

CSCI-GA.3033-111: Protein Design HW2

Christine Wu (email: cw4459@nyu.edu)

Oct 27 2024

1 Explore the AlphaFold2, AlphaFold3, and Rosetta Fold colab notebooks

Solution: We have AlphaFold2, AlphaFold3, and Rosetta Fold colab notebooks present as followed: [AlphaFold2](#), [AlphaFold3](#), and [Rosetta Fold](#).

2 Predict some protein structures. Choose at least 3 proteins from [Lysozyme (P00698), Cytochrome C (P99999), Myoglobin (P02144), GFP (P42212), Beta-Lactamase (P62593), Thioredoxin (P0AA25), Ribonuclease A (P61823)], and use all three methods to generate structures for them (Sequence can be obtained on uniprot).

Solution: We first picked 4 proteins from the list, which are Lysozyme (P00698), Cytochrome C (P99999), Myoglobin (P02144), GFP (P42212). Then use [UniProt](#) to find their sequences. Lastly, use the three notebooks to generate the structures. We have the results as the following:

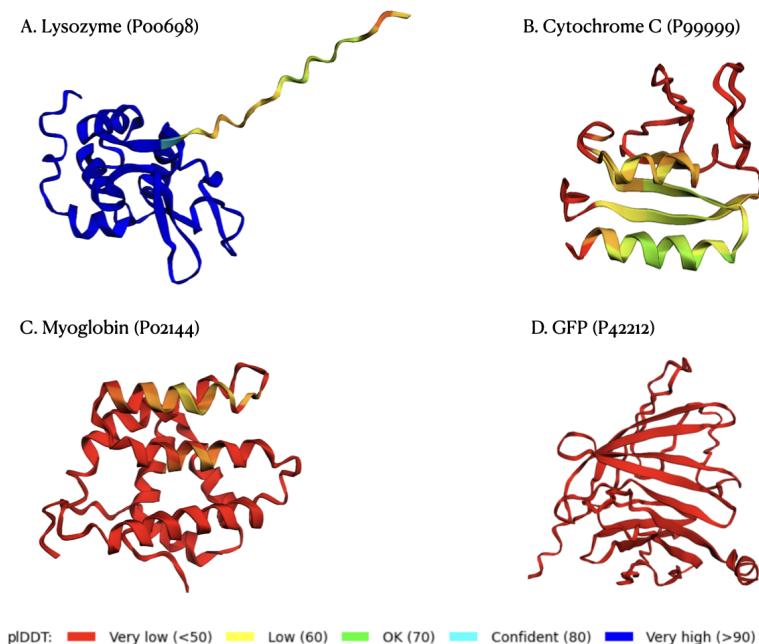


Figure 1: AlphaFold2 Predictions

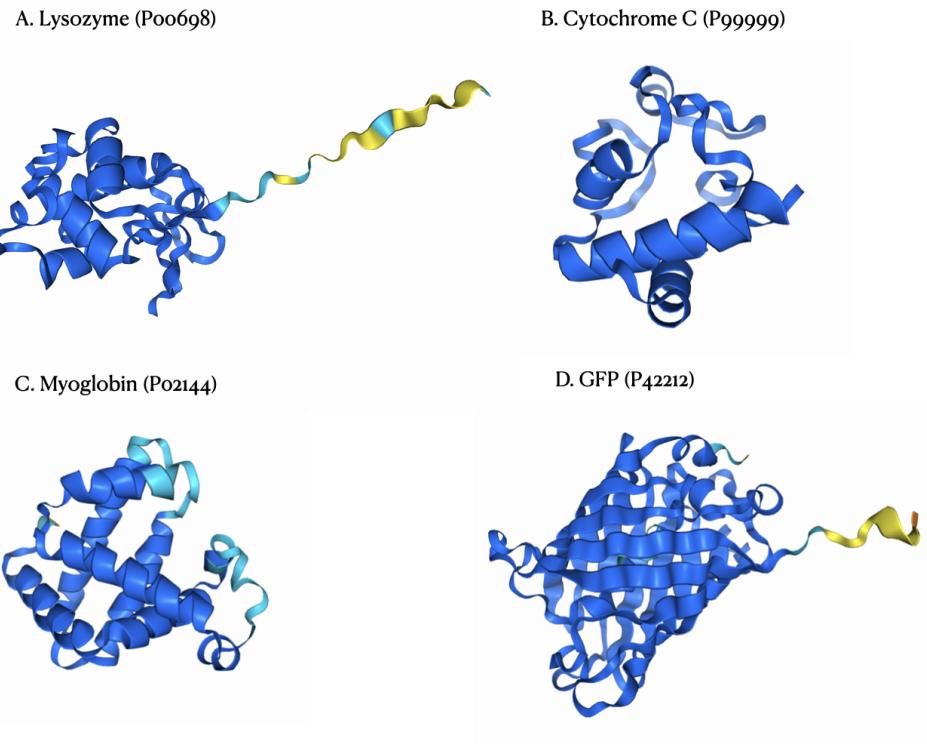


Figure 2: AlphaFold3 Predictions

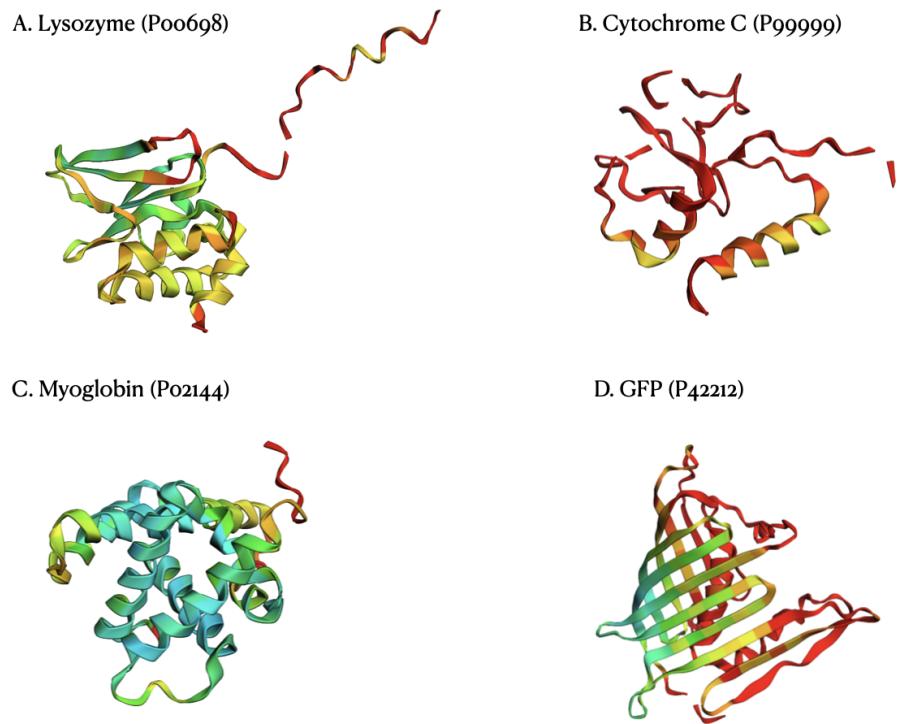


Figure 3: Rosetta Fold Predictions

The sequences we used are:

- **Lysozyme (P00698)** ([Lysozyme \(P00698\)](#)):

MRSLLILVLCFLPLAALGKVGRCELAAMKRHGLDNYRGYSLGNWVCAAKFESNFNTQATNRNTDGSTDYGILQINSRWW
CNDGRTPGSRNLNCIPCSALLSSDITASVNCACKIVSDNGMNAWVAWRNRCKGTDVQAWIRGCRL

- **Cytochrome C (P99999)** ([Cytochrome C \(P99999\)](#)):

MGDVEKGKKIFIMKCSQCHTVEKGKHKTGPNLHGLFGRKTGQAPGYSYTAANKNKGIIWGEDTLMEYLENPKKYIPGKTM
IFVGIIKKKEERADLIAYLKKATNE

- **Myoglobin (P02144)** ([Myoglobin \(P02144\)](#)):

MGLSDGEWQLVNVWGKVEADIPGHGQEVLIRLFKGHPETLEKFDFKFKHLKSEDEMKAEDLKKHGATVLTALGGILKKKG
HHEAEIKPLAQSHATKHKIPVKYLEFISECIIQVLQSKHPGDFGADAQGAMNKALELFRKDMASNYKELGFQG

- **GFP (P42212) Myoglobin (P02144) (GFP (P42212))**:

MSKGEELFTGVVPILVELDGDVNGHKFSVSGESEGDAKYGKLTLKFICTTGKLPVPWPTLVTTFSYGVQCFSRYPDHMKQH
DFFKSAMPEGYVQERTIFFKDDGNYKTRAEVKFEGLDTLVRRIELKGIDFKEDGNILGHKLEYNNNSHNVYIMADKQKNGIK
VNFKIRHNIEDGSVQLADHYQQNTPIGDPVLLPDNHYLSTQSALKDPNEKRDHMVLLEFVTAAGITHGMDELYK

3 Take note of the pLDDT charts created. What type of region did each respective model predict well? Why might the models have their respective strengths and weaknesses?

Solution: We have the following pLDDT charts presented. For alphafold3, we have our own python code to generated the plots with the downloaded data by using the `fold_2024_10_27_16_37_full_data_0`, `fold_2024_10_27_16_37_full_data_1`, `fold_2024_10_27_16_37_full_data_2`, `fold_2024_10_27_16_37_full_data_3`, `fold_2024_10_27_16_37_full_data_4`. As in the webpage, there were no pLDDT plots presented. Moreover, we also have the Predicated aligned error (PAE) diagrams. The following are the results:

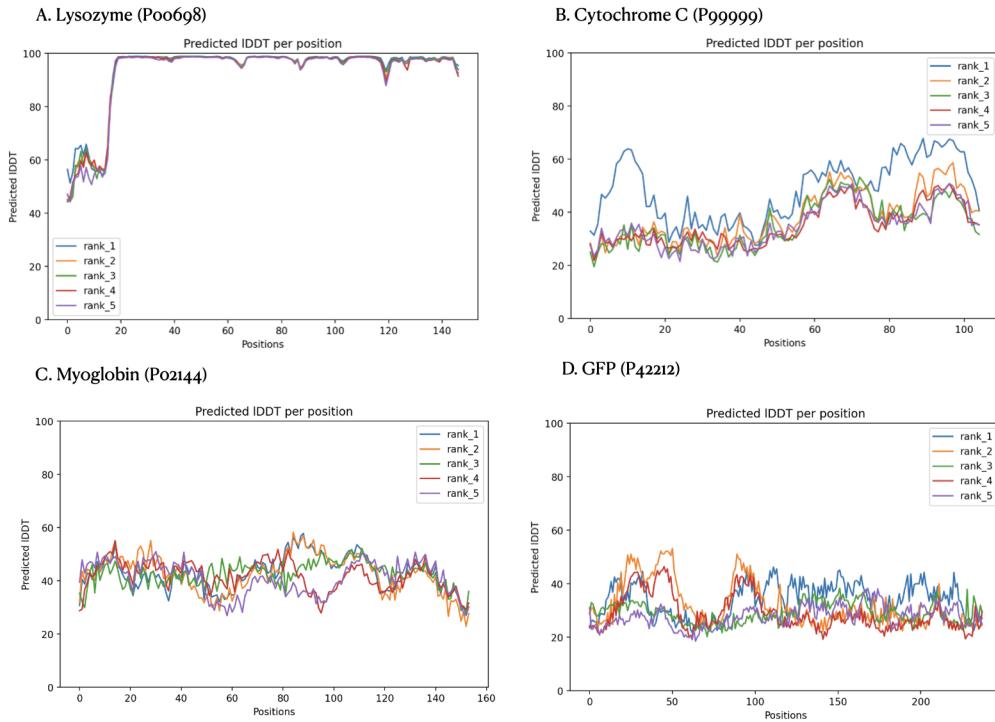


Figure 4: AlphaFold2 pLDDT charts

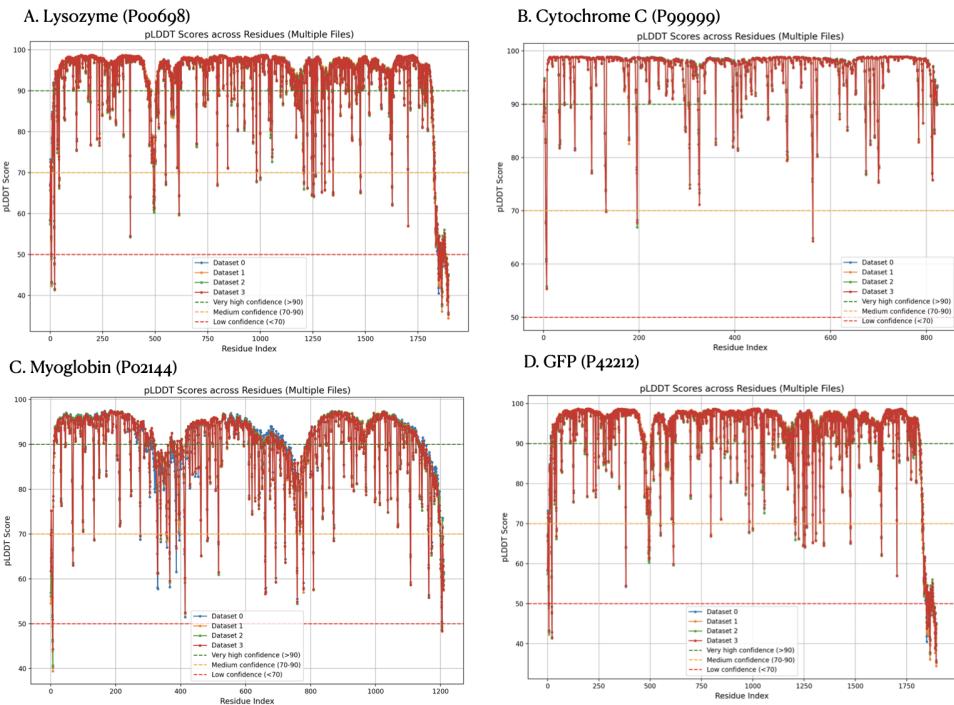


Figure 5: AlphaFold3 pLDDT charts

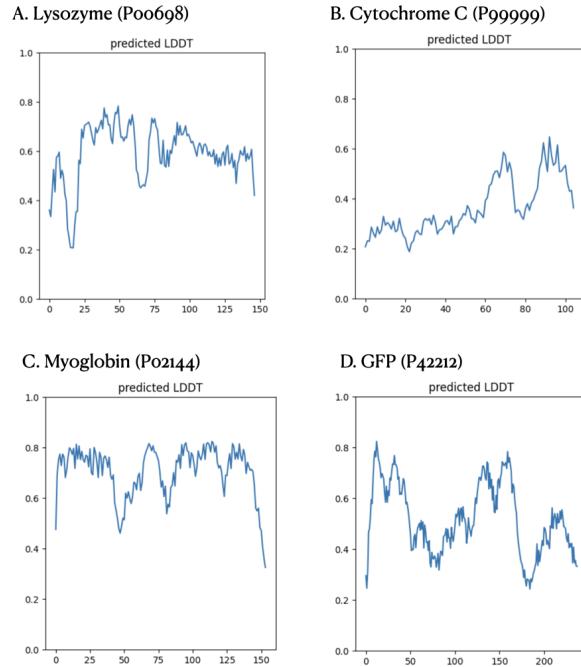


Figure 6: Rosetta Fold pLDDT charts

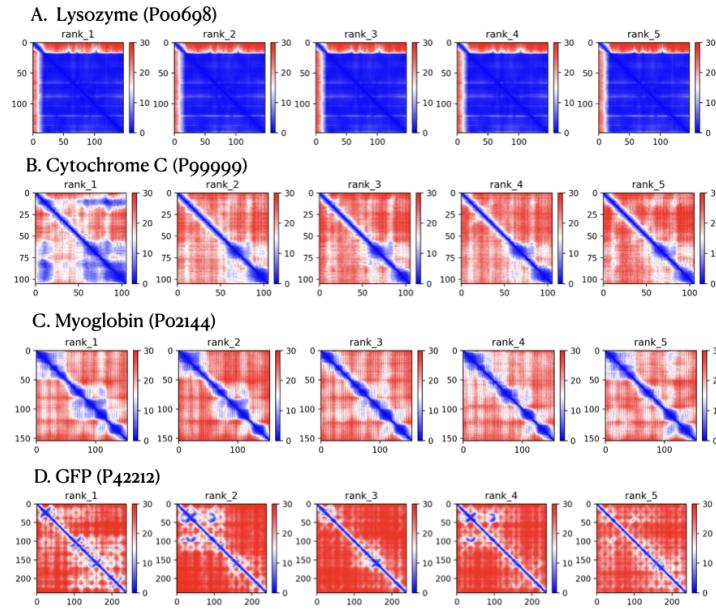


Figure 7: Alphafold2 PAE charts

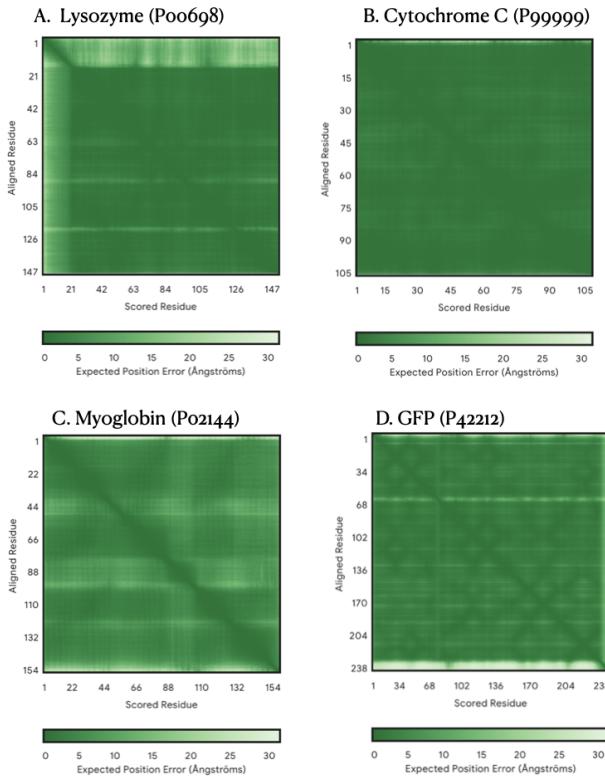


Figure 8: Alphafold3 PAE charts

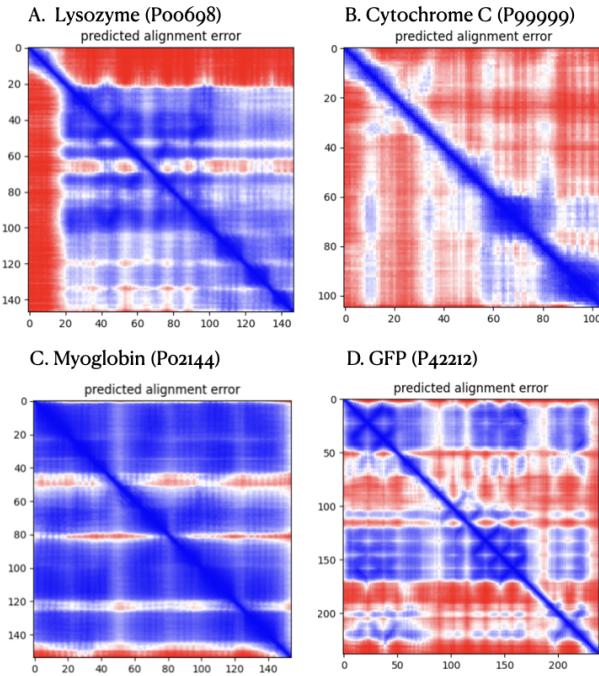


Figure 9: Rosetta Fold PAE charts

Summary: In order to tell which region each respective model predicted well, we can combine the answers from question 1 and the pLDDT and PAE plots from question 2. From question 1, we see the score visually by the bar indicating Regions with $p\text{LDDT} > 90$, Regions with $p\text{LDDT} 70\text{--}90$, and Regions with $p\text{LDDT} < 70$:

- **Alphafold2 predictions:** Proteins like Lysozyme (P00698) and Myoglobin (P02144) have high pLDDT confidence across their structured regions, indicated by blue in the models. These proteins have well-defined secondary structures, like alpha-helices and beta-sheets, which are predicted with high confidence. Cytochrome C (P99999) and GFP (P42212), however, display a combination of high- and low-confidence regions, with flexible loops and unstructured areas marked in yellow and red, indicating uncertainty in these regions. The pLDDT charts (Figure 4) confirm high confidence for structured regions in Lysozyme and Myoglobin, while Cytochrome C and GFP exhibit greater variability, particularly in loops or disordered regions. In the PAE diagrams (Figure 7), Alphafold2 shows low predicted alignment error (in blue) for secondary structures in Lysozyme and Myoglobin, while Cytochrome C and GFP display increased alignment errors (in red) in flexible loop regions, reflecting the model's difficulty in predicting these disordered areas accurately.
- **Alphafold3 predictions:** The Alphafold3 models demonstrate improved confidence across all proteins, with most regions, including well-structured secondary elements like alpha-helices and beta-sheets, colored blue, indicating pLDDT scores above 90. The pLDDT charts (Figure 5) also show consistently high confidence, with very few regions below pLDDT 70, suggesting greater predictive accuracy compared to Alphafold2. In the PAE diagrams (Figure 8), Alphafold3 maintains low alignment error across structured regions for Lysozyme and Myoglobin and shows improved alignment for Cytochrome C and GFP as well, although some loop regions still exhibit moderate alignment error. This suggests that Alphafold3 is better equipped to handle both structured secondary elements and flexible regions in complex proteins.
- **RosettaFold predictions:** RosettaFold predictions display more variation in confidence, particularly for Cytochrome C and GFP, where lower-confidence regions are indicated by red and yellow, highlighting difficulties in modeling complex and unstructured regions. While RosettaFold can reliably predict

well-defined secondary structures, like alpha-helices in Lysozyme and Myoglobin, it struggles more with regions that lack stable secondary structure, as seen in Cytochrome C and GFP. The pLDDT charts (Figure 6) show greater variability in confidence levels compared to the AlphaFold models, especially in flexible or disordered regions, where pLDDT values frequently fall below 0.5. The PAE diagrams (Figure 9) reveal the most alignment error in RosettaFold predictions, especially for proteins with flexible regions, such as Cytochrome C and GFP, indicating that RosettaFold struggles more with reliably predicting disordered regions compared to AlphaFold2 and AlphaFold3.

Overall, AlphaFold2 and AlphaFold3 excel in predicting well-structured proteins due to their transformer-based architectures and the use of multiple sequence alignments. However, they both struggle with disordered regions and protein complexes (the latter issue being mitigated in AlphaFold3). RosettaFold, while flexible and good for novel folds, may not have the same level of accuracy for complex structures and shows more variability in predictions, particularly for dynamic or disordered regions.

The differences in performance between AlphaFold2, AlphaFold3, and RosettaFold arise from the core methodologies behind each model. AlphaFold2 uses deep learning and evolutionary information from multiple sequence alignments (MSAs) to predict protein structures, making it particularly strong at modeling well-structured proteins. AlphaFold3 builds on AlphaFold2's architecture but adds improved all-atom modeling, which allows it to handle more complex tasks like protein-ligand and protein-DNA interactions, making it more versatile. On the other hand, RosettaFold relies heavily on co-evolutionary signals and flexible sampling techniques, which makes it effective for predicting novel folds, especially when evolutionary data is available, but it struggles more with dynamic and disordered regions, showing more variability in confidence. These differences reflect the models' unique balance between machine learning, evolutionary data, and structural sampling.

```

1 ## The code for alphafold3's pLDDT:
2 import json
3 import matplotlib.pyplot as plt
4 import numpy as np
5
6 def extract_pplddt_from_json(json_file):
7     # Load the JSON file
8     with open(json_file, 'r') as f:
9         data = json.load(f)
10
11    # Extract pLDDT scores from the 'atom_pplddts' field
12    pplddt_scores = data.get('atom_pplddts', [])
13
14    if not pplddt_scores:
15        print(f"No pLDDT scores found in {json_file}.")
16        return None
17
18    return pplddt_scores
19
20 def plot_pplddt(all_scores, labels):
21     # Plot the pLDDT scores for multiple data sets
22     plt.figure(figsize=(12, 8))
23
24     for scores, label in zip(all_scores, labels):
25         plt.plot(np.arange(len(scores)), scores, marker='o', linestyle='--', markersize=3, label=label)
26
27     plt.title('pLDDT Scores across Residues (Multiple Files)', fontsize=14)
28     plt.xlabel('Residue Index', fontsize=12)
29     plt.ylabel('pLDDT Score', fontsize=12)
30
31     # Add confidence threshold lines
32     plt.axhline(y=90, color='green', linestyle='--', label='Very high confidence (>90)')
33     plt.axhline(y=70, color='orange', linestyle='--', label='Medium confidence (70-90)')
34     plt.axhline(y=50, color='red', linestyle='--', label='Low confidence (<70)')
35
36     plt.legend()
37     plt.grid(True)
38     plt.show()
39

```

```

40 # List of file paths
41 json_files = [
42     'fold_2024_10_27_16_37/fold_2024_10_27_16_37_full_data_0.json',
43     'fold_2024_10_27_16_37/fold_2024_10_27_16_37_full_data_1.json',
44     'fold_2024_10_27_16_37/fold_2024_10_27_16_37_full_data_2.json',
45     'fold_2024_10_27_16_37/fold_2024_10_27_16_37_full_data_3.json'
46 ]
47
48 json_files = [
49     'fold_2024_10_27_16_37-2/fold_2024_10_27_16_37_full_data_0.json',
50     'fold_2024_10_27_16_37-2/fold_2024_10_27_16_37_full_data_1.json',
51     'fold_2024_10_27_16_37-2/fold_2024_10_27_16_37_full_data_2.json',
52     'fold_2024_10_27_16_37-2/fold_2024_10_27_16_37_full_data_3.json'
53 ]
54
55 json_files = [
56     'fold_2024_10_27_16_37-3/fold_2024_10_27_16_37_full_data_0.json',
57     'fold_2024_10_27_16_37-3/fold_2024_10_27_16_37_full_data_1.json',
58     'fold_2024_10_27_16_37-3/fold_2024_10_27_16_37_full_data_2.json',
59     'fold_2024_10_27_16_37-3/fold_2024_10_27_16_37_full_data_3.json'
60 ]
61
62 json_files = [
63     'fold_2024_10_27_16_37-4/fold_2024_10_27_16_37_full_data_0.json',
64     'fold_2024_10_27_16_37-4/fold_2024_10_27_16_37_full_data_1.json',
65     'fold_2024_10_27_16_37-4/fold_2024_10_27_16_37_full_data_2.json',
66     'fold_2024_10_27_16_37-4/fold_2024_10_27_16_37_full_data_3.json'
67 ]
68
69 all_plddt_scores = []
70 labels = ['Dataset 0', 'Dataset 1', 'Dataset 2', 'Dataset 3']
71
72 # Extract pLDDT scores for all files
73 for json_file in json_files:
74     plddt_scores = extract_plddt_from_json(json_file)
75     if plddt_scores:
76         all_plddt_scores.append(plddt_scores)
77
78 # Plot the results
79 if all_plddt_scores:
80     plot_plddt(all_plddt_scores, labels)

```

4 Find each protein's experimentally derived structure. Load the experimental and predicted structures into PyMol. Align them, and note any differences in structures.

Solution: We used the experimental structures, listed as: 2CDS for Lysozyme (P00698), 2N9I for Cytochrome C (P99999), 3RGK for Myoglobin (P02144), 2G16 for GFP (P42212). The results and the summary are the following:

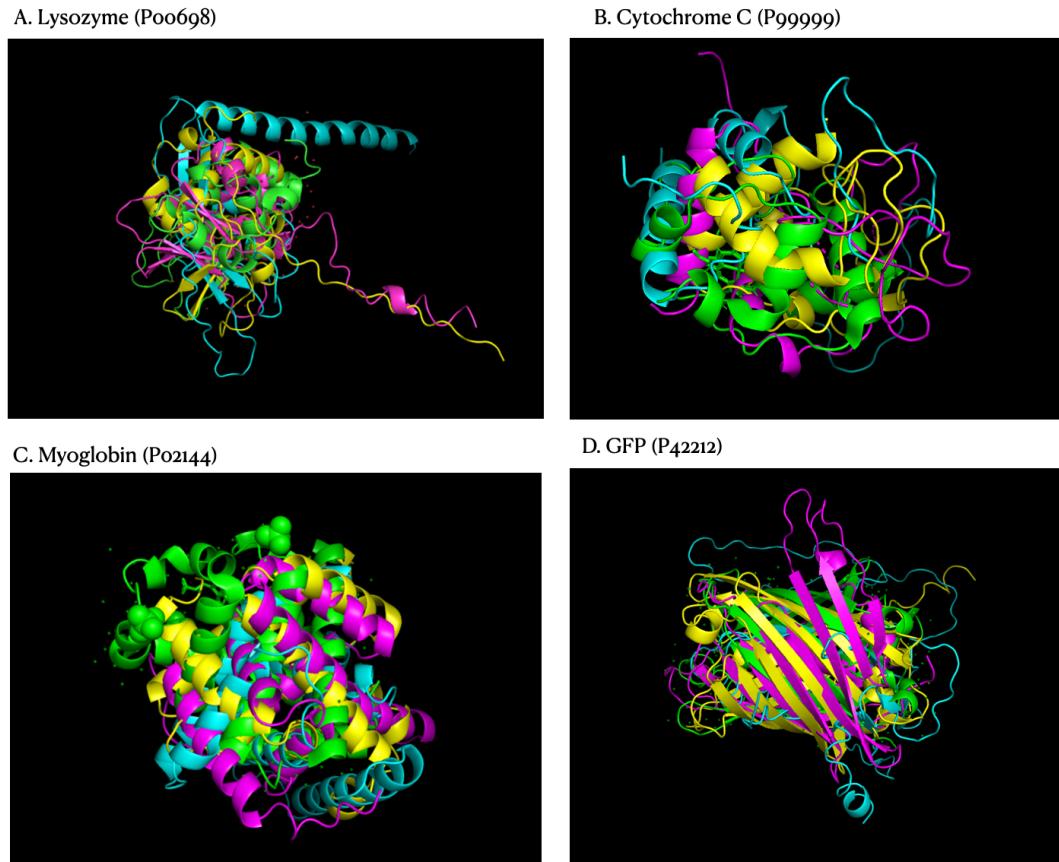


Figure 10: Experimental: green, Alpha2: blue, rosetta: pink, Alpha3: yellow.

Summary: For structural differences, Cytochrome C shows strong similarity across both the experimental and predicted structures, standing out compared to the other three proteins. In Lysozyme, the predicted structures differ slightly, especially in the extended tail region present in the predictions but absent in the experimental structure. For Myoglobin, the experimental structure includes an ion and a bound molecule, missing in all the predictions, though the core structural elements align well. Lastly, for GFP, while all models accurately predict the beta-sheet formation, the experimental structure's beta-sheets are more compact and tightly packed compared to the more extended forms seen in the predictions.

5 Restriction Enzyme ECORI binds to the DNA recognition site GAATTC. Generate structures for ECORI using Alphafold2, and Alphafold3. This is an opportunity to explore Alphafold3's capabilities as an all-atom capable model.

Solution: We picked ECORI from here: [ECORI](#), and we obtained the alphafold2 and alphafold3 results as the following:

A. Alphafold2 Prediction



B. Alphafold3 Prediction

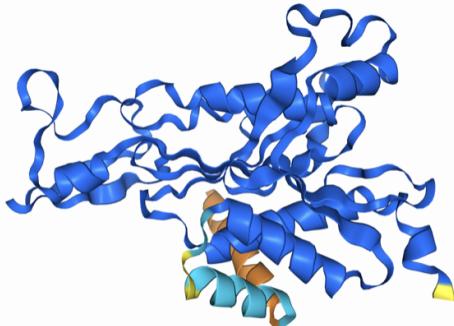


Figure 11: ECORI Structures

Summary: The comparison between Alphafold2 and Alphafold3 highlights the advantages of Alphafold3's all-atom capabilities. While Alphafold2 provides a reliable prediction of the protein's backbone and general fold, Alphafold3 offers a more refined, detailed atomic-level model, including accurate side-chain positioning. This added detail is particularly important for understanding specific molecular interactions, such as how ECORI binds to and cleaves DNA. The all-atom features of Alphafold3 allow for more precise modeling of protein-DNA complexes and active site interactions, making it a superior tool for studying enzymes and their mechanisms in finer detail.

6 Generate multiple different structures with Alphafold3 by prepending and appending nucleotides, both including and not including the complementary DNA recognition site. Is there a difference in Alphafold3's prediction when DNA is present?

Solutions: We will try three things here:

- Alphafold3 with just the ECORI sequence (no DNA).
- Alphafold3 with ECORI + GAATTC prepended or appended, since ECORI binds to the sequence GAATTC.
- Alphafold3 with ECORI + a random DNA sequence prepended or appended, and we picked TGCATG, and we do it with prepended.

The following is the result:

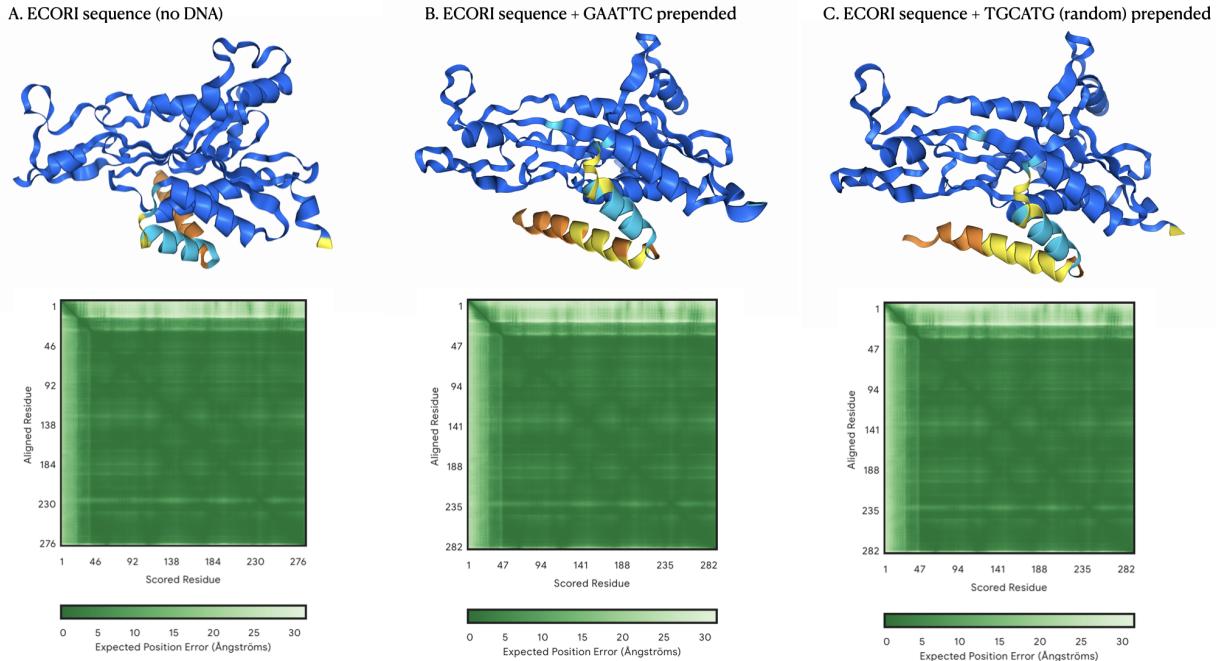


Figure 12: Alphafold3 Results

Summary: The predicted structures of ECORI, with and without DNA, show that the overall fold of the enzyme remains stable across all conditions. In all three cases—no DNA, GAATTCT prepended, and TGCATG (random DNA) prepended—the core alpha-helical structure of ECORI is maintained. The pLDDT charts demonstrate high confidence in Alphafold3’s predictions, with minimal changes in confidence across the different inputs. This indicates that the protein’s structure is consistently well predicted, regardless of the presence or type of DNA sequence.

However, there are subtle differences between the predictions when the recognition site GAATTCT is present compared to the random sequence. These variations, though slight, may suggest that Alphafold3 predicts some interaction between ECORI and its specific DNA substrate. The random DNA sequence (TGCATG) does not appear to induce any significant structural changes, indicating that non-specific DNA does not strongly affect Alphafold3’s predictions in this context.

Overall, while the predicted structures are similar, the inclusion of the recognition sequence hints at potential protein-DNA interactions, even if they are not captured dramatically in these predictions. The high confidence in all cases supports the robustness of Alphafold3 in predicting ECORI’s general structure, though capturing fine-grained binding interactions might require further refinement.

7 Look into Restriction Enzyme Biology a little deeper. What else would you propose adding to the Alphafold3 input to sharpen the prediction?

Solution: In order to further sharpen Alphafold3’s predictions for ECORI, several biological elements could be added to the input:

- One important factor is the inclusion of cofactors like Mg^{2+} ions, which are essential for ECORI’s enzymatic activity. Incorporating these metal ions in the prediction process could help Alphafold3 capture the enzyme’s active conformation more accurately, as Mg^{2+} is crucial for DNA cleavage and could induce specific structural rearrangements.

- Another potential improvement would be to better model DNA-protein interactions by allowing AlphaFold3 to account for the binding between ECORI and its recognition site, GAATTC. While AlphaFold3 is capable of predicting protein-DNA complexes, refining its ability to capture the dynamics of specific binding interactions could lead to more accurate predictions. This is especially relevant for restriction enzymes, where binding to DNA often induces conformational changes necessary for catalysis.
- Lastly, incorporating environmental factors such as pH and temperature could also enhance the predictive accuracy. These factors can affect protein folding and stability, particularly for enzymes that function under specific conditions. Although AlphaFold3 does not currently model such factors, future versions could benefit from simulating proteins in environments that mimic their physiological conditions.