn, M. D. & Rigden, D. J. (2014). *Acta Crystallogr D* **70**, 1174-1174.
- Bouatta, N., Sorger, P. & AlQuraishi, M. (2021). *Acta Crystallographica Section D: Structural Biology* **77**, 982-991.
- Bunkoczi, G., Echols, N., McCoy, A. J., Oeffner, R. D., Adams, P. D. & Read, R. J. (2013). *Acta Crystallographica Section D-Structural Biology* **69**, 2276-2286.
- Chen, V. B., Arendall, W. B., Headd, J. J., Keedy, D. A., Immormino, R. M., Kapral, G. J., Murray, L. W., Richardson, J. S. & Richardson, D. C. (2010). *Acta Cryst. D* **66**, 12-21.
- Cronk, J. D., Endrizzi, J. A., Cronk, M. R., O'Neill, J. W. & Zhang, K. Y. (2001). *Protein Science* **10**, 911-922.
- Cronk, J. D., Rowlett, R. S., Zhang, K. Y., Tu, C., Endrizzi, J. A., Lee, J., Gareiss, P. C. & Preiss, J. R. (2006). *Biochemistry* **45**, 4351-4361.
- Das, R. & Baker, D. (2008). *Annu Rev Biochem* **77**, 363-382.
- Echols, N., Morshed, N., Afonine, P. V., McCoy, A. J., Miller, M. D., Read, R. J., Richardson, J. S., Terwilliger, T. C. & Adams, P. D. (2014). *Acta Cryst. D* **70**, 1104-1114.
- Emsley, P., Lohkamp, B., Scott, W. G. & Cowtan, K. (2010). *Acta Crystallogr D* **66**, 486-501.
- Evans, G., Axford, D., Waterman, D. & Owen, R. L. (2011). *Crystallogr. Rev.* **17**, 105-142.
- Evans, P. & McCoy, A. (2008). *Acta Crystallographica Section D-Structural Biology* **64**, 1-10.
- Evans, P. R. & Murshudov, G. N. (2013). *Acta Cryst. D* **69**, 1204-1214.
- Fuchs, M. R., Bhogadi, D. K., Jakoncic, J., Myers, S., Sweet, R. M., Berman, L. E., Skinner, J., Idir, M., Chubar, O., McSweeney, S. & Schneider, D. K. (2016). *AIP Conf. Proc.* **1741**, 0300026.
- Guo, G. R., Fuchs, M. R., Shi, W. X., Skinner, J., Rerman, E., Ogata, C. M., Hendrickson, W. A., McSweeney, S. & Liu, Q. (2018). *Iucrj* **5**, 238-246.
- Hendrickson, W. A. (2014). *Q Rev Biophys* **47**, 49-93.
- Hungler, A., Momin, A., Diederichs, K. & Arold, S. T. (2016). *J Appl Crystallogr* **49**, 2252-2258.
- Jumper, J., Evans, R., Pritzel, A., Green, T., Figurnov, M., Ronneberger, O., Tunyasuvunakool, K., Bates, R., Žídek, A. & Potapenko, A. (2021). *Nature* <https://doi.org/10.1038/s41586-41021-03819-4>.
- Kagawa, W., Sagawa, T., Niki, H. & Kurumizaka, H. (2011). *Acta Crystallographica Section D-Structural Biology* **67**, 1045-1053.
- Keegan, R. M., McNicholas, S. J., Thomas, J. M. H., Simpkin, A. J., Simkovic, F., Uski, V., Ballard, C. C., Winn, M. D., Wilson, K. S. & Rigden, D. J. (2018). *Acta Crystallographica Section D-Structural Biology* **74**, 167-182.
- Langer, G., Cohen, S. X., Lamzin, V. S. & Perrakis, A. (2008). *Nat Protoc* **3**, 1171-1179.
- Laskowski, R. A., Macarthur, M. W., Moss, D. S. & Thornton, J. M. (1993). *J. Appl. Cryst.* **26**, 283-291.
- Liu, Q. & Hendrickson, W. A. (2015). *Curr Opin Struc Biol* **34**, 99-107.
- Liu, Q., Liu, Q. L. & Hendrickson, W. A. (2013). *Acta Crystallogr D* **69**, 1314-1332.
- Long, F., Vagin, A. A., Young, P. & Murshudov, G. N. (2008). *Acta Crystallographica Section D-Structural Biology* **64**, 125-132.
- Mccoy, A. J., Grosse-Kunstleve, R. W., Adams, P. D., Winn, M. D., Storoni, L. C. & Read, R. J. (2007). *J Appl Crystallogr* **40**, 658-674.
- Navaza, J. (2001). *Acta Crystallogr D* **57**, 1367-1372.
- Niedzialkowska, E., Gasiorowska, O., Handing, K. B., Majorek, K. A., Porebski, P. J., Shabalin, I. G., Zasadzinska, E., Cymborowski, M. & Minor, W. (2016). *Protein Science* **25**, 720-733.
- Panjikar, S., Parthasarathy, V., Lamzin, V. S., Weiss, M. S. & Tucker, P. A. (2009). *Acta Crystallographica Section D-Structural Biology* **65**, 1089-1097.
- Rossmann, M. G. (1990). *Acta Crystallogr A* **46**, 73-82.
- Schneider, D. K., Shi, W., Andi, B., Jakoncic, J., Gao, Y., Bhogadi, D. K., Myers, S. F., Martins, B., Skinner, J. M. & Aishima, J. (2021). *J Synchrotron Radiat* **28**, 650-665.

- Simpkin, A. J., Simkovic, F., Thomas, J. M. H., Savko, M., Lebedev, A., Uski, V., Ballard, C., Wojdyr, M., Wu, R., Sanishvili, R., Xu, Y. B., Lisa, M. N., Buschiazzi, A., Shepard, W., Rigden, D. J. & Keegan, R. M. (2018). *Acta Crystallographica Section D-Structural Biology* **74**, 595-605.
- Takemaru, L., Guo, G. R., Zhu, P., Hendrickson, W. A., McSweeney, S. & Liu, Q. (2020). *J Appl Crystallogr* **53**, 277-281.
- Tunyasuvunakool, K., Adler, J., Wu, Z., Green, T., Zielinski, M., Židek, A., Bridgland, A., Cowie, A., Meyer, C. & Laydon, A. (2021). *Nature* 1-9.
- Vagin, A. & Lebedev, A. (2015). *Acta Crystallogr A* **71**, S19-S19.
- Vagin, A. & Teplyakov, A. (1997). *J Appl Crystallogr* **30**, 1022-1025.
- Vagin, A. & Teplyakov, A. (2000). *Acta Crystallographica Section D-Structural Biology* **56**, 1622-1624.
- Waterman, D. G., Winter, G., Gildea, R. J., Parkhurst, J. M., Brewster, A. S., Sauter, N. K. & Evans, G. (2016). *Acta Cryst. D* **72**, 558-575.
- Zhu, P., Yu, X. H., Wang, C., Zhang, Q. F., Liu, W., McSweeney, S., Shanklin, J., Lam, E. & Liu, Q. (2020). *Nat Commun* **11**,

Figure Captions

Figure 1. Schematic workflow of sequence-independent crystallographic phasing using AlphaFold-predicted *E. coli* structures. A total number of 4363 AlphaFold-predicted structures were downloaded from the AlphaFold structure database. After filtering based on protein sequence length, 4175 structures were selected for molecular replacement using MOLREP. The output candidate solutions were sorted based on TR/sig and the AlphaFold structure with the highest TR/sig peak height was selected for focused molecular replacement and downstream model building and refinement.

Figure 2. Alphafold structure for phasing *E. coli*. YncE. (a) Histogram of rotation and translation peaks. (b) Progressive molecular replacement while searching for four molecules in a.u. (c) Refined YncE structure. (d) Comparison of the refined structure with the AlphaFold structure.

Figure 3. AlphaFold structure for phasing *E. coli*. YadF. (a) Histogram of rotation and translation peaks. (b) Refined YadF structure. (c) Comparisons with the AlphaFold structure. (d) Active-site structure. Residues interacting with the zinc site are shown as sticks. Bijvoet difference Fourier map for anomalous scatterers were shown as magenta isomeshes contoured at 3σ . As a comparison, AlphaFold-predicted structure is shown in gray.

Table 1. Data collection and refinement statistics

Data collection	YadF/P61517	YncE/P76116
Beamline	FMX (17-ID-2, NSLS-II)	FMX (17-ID-2, NSLS-II)
Wavelength (Å)	1.891	0.979
Space group	P4 ₂ 2 ₁ 2	P2 ₁
Cell dimensions a,b,c (Å) β (°)	67.52, 67.52, 85.25	53.17, 147.27, 96.90, 104.4
Solvent content (%)	43.0	51.8
Bragg spacings (Å)	36-2.3 (2.36-2.3)	50-2.5 (2.56-2.5)
Total reflections	222819	134117
Unique reflections ¹	9286 (665)	47818 (3604)
Completeness (%)	100.0 (100.0)	95.9 (97.3)
<I/σ(I)>	9.9 (2.2)	7.3 (2.1)
R _{merge}	0.258 (0.912)	0.087 (0.048)
Multiplicity	24.0 (18.8)	2.8 (2.8)
CC _{1/2} (%)	99.5 (81.2)	98.1 (97.3)
Refinement		
Resolution (Å)	2.3	2.5
No. reflections	16710	87600
R _{work} /R _{free}	0.203/0.241	0.236/0.256
No. atoms	1756	10140
Wilson B (Å ²)	31.0	30.7
Average (Å ²)	40.2	46.6
R.m.s deviations		
Bond length (Å)	0.002	0.002
Bond angle (°)	0.414	0.521
PDB code	XXXX	XXXX

¹Values in parentheses are for the highest resolution range.

Table 2. Comparison of YncE/P76116 structure with PDB structures listed under UNIPROT entry P76116.

	P76116/XXXX (this work)	3VGZ	3VH0
Space group	P2 ₁	C222 ₁	I4 ₁
Resolution (Å)	2.5	1.7	2.9
Number of chains	4	4	4
Cell dimensions			
a,b,c (Å)	a=53.2	a=119.2	a=171.2
α, β, γ (°)	b=147.3 c=96.9 β=104.4	b=139.3 c=173.7	c=177.2
RMSD vs P76116 (Å)	-	0.44	0.37
Crystallization conditions	0.2 M Li ₂ SO ₄ , 0.1 M MES, pH 6.0, 20% PEG4000	0.1 M sodium acetate, pH 4.4, 0.2 M (NH ₄) ₂ SO ₄ , 25% PEG 4000	0.1 M trisodium citrate pH 5.6, 2% tacsimate, pH 5.0, 16% PEG 3350

Table 3. Comparison of YadF/P61517 structures with PDB structures listed under UNIPROT entry P61517.

	P61517/XXXX (this work)	1I6O	1I6P	1T75	2ESF	4ZNZ
Space group	P4 ₂ 2 ₁ 2	P4 ₃ 22	P4 ₂ 2 ₁ 2	P4 ₃ 2 ₁ 2	P4 ₃ 22	P4 ₂ 2 ₁ 2
Resolution (Å)	2.3	2.2	2.0	2.5	2.25	2.7
Number of chains	1	2	1	4	2	1
Cell dimensions						
a,b,c (Å)	a=67.5	a=81.2	a=68.5	a=110.4	a=82.7	a=67.9
α, β, γ (°)	c=85.3	c=162.1	c=85.9	c=162.5	c=162.2	c=84.9
RMSD vs P61517 (Å)	-	0.77	0.35	0.78	0.79	0.45
Crystallization conditions	0.1 M sodium cacodylate, pH 6.8, 1.8 M (NH ₄) ₂ SO ₄	0.1 M MES pH 6.3, 1.6-1.8 M (NH ₄) ₂ SO ₄ , 4% PEG 400	0.1 M MES pH 6.3, 1.6-1.8 M (NH ₄) ₂ SO ₄	PEG 3000, pH 4.5	0.1 M MES, pH 6.5, 1.65 M (NH ₄) ₂ SO ₄ , 4% PEG 400	0.1 M Bis-Tris Propane, 60% Tacsimate, pH 7.0

Figure 1.

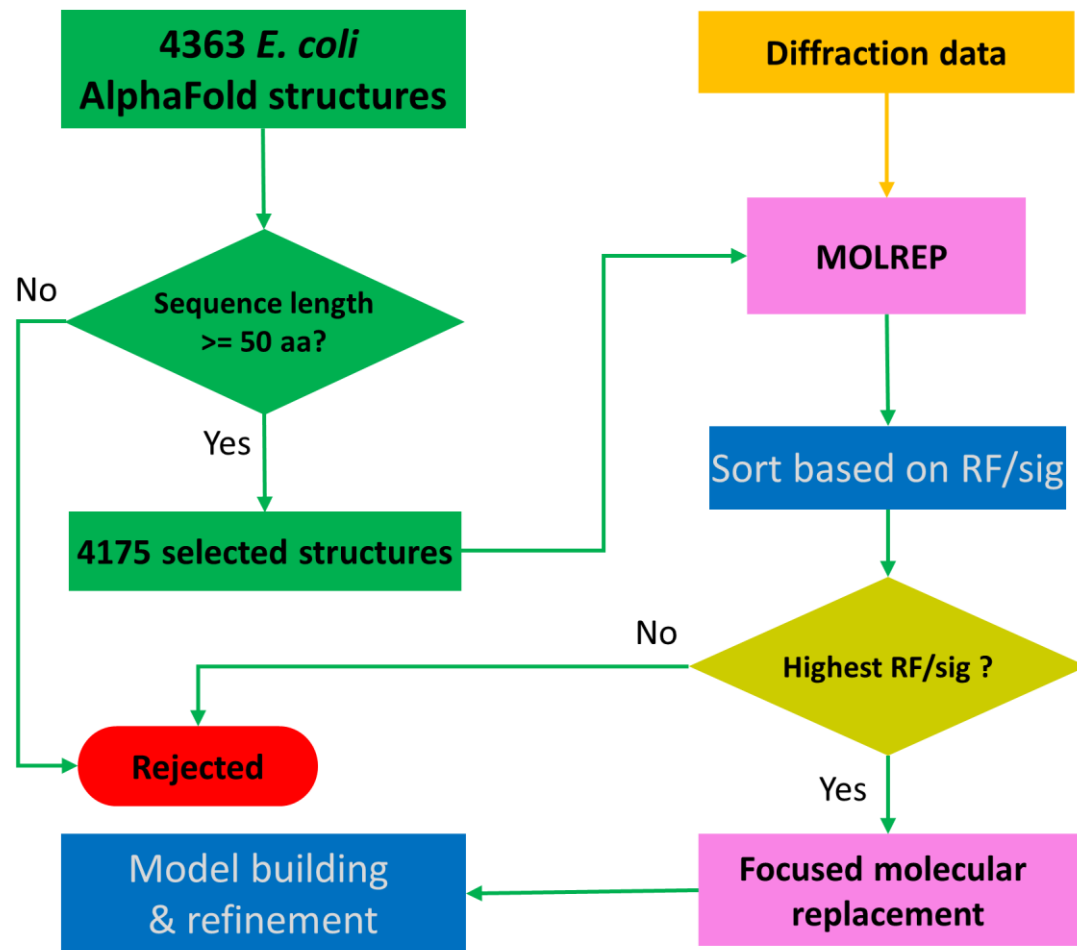


Figure 2.

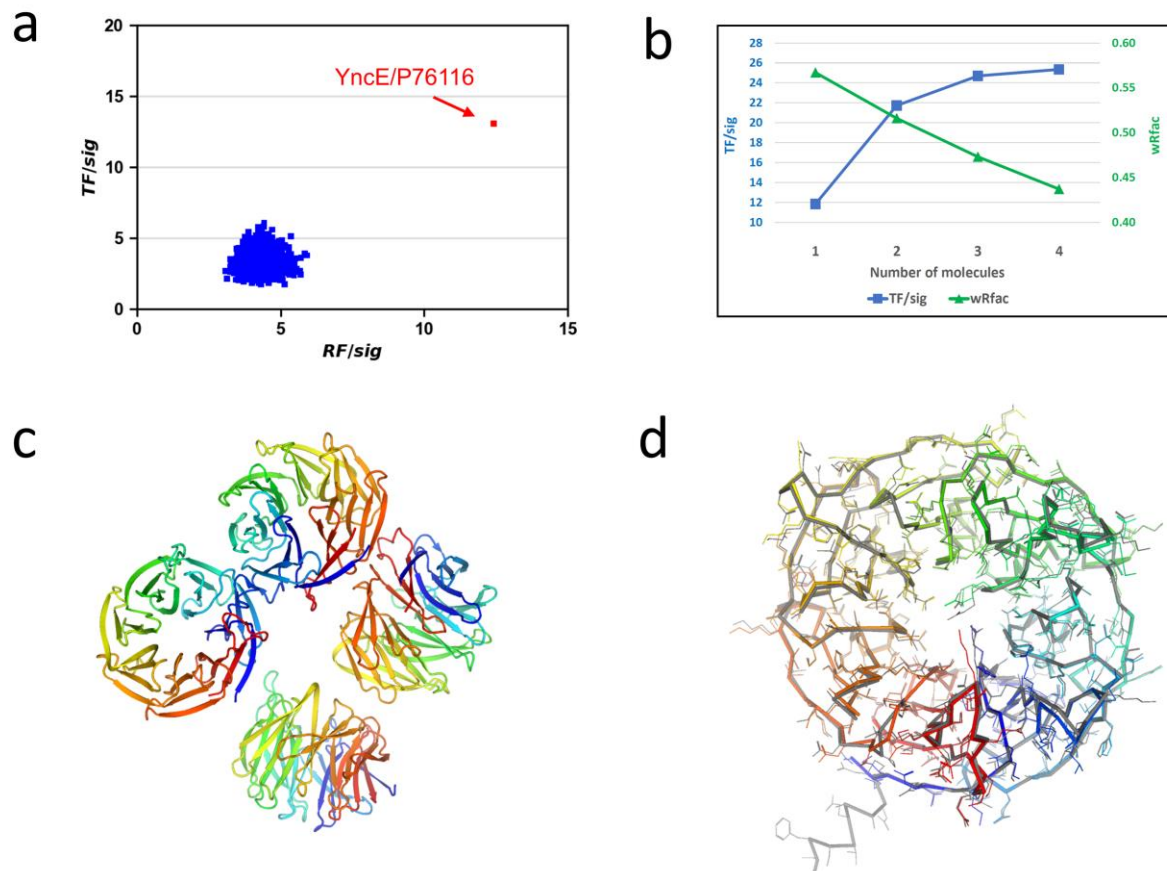


Figure 3.

