
AlphaFold: A beginner's guide and in-depth exploration of the revolutionary AI tool and its inner workings.
Shravan Saranyan

Table of contents:

Abstract	2
Purpose of AlphaFold	2
Creation of AlphaFold	3
How AlphaFold works	4
Evoformer Stage	5
Structural Stage	6
How to use and run AlphaFold	7
Block 1 - Input Protein Sequences	8
Block 2 - MSA Options	8
Block 3 - Advanced Settings	9
Block 4 - Run Prediction	11
Block 5 - Display 3D Structure and Plots	14
Limitations and Potential Advancements	17
Limitations of AlphaFold	17
Potential Advancements for AlphaFold	18
Citations and Additional Resources	22

Abstract

AlphaFold is a very famous example of a machine learning program that was used to solve a generational problem in figuring out the 3D structure of a protein given its sequence, which it achieves to some degree. However, the program is often misunderstood, and for a long time wasn't very accessible to many researchers and the general public, until Colabfold was created for anyone to use AlphaFold in an easy, user-friendly format. While accessible, it can still be difficult for many researchers to navigate as the tool uses complicated techniques and jargon, so this paper aims to explain why AlphaFold was created and what problem it solves, how it works, a deep dive into the code of it all, and what it can and can't do, as well as how to potentially make it better. This paper aims to shine light and give insight into the world of AlphaFold in a vernacular that anyone can understand and interpret.

Purpose of AlphaFold

AlphaFold is an AI (artificial intelligence) program powered by a neural network developed by Google's DeepMind to *predict protein structures* [1,5]. It is entirely open source and is used to predict the 3D structure of proteins from an amino acid chain, below is a quick summary of this process (Fig. 1).

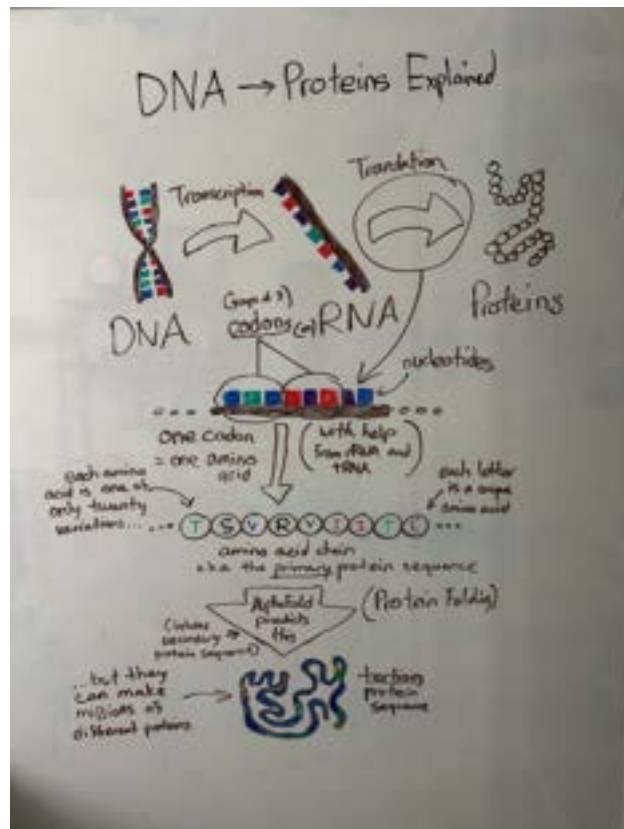


Figure 1. **DNA to Proteins Explained.** A refresher on how cells transcript and translate DNA into folded 3D proteins.

While DNA acts as the blueprint of genetic information for our cells, proteins are the core molecules executing the processes and functions dictated by our DNA. The journey from DNA to proteins involves two key steps: transcription (DNA to mRNA) and translation (mRNA to protein). During translation, our cells read the mRNA and segment it into groups of three letters called codons. Each codon is then translated into one of twenty different amino acids, which are linked together in a chain. Following the conversion of the RNA sequence into a chain of amino acids, the chain detaches and folds into a protein in a highly complex process, resulting in a 3D tertiary protein structure (see Fig. 5).

The behavior of proteins is often guided by their 3D structure and shape, which dictates how they engage and coordinate with other proteins in the body [2]. Understanding a protein's structure is critical for discerning its function and potential relevance in research and medicine. For example, in drug discovery and development, knowing the shape of a viral protein can enable researchers to design drugs to neutralize its function more effectively [3, 13]. Accurate structural predictions can also help scientists see how proteins malfunction or misfold to cause several proteinaceous diseases such as Parkinson's and Alzheimer's [4, 13]. Thus, a novel, cost-effective means to accomplish this became all the more valuable, which is what AlphaFold aims to achieve.

Creation of AlphaFold

Computational methods for predicting accurate 3D protein structures have usually focused on physical interactions or evolutionary history. Using the physical interactions approach involves a lot of molecular physics and thermodynamics and is generally impractical. Using the evolutionary history of protein sequences is a more practical and efficient alternative. To determine the protein structure, this method analyzes the evolutionary history of the protein, similarities to previously solved proteins, and the relationships between similar pairs of amino acids (residues) that have evolved across different species (known as pairwise evolutionary correlations) [5, 13]. This method greatly benefits from experimental protein structures deposited in the Protein Data Bank (PDB) [30], the genome sequencing of more species to derive more protein sequences, and better deep learning neural network interpretations to interpret these protein structures more efficiently [6].

The evolutionary history method proved viable with the creation of AlphaFold, the neural network behind the first computational approach to predict protein structures with high accuracy. AlphaFold structures have a median backbone (core of the protein) accuracy of within 0.96 angstroms (an angstrom is 0.1 nm, used to measure atomic structures); the next best model has

an accuracy of within 2.8 angstroms, and it also models highly accurate side chains. The overall accuracy of AlphaFold is within 1.5 angstroms compared to 3.5 angstroms of the nearest competitor; it is also highly scalable to very long proteins [5] (Fig. 2).

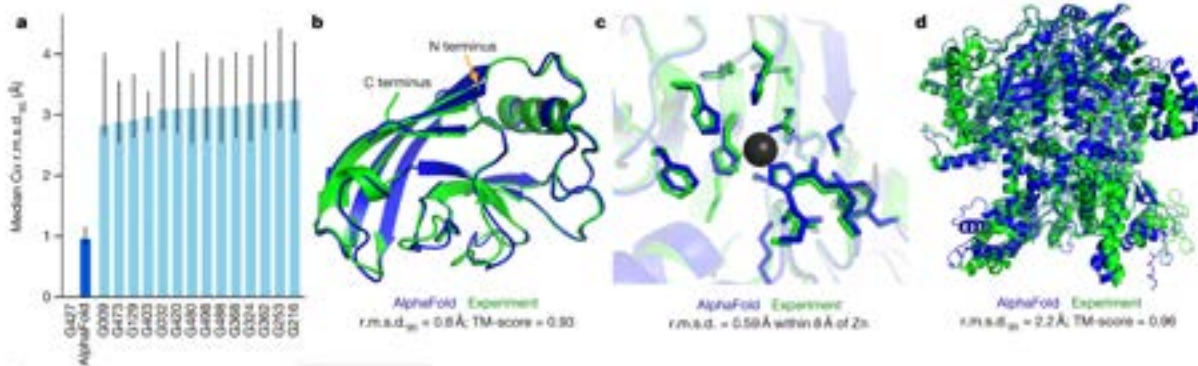


Figure 2. **Examples of AlphaFold Accuracy.** Above is the median backbone accuracy of AlphaFold (G427) in angstroms compared to other models (G009-G216) and some example outputs (a, b, c, d). [5].

How AlphaFold works

A lot of the processes that occur below in AlphaFold are powered by neural networks, which are a kind of machine learning model that mimics how the human brain works (see Fig. 3). The fundamental unit of neural networks are neurons, which receive an input, process it, and give an output. Several of these neurons are grouped into columns, known as layers. The neural network contains an input layer, an output layer, and hidden layers which perform computations. Each neuron in a layer will feed its output as an input to the neurons in the next layer. Still, the importance of each input the neurons in the next layer receive is determined by the weight or strength of the connection the neurons have; the stronger the weight, the more influential the input is. There are also additional biases in the network, which are independent parameters added to the inputs of neurons to adjust the output, along with the weights from the previous neurons. The neural network is then asked to provide a desired output, and it tweaks itself with all of its variables until its output is constantly close to the desired output [7, 13]. By doing this training to maximize the confidence and accuracy of its protein model, AlphaFold can consistently pump out accurate protein structures confidently given the input of a protein sequence.

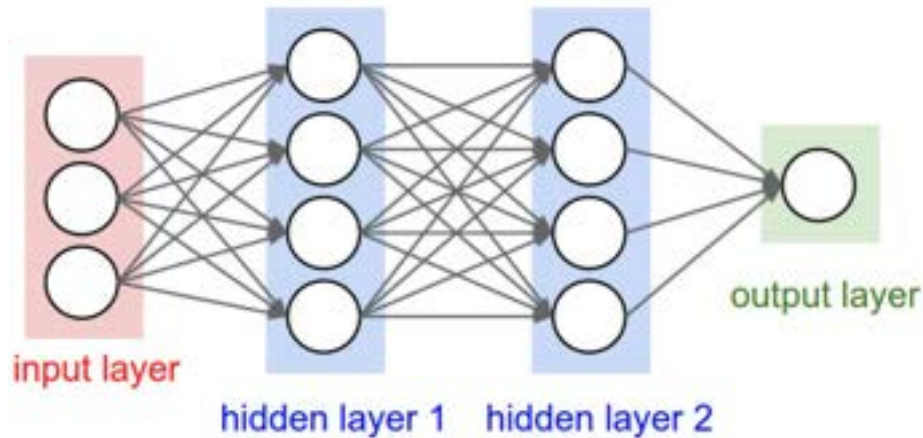


Figure 3. **Neural Network Diagram.** In the example above, the layers are labeled, the dots are the neurons, and the arrows are the weights from the previous layers. [8].

The remarkable thing about neural networks is that we don't know what the inside of them looks like; the network itself creates its own architecture. We humans merely give a desired example output some example inputs and let the network figure out how to best train itself to provide the desired output given an input [7]. This is why neural networks are often extremely complicated and convoluted. There are many different variations of neural networks, including networks that monitor their progress layer by layer, cycle back to previous layers, and use multiple dimensions to map spatial features [9, 13]. AlphaFold includes all of these different kinds of neural networks at some stage or another, and they are the reason AlphaFold can do what it does.

The AlphaFold network comprises of two main stages:

Evoformer Stage

In this main stage, the MSA (Multiple Sequence Alignments), and pairwise features of the input protein are processed through repeated layers of Evoformer, a novel neural network block [5] (Fig. 4a).

1. MSA—Multiple Sequence Alignments are protein sequences similar to the desired protein sequences that *are all aligned together*. They provide key information about evolutionary history and conservation.
2. Pairwise Features - The relationships between *any pair of two residues* (what's left of amino acids when they join into a chain) within the sequence contained within a Pair representation.

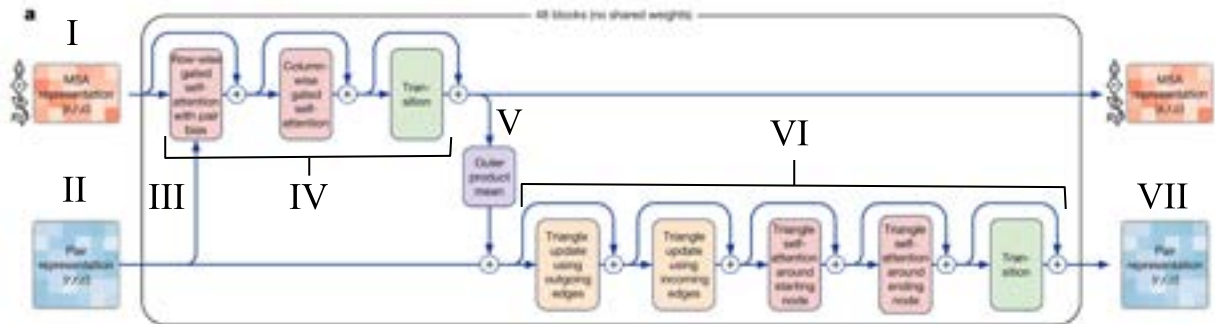
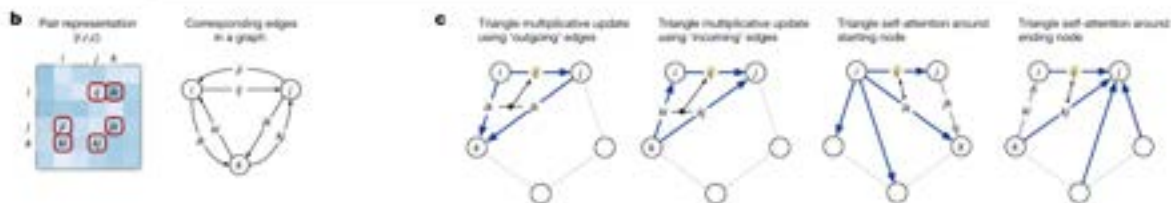


Figure 4a. **Architectural Details of Evoformer Stage.** A visual representation of the Evoformer Stage of AlphaFold, contains seven major steps (I-VII). [5].

In every single block (layer), the MSA representation (I) is first refined using the Pair representation containing the Pairwise Features (II), and then it updates the Pair representation. The MSA representation refines itself by using the Pairwise Features (III) to find evolutionary information and context and then transitioning itself to absorb this information (IV). Then, it provides this information back to the Pair representation (V), which updates its Pairwise Features by graphing them as triangles and adjusting each edge and node to minimize interference or error (VI, see Fig. 4b). This process is repeated iteratively 48 times, each time gaining more evolutionary context and becoming more accurate. The end product is the original protein sequence (amino acid chain) and the Pair Representation, which now contains all the evolutionary context and information about the protein that can be taken advantage of in the next stage in the form of relations between residues (VII) [5, 13].

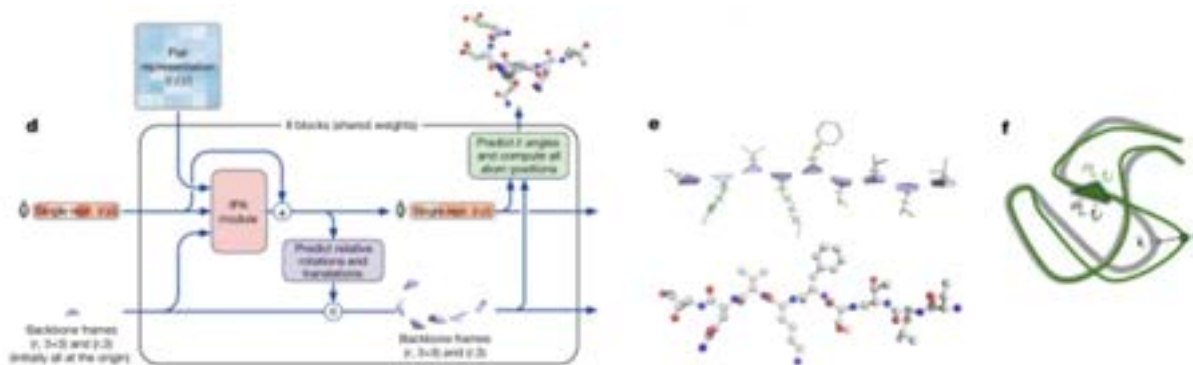


Figures 4b-c. **Pair Representation Details.** Above is a visual representation of how select pairs of residues are weighted in a pair representation and shown as a graph representation (b), as well as how it self-updates the graph to minimize interference (c). [5].

Structural Stage

The structural stage takes in the refined pair representation of every possible pair of residues from the previous stage and the original protein sequence, then determines the 3D structure of the protein and generates it (Fig. 4d). The main body of the protein, known as the basic backbone structure, is represented as a series of independent rotations and translations to a frame (known as a residue gas representation, shown in Fig. 4e) for each residue in the

sequence. The rotations and translations that iteratively update the frame are determined by the IPA module (which takes in the protein sequence and pair representation, shown in d) and then are enacted by an update operation. Initially, when modeling, the rotations and translations follow certain prioritizations, but the complex rules behind the geometry of peptide chains (protein chains) are ignored to allow for the specific refinement of each part of the chain and then are factored at the end using a violation loss term. After every block has finished, the stage uses the protein sequence and fully updated backbone frames to create a predicted 3D molecular model of the folded protein and provides the accuracy of the predicted structure compared to the true structure if known (Fig. 4c) [5, 13].



Figures 4d-f. **Architectural Details of Structural Stage.** Above is the Architectural Details of the Structural Stage of AlphaFold (d), example of the backbone frames (e), and the algorithm to check structural accuracy(f). [5].

AlphaFold predicts two kinds of protein structures, **monomers**, and **multimers**. Monomers are proteins folded from a single protein chain, whereas multimers are composed of several protein chains that fold together into a larger structure. Overall, AlphaFold inputs the primary protein structure, which is the amino acid chain, and during the structural stage, the model predicts the tertiary protein structure, which is either given as the output (if monomer) or combined in an additional stage with the other tertiary protein structures for each protein (amino acid) chain to form a quaternary protein structure (if multimer) [13] (see Fig. 5).

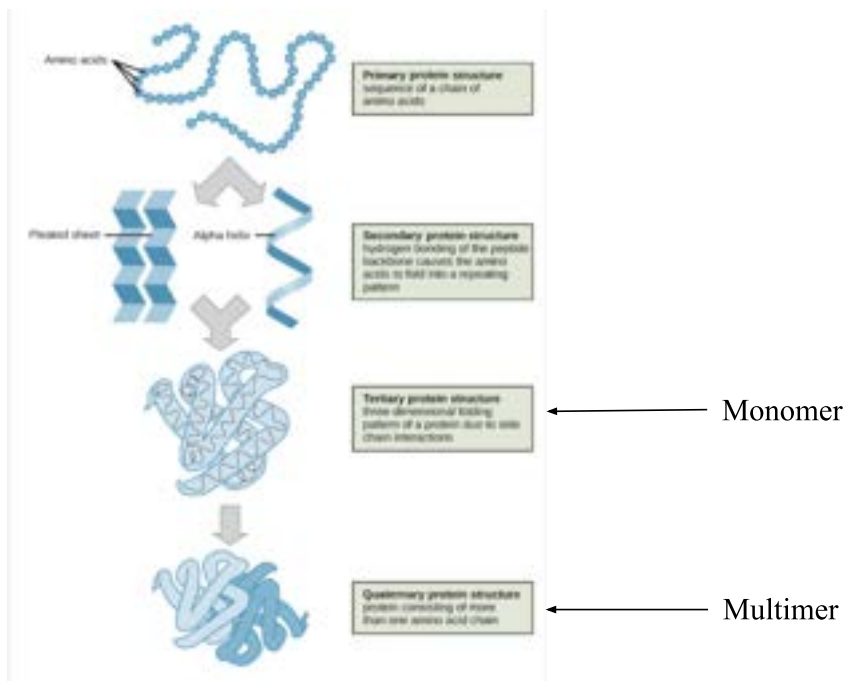


Figure 5. **Four Levels of Protein Structure.** Above are the four stages of protein structure as the protein folds from a one-dimensional chain into a three-dimensional structure, and which stage AlphaFold predicts as its output depending on the protein type (monomer/multimer). [10].

How to use and run AlphaFold

AlphaFold's code can be found on GitHub and run locally or on Google Colab. Google Colab is a tool that Google recently created that runs code on Google's servers. It's free to use, beginner-friendly, and generally more efficient than running code locally on a computer [11].

Here is the Colab file for AlphaFold [15]:

<https://colab.research.google.com/github/sokrypton/ColabFold/blob/main/AlphaFold2.ipynb>

And here is the Github repo that contains the code [16]:

<https://github.com/sokrypton/ColabFold>

We will now go block by block through the code in two stages. First, we will look at the Colab user interface and explain what each stage does and its settings. Then, we will delve deeper into the actual code of AlphaFold below to see how these processes are expressed and run.

There are several stages of the code in the Colab version; they can be divided as such:

- Block 1 - Input Protein Sequences
- Block 2 - MSA options

- Block 3 - Model and other Advanced settings
- Block 4 - Run Prediction
- Block 5 - Display 3D structure

Block 1 - Input Protein Sequences

In this block, the desired protein sequence is inputted, and a test run is started (Fig. 6a).

The screenshot shows a Jupyter Notebook cell with the following content:

```

query_sequence: PIAQIHLEGRSDEQKETLREVSEASRSLSAPLTSVVRVITEMAKDHFQGGELASK
+ Use | to specify inter-protein chainbreaks for modeling complexes (supports homo- and hetero-oligomers). For example PL-SK-PL-SK for a homodimer
jobname: test
num_relax: 0
+ specify how many of the top ranked structures to relax using amber
template_mode: none
+ none = no template information is used. pdb100 = detect templates in pdb100 (see notes). custom = upload and search own templates (PDB or mmCIF format, see notes)
Show code
-----
jobname: test_a5e17
sequence: PIAQIHLEGRSDEQKETLREVSEASRSLSAPLTSVVRVITEMAKDHFQGGELASK
length: 58
  
```

Figure 6a. **Block 1**. Screenshot of Block 1 in AlphaFold2.ipynb (Colabfold). [15].

The **query_sequence** is where the desired protein is inputted. Each letter stands for one amino acid residue [18]. The sequence is formatted to capitalize and remove spaces. A colon should be used to separate each protein chain within a complex (multimer) [15].

jobname is used to create a name for the test run, a hash is created using the inputted query_sequence and is tacked onto this string to create a unique identifier for the test run; in this example, “test_a5e17”

num_relax determines how many “top ranked structures” go through an additional refinement process called **AMBER**, which uses molecular dynamics to create a more realistic and accurate model [12, 13].

template_mode is used to find templates if asked. Templates are similar protein structures that have already been determined and can be used to help predict the structure [15, 19].

The output of this block is the identifier of the test run (jobname), the formatted sequence, and the length of the sequence in residues [13, 15].

Block 2 - MSA Options

In this block, the options for the Multiple Sequence Alignments (MSA) mode and pairing mode are configured (Fig. 6b).

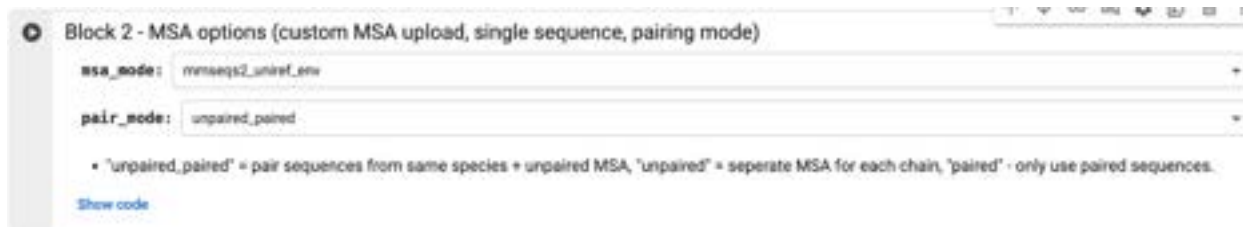


Figure 6b. **Block 2**. Screenshot of Block 2 in AlphaFold2.ipynb (Colabfold). [15].

The **msa_mode** determines what kind of MSA is used, the three main modes are:

mmseqs2_uniref(_env): These modes use a computer-generated MSA derived from the mmseqs2 server [11, 14], and are the default option.

custom: This mode uses a custom user-inputted MSA, and is useful as a more advanced option for testing certain MSAs or aiming for more customized results [11].

single_sequence: This mode only uses the original sequence and no others, and is usually used for benchmarking MSAs [13, 15].

The **pair_mode** determines if the model should use an unpaired MSA, paired MSA, or both [15]:

unpaired: Uses a regular, unpaired MSA, which identifies similar sequences for *only one chain*, and is used to find intra-chain coevolutionary information; is used by default for monomers [22, 23].

paired: Uses a paired MSA, which identifies similar sequences for *multiple chains*, and is used to find inter-chain coevolutionary information; is used by default for homo-multimers (complexes of identical chains) [22, 23].

unpaired_paired: Uses both an unpaired and paired MSA to gain coevolutionary information for both each individual chain and the entire complex as a whole; is used by default for hetero-multimers (complexes of differing chains) [22, 23].

Note: When using unpaired_paired, it will default the pair_mode to whatever is dictated by the model_type, which itself is dependent on whether the sequence is a monomer, homo-monomer, or hetero-monomer [23].

Block 3 - Advanced Settings

These settings for the model are advanced, optional, and mostly self-explanatory. The most important and interesting ones are covered below (Fig. 6c).

Block 3 - Advanced settings

model_type:

- if auto selected, will use `alphafold2_ptm` for monomer prediction and `alphafold2_multimer_v3` for complex prediction. Any of the `model_types` can be used (regardless if input is monomer or complex).

num_recycles:

- if auto selected, will use `num_recycles=20` if `model_type=alphafold2_multimer_v3`, else `num_recycles=3`.

recycle_early_stop_tolerance:

- if auto selected, will use `tol=0.5` if `model_type=alphafold2_multimer_v3` else `tol=0.0`.

relax_max_iterations:

- max amber relax iterations, 0 = unlimited (AlphaFold2 default, can take very long)

pairing_strategy:

- greedy = pair any taxonomically matching subsets, complete = all sequences have to match in one line.

Sample settings

- enable dropouts and increase number of seeds to sample predictions from uncertainty of the model.
- decrease `max_msa` to increase uncertainty

max_msa:

num_seeds:

use_dropout:

Save settings

save_all:

save_recycles:

save_to_google_drive:

- if the `save_to_google_drive` option was selected, the result zip will be uploaded to your Google Drive

dpi:

- set dpi for image resolution

Figure 6c. **Block 3.** Screenshot of Block 3 in AlphaFold2.ipynb (Colabfold). [15].

The **model_type** is very important, and determines which version of AlphaFold to use [11, 15] (the pros and cons of each as well as an example are outlined in this resource [21]).

num_recycles determines how many times the program is recycled, i.e., when the output of AlphaFold is fed back in for refinement, by default it happens three times [5, 11, 13, 15].

relax_max_iterations limits the amount of iterations AMBER can run for if AMBER is activated, 0 means no limit is imposed and it can take very long to run [15].

The **pairing_strategy** of the MSA can be either greedy or complete, complete means that every single pair of residues is checked, while greedy means that only selective, informative pairings are chosen. Greedy is the default and is more efficient, but complete is more accurate [13, 15].

max_msa gives the option to limit the size of the MSA for the sake of efficiency while trading off accuracy [13, 15].

use_dropout gives the option to use Dropout, a technique where neurons are randomly removed during the neural network stages to prevent overfitting, when the model becomes too

comfortable in its current state and refuses to change and evolve. The randomness of this and other parts of the model are determined by random seeds (**num_seeds**); more seeds prevent variance from the random initialization of each seed [13, 15, 17].

Block 4 - Run Prediction

This block is where AlphaFold is actually run, so we will look at the code and show the most relevant processes that occur within the code, and where everything ties in. This will include both the code from the ColabFold (Fig. 6d-e) and the code from GitHub (Fig. 7) [15, 16].

```
# Imports required libraries and functions for ColabFold
import sys
import warnings
warnings.simplefilter(action='ignore', category=FutureWarning)
from Bio import BiopythonDeprecationWarning
warnings.simplefilter(action='ignore', category=BiopythonDeprecationWarning)
from pathlib import Path
from colabfold.download import download_alphafold_params, default_data_dir
from colabfold.utils import setup_logging
from colabfold.batch import get_queries, run, set_model_type
from colabfold.plot import plot_msa_v2

import os
import numpy as np
```

Figure 6d. **Imports for Block 4.** Screenshot of required libraries and functions for ColabFold being imported in AlphaFold2.ipynb (Colabfold). [15, 16].

The portion above (Fig. 6d) imports several modules from the colabfold libraries, the key module that runs AlphaFold is *colabfold.batch.run* [20], which calls the run() function from the file batch.py within the colabfold folder. Note: Several other libraries (such as get_queries) are run separately and inputted into the run() function as parameters rather than natively [16, 20].

```
results = run(
    queries=queries,
    result_dir=result_dir,
    use_templates=use_templates,
    custom_template_path=custom_template_path,
    num_relax=num_relax,
    msa_mode=msa_mode,
    model_type=model_type,
    num_models=5,
    num_recycles=num_recycles,
    relax_max_iterations=relax_max_iterations,
    recycle_early_stop_tolerance=recycle_early_stop_tolerance,
    num_seeds=num_seeds,
    use_dropout=use_dropout,
    model_order=[1,2,3,4,5],
    is_complex=is_complex,
    data_dir=Path("."),
    keep_existing_results=False,
    rank_by="auto",
    pair_mode=pair_mode,
    pairing_strategy=pairing_strategy,
    stop_at_score=float(100),
    prediction_callback=prediction_callback,
    dpi=dpi,
    zip_results=False,
    save_all=save_all,
    max_msa=max_msa,
    use_cluster_profile=use_cluster_profile,
    input_features_callback=input_features_callback,
    save_recycles=save_recycles,
    user_agent="colabfold/google-colab-main",
)
```

Figure 6e. **Run function in Block 4.** Screenshot of the run function [20]; which is being called in AlphaFold2.ipynb (Colabfold). [15, 20].

The run() function takes in a lot of parameters (Fig. 6e), including the query sequence (Block 1) and all the MSA and advanced settings (Blocks 2-3). The run function within AlphaFold then outputs five 2D predicted structures and ranks them from most to least accurate [15]. The function behaves like a main function within Colabfold (although there is another main() function within the GitHub version), and as such, offloads most of its processes to other functions that it calls within itself or takes in as inputs [16, 20]. The most important functions that directly relate to the prediction itself are shown below in order of where they are defined in the code [20]:

predict_structure() (Lines 324-542) - This function contains the heart of AlphaFold; it takes the sequence and features from the MSA as inputs and outputs the predicted structure. It generates a random seed (more if requested), loads in the AlphaFold models, and the parameters that control how they behave, and preprocesses the sequence, MSA, etc. After this, it takes in all of these and predicts a structure, rates the model confidence of each part of the structure, and outputs it as a protein object. Finally it reranks the models based on the predicted confidence [13, 20].

```

422     # predict
423     result, recycles = \
424     model_runner.predict(input_features,
425                          random_seed=seed,
426                          return_representations=return_representations,
427                          callback=callback)
428
429     prediction_times.append(time.time() - start)
430
431     #####
432     # parse results
433     #####
434
435     # summary metrics
436     mean_scores.append(result["ranking_confidence"])
437     if recycles == 0: result.pop("tol",None)
438     if not is_complex: result.pop("iptm",None)
439     print_line = ""
440     conf.append({})
441     for x,y in [{"mean_plddt", "pLDDT"}, {"ptm", "pTM"}, {"iptm", "ipTM"}]:
442         if x in result:
443             print_line += f" {y}={result[x]:.3g}"
444             conf[-1][x] = float(result[x])
445     conf[-1]["print_line"] = print_line
446     logger.info(f"[tag] took {prediction_times[-1]:.1f}s ({recycles} recycles)")
447
448     # create protein object
449     final_atom_mask = result["structure_module"]["final_atom_mask"]
450     b_factors = result["plddt"][:, None] * final_atom_mask
451     unrelaxed_protein = protein.from_prediction(
452         features=input_features,
453         result=result,
454         b_factors=b_factors,
455         remove_leading_feature_dimension=("multimer" not in model_type))

```

Figure 7. **Portion of the predict_structure() function.** Screenshot of a relevant portion of the predict_structure() function within the run() function. [15, 20].

The figure above (Fig. 7) shows the prediction (422-429), model confidence (435-446), and creation of protein objects (449-455) [20].

pair_sequences(): Pairs the query sequence(s) and a paired MSA together into a larger MSA, only used for multimers [13, 20].

pad_sequences(): Aligns the sequences within an unpaired MSA and formats them to make sure the sequences all have the same length by adding spaces, which are only used for an unpaired MSA (usually monomers) [13, 20].

get_msa_and_templates(): Finds and obtains the unpaired MSA, paired MSA, and templates (if chosen) that match the query_sequence(s) [13, 20].

build_monomer_feature(), build_multimer_feature(): These functions parse through the unpaired MSA and paired MSA, respectively, and create features for each depending on the MSA pair type. Features include key information derived from the MSA, the sequences within, and the MSA itself [13, 20].

process_multimer_features(): Additional processing for the multimer features, as they are more complex [13, 20].

pair_msa(): Takes in several unpaired MSA for each chain in a multimer and combines them into a paired MSA that contains additional features that describe how the chains link together and form a system [13, 20].

generate_input_features(): Creates sufficient paired and unpaired features using some of the functions above to input into the structure prediction [13, 20].

unserialize_msa(): Loads in an MSA stored in a FASTA format (text-based way to represent protein sequence data in an efficient and human-interpretable way) and converts it back to a machine-usable format for analysis and processing [13, 20, 24].

msa_to_str(): Adds an unpaired or paired MSA to a FASTA string (i.e. serializes the MSA) [13, 20].

Below is a visual graph of these functions ordered by when they are run in the program [20].

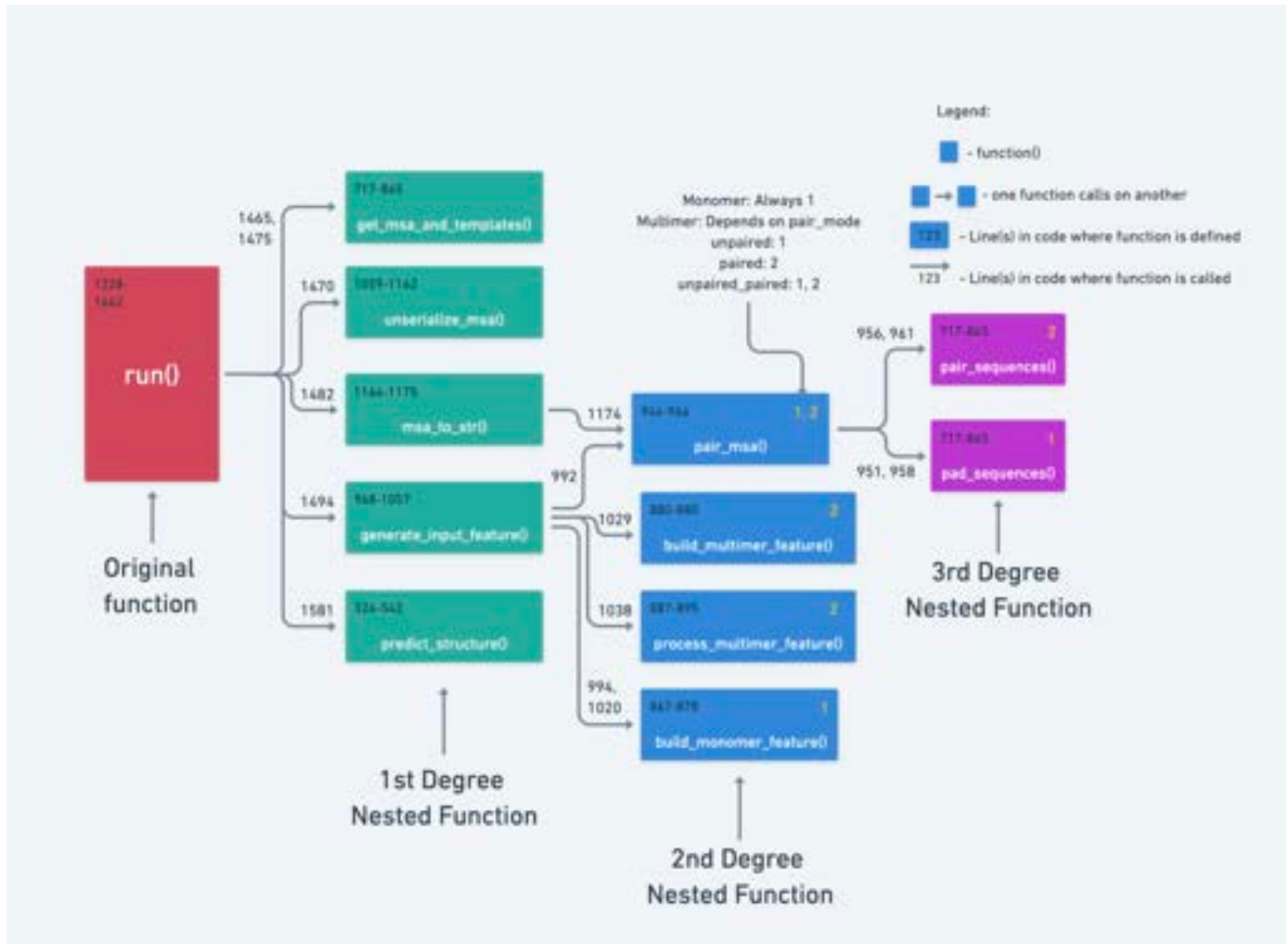


Figure 8. **Visualization of how the run() function is executed.** A visualization of how the run() function is executed, containing information about where each function is defined and called with line numbers and additional info on how the functions are structured when executed within the run function code on Github.

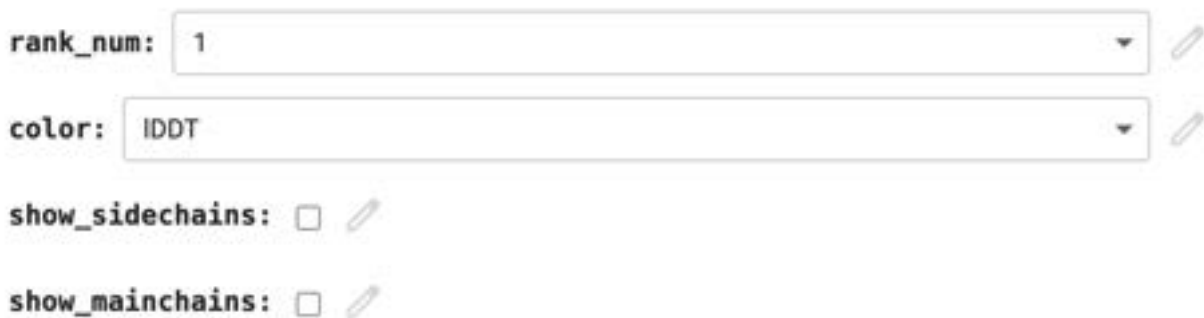
To read this graph (Fig. 8), start from run() and follow the arrows from top to bottom, left to right, and from box to box (function to function). The program runs these functions in this order: each colored box represents how far each function is called respective to run() (1st, 2nd, 3rd), the numbers inside the box (upper left) represent the lines of code where the function is defined, and the numbers adjacent to the arrows represent the lines of code where the function is called upon by the previous function. Additionally, the yellow numbers show which functions are run depending on the pair_mode (which in turn is dependent on monomer vs multimer) [22, 23].

To summarize what the run() function does in order, it gets the desired MSA(s) and templates (optional), converts the MSA(s) into a machine-readable format (unserialized them), then processes them (depending on monomer vs multimer it pairs or pads them) and serializes it

back into text. Then, the input features are created, the MSA(s) are processed again using the same method as before, and specific features are built and processed to input specifically into AlphaFold (again, depending on monomer vs multimer). Finally, the protein structure is predicted using the features derived from the MSA(s) [15, 20].

Block 5 - Display 3D Structure and Plots

The final portion of the Colab version of AlphaFold involves displaying the predicted 3D structure of the protein (Fig. 9-10c) [15].



The screenshot shows four configuration options for displaying the 3D structure:

- `rank_num:` A dropdown menu with the value `1` selected.
- `color:` A dropdown menu with the value `IDDT` selected.
- `show_sidechains:` A checkbox that is currently unchecked.
- `show_mainchains:` A checkbox that is currently unchecked.

Figure 9. **Block 5.** Screenshot of Block 5 in AlphaFold2.ipynb (Colabfold). [15].

During the structure prediction from the previous block, the model generates five predicted structures (by default) and ranks them 1 to 5, most to least accurate. **rank_num** simply states which one to use; by default, it uses “1” (see Fig. 9).

The **color** paints the protein structure based on the confidence of the model (IDDT), the different chains present (chain), or aesthetically (rainbow). The most practical of these options is IDDT, which is the default and shown below (in Fig. 10a):

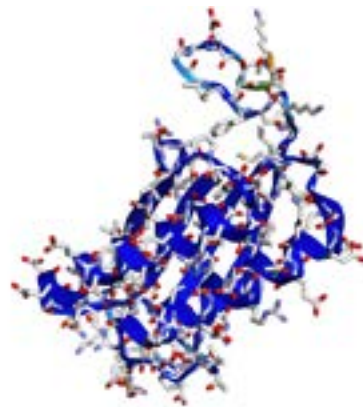


pLDDT: ■ Very low (<50) ■ Low (60) ■ OK (70) ■ Confident (80) ■ Very high (>90)

Figure 10a. **Default 3D Predicted Protein Structure.** Screenshot of the 3D protein structure predicted by AlphaFold2.ipynb (Colabfold). [15].

The model above (Fig. 10a) is very confident, and as such is mostly blue, except for the ends, which are more orange, as the model is not very confident in that section. This 3D structure can be endlessly rotated in all axes to see the full structure, as well as zoomed in.

show_sidechains and **show_mainchains** are options that display the additional molecular chains that come along with the backbone structure predicted by AlphaFold; see Fig. 10b.



pLDDT: ■ Very low (<50) ■ Low (60) ■ OK (70) ■ Confident (80) ■ Very high (>90)

Figure 10b. **3D Predicted Protein Structure with sidechains and mainchains.** Screenshot of the 3D protein structure predicted by AlphaFold2.ipynb (Colabfold) with sidechains and mainchains displayed as well. [15].

The model above (Fig. 10b) is more accurate in terms of how the protein would appear, but it does look very messy and makes it harder to see the more relevant structure, so these modes are usually turned off.

The Colab version of AlphaFold additionally creates plots for bioinformatics purposes (Fig. 10c).

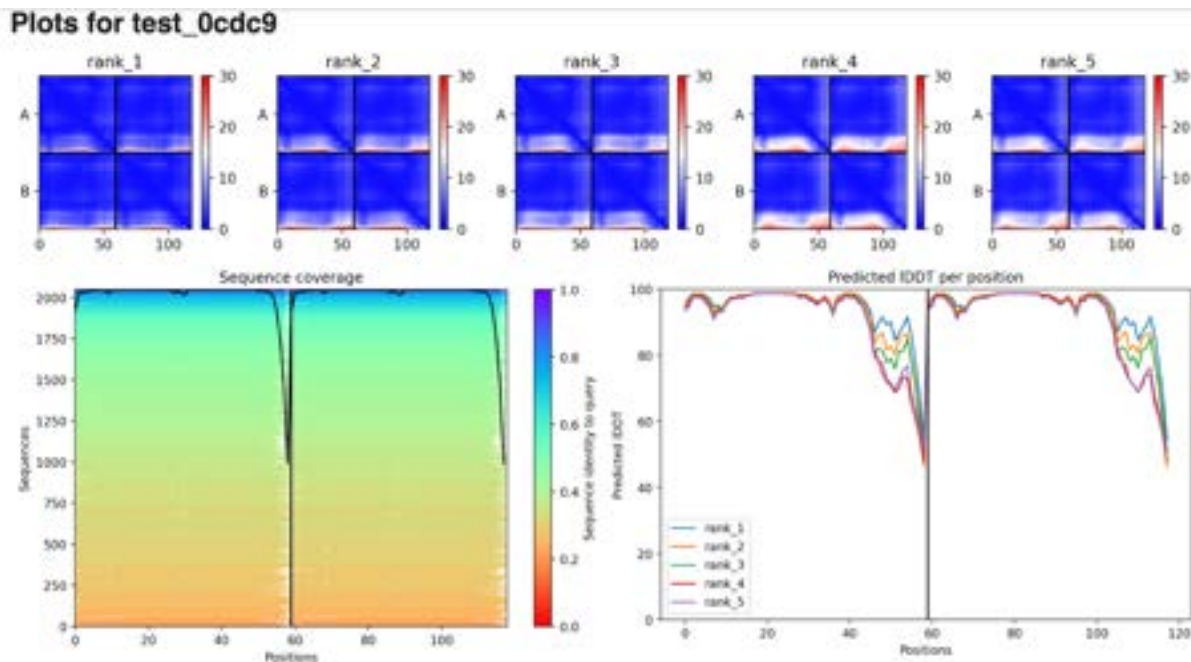


Figure 10c. **Assorted plots for the model.** Screenshot of the plots generated by AlphaFold2.ipynb (Colabfold) for the test model. [15].

The square graphs show the accuracy and correlation of the pairwise representations for every ranked model, the bottom two graphs show MSA coverage (the bluer the line the more related it is to the original) and the confidence level by position of the model for every ranked model [15].

Limitations and Potential Advancements

Despite all of the great advances pioneered by AlphaFold in the fields of molecular biology and artificial intelligence, the tool is still very new and has plenty of limitations, so this section will cover some of these limitations as well as a potential improvements and ideas that could make AlphaFold more accessible and a better tool.

Limitations of AlphaFold

AlphaFold sometimes struggles with predicting new or unusual proteins, known as “orphans,” because they have lesser close relatives, which is how it derives evolutionary information to predict protein structures. These protein predictions result in low-accuracy predictions with low confidence scores (red/orange in IDDT), and they can be made far worse if the “orphan”

proteins have little to no related structures in the PDB (Protein Data Bank). This is why it's essential to update the PDB with as many structures as possible, as AlphaFold excels at predicting new proteins if there are enough related sequences in the PDB to produce a quality MSA for prediction [25, 30].

Additionally, it struggles with high variation within proteins, whether that be regular point mutations within DNA that are then translated to altered residues within the protein or proteins with inherently high variation within sequences, such as antibodies. This is because of both a lack of data on the effect of variations within proteins and AlphaFold's focus on long-term evolutionary changes rather than physical forces that cause short-term changes such as radiation, mutation, etc [25].

While performing their functions, proteins often change their structure to adapt their function (a great example is the *complement system* [26]), but only a few structures have their potential conformations in the PDB (which AlphaFold uses for reference. By default, AlphaFold only predicts static structures, but researchers have found ways to trick AlphaFold to predict different versions of the proteins) [25].

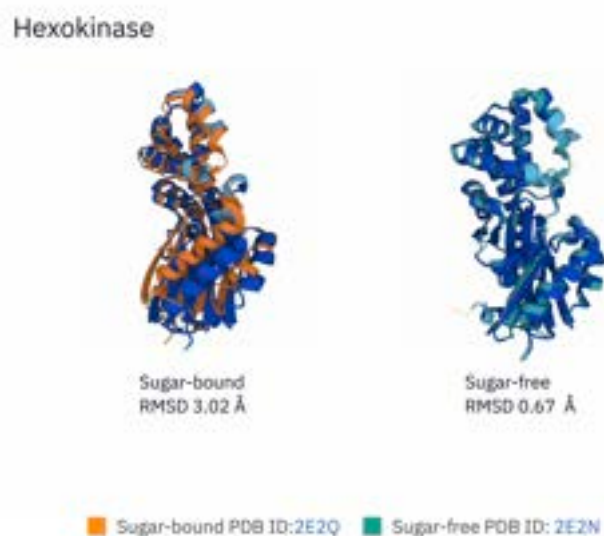


Figure 11. **Accuracy of Sugar-bound vs Sugar-free versions of Hexokinase.** Screenshot of the plots generated by AlphaFold2.ipynb (Colabfold) for the test model. [25].

Above (Fig. 11) is an example of how a protein can conform to a completely different structure when bonding with a different molecule. Here, when the protein Hexokinase (right) bonds with a sugar molecule, it conforms to a sugar-bound structure (left), and when AlphaFold predicts both, it has a lot more trouble predicting the sugar-bound structure (within 3.02 angstroms to within 0.67 angstroms, the lower, the more accurate) [25].

There are also several functions within the process of protein folding and interactions that AlphaFold wasn't designed to do [25]. These include but are not limited to:

- Other molecules that interact (but don't bond) with proteins, such as ions, nucleic acids (DNA, RNA), and other non-proteins that work closely with proteins and may provide key context on the protein's function.
- Any modifications that happen to the protein after translation.
- Free nucleic acids (unfolded DNA or RNA that didn't translate into proteins).
- Any ions or other molecules (ligands) that bond post-translation (AlphaFold usually predicts them when absent, a false positive).
- Anything to do with the Membrane Plane (a fluid membrane made up of fat lipids, complex molecules, and proteins within cells), including a protein's relative position to it (if the protein spans through it).

Potential Advancements for AlphaFold

In this final section, we will cover examples of potential advancements that can be made to AlphaFold, ranging from small (portion of AlphaFold), to medium (AlphaFold as a whole), to large (the field as a whole).

Small advancements:

When trying to input a protein into AlphaFold, scientists often go onto the *NIH library* [27], and search for proteins, looking for the protein sequence listed at the bottom. The problem is the sequence usually looks like the sequence below (Fig. 12a).

```
ORIGIN
  1 mttcsrqrfts sssmkgscgi gggiggssr issvlaggsc rapstygggl svssrfssgg
  61 acglgggygg gfsssssfgs gfgggygggl gagfggglga gfgggfaggd gllvgsekv
 121 mqnldrlas yldkvralee anadlevkir dwyqqrqse ikdyspyfkt iedlrnkiaa
 181 atienaqpil qidnarlaad dfrtkyehel alrqtveadv nglrrvldel tlartdlemq
 241 ieglkelay lrknheeml alrgqtggdv nvemdaapgv dlsrilnemr dqyeqmaekn
 301 rrdawtfls kteelnkeva snsolvqssr sevtelrrvl qgleielqsq lsmkaslens
 361 leetkrycm qlsqiqglig sveeqlaqlr cemeqqsqey qilldvktrl eqeiatyrll
 421 legedahlss qqasqgsyss revftsssss ssrqtrpilk eqssssfsqg qss
```

Figure 12a. **Example of Protein sequence from NIH library.** Screenshot of how a protein sequence would appear in the NIH library. [27].

The sequence above (Fig. 12a) isn't compatible with the input the way it is formatted. This problem can easily be fixed by adding the code below (Fig. 12b), which removes the spaces (already done in the current version [15]) and then numbers and capitalizes the sequence.


```
# remove whitespaces
query_sequence = "".join(query_sequence.split())
# Allows for spaces and numbers in query sentence by removing them
query_sequence = ''.join([i for i in query_sequence if not i.isdigit()])
query_sequence = query_sequence.upper()
```

Figure 12b. **Code snippet that fixes the formatting issue.** Screenshot of a code snippet that would fix the formatting issue with the NIH library (the two lines below the second comment). [27].

This simple snippet of code (Fig. 12b), when added, saves scientists a lot of time and effort from removing the digits and capitalizing manually and more seamlessly integrates the library of protein sequences and the program.

Medium advancements:

There are several areas where AlphaFold could be advanced, from the code to the interface to the content of the program. Here are some potential improvements below:

(Note: These are all just ideas, which may vary in actual feasibility, some may or may not work)

- The ColabFold program takes a long time to run, especially for multimers and complex proteins. There are modes within the advanced settings of ColabFold that reduce the runtime [15] (such as reducing the size of the MSA or reducing the number of ranks generated), but all of those modes have to be inputted manually, and for a novel protein, it can be hard to tell the correct tradeoff between accuracy and efficiency. So potentially, a time-saving mode could be created that would automatically tweak specific modes to reduce the runtime if it deems the tradeoff of accuracy to be worth it; this automatic mode could be powered by some neural network to run the ColabFold program with more efficiency and manage the tradeoff of accuracy.
- Instead of manually inputting a query sequence, the ColabFold program could have a search function that saves the user from having to search and find a protein sequence; they would only have to type the name of a protein, and it would predict the most common version (For searching for particular proteins the specific tag used in the NIH library [27] would have to be used)
- Because Google Colab uses a different computer when a runtime is restarted, any one-time-only downloads that AlphaFold requires must be repeated, which can take a lot of time (if using both amber and templates, 2 minutes are used just to download their

libraries, and another 90 seconds are used to download the specific AlphaFold model (ptm) [21]). This problem could be fixed by potentially working with Google and downloading all of these files onto some of their machines and assigning those specific machines to ColabFold users to save time on downloads.

- Additionally, restructuring the order in which ColabFold is run by placing any user-inputted parameters (`model_type`, `pair_mode`, `template_mode`, etc.) nearer to the end would make the turnaround for quick runs of the model with slight changes to certain parameters (for purposes of experimentation) far more efficient. Anything unchanged within ColabFold (such as downloading libraries, defining functions, etc.) could be run only once, so subsequent runs of ColabFold could be processed far quicker.
- Looking at the `model_type` options [21], some models had better overall accuracy (**ptm**), and some models had better point accuracy (**multimer_v2**, **deepfold_v1**). By using two algorithms, one to determine where each model does best and focus that model on what section it does best, and another one to stitch each piece of the protein model together, a “Frankenstein” model can be created with the best attributes from every available model. This model would be extremely time-intensive (measures like the ones above could be used to save time) but could potentially be more confident and accurate than its parts.
- There are *three general types of tertiary structures* (what AlphaFold outputs for monomers): **globular**, **fibrous**, and **membrane** [29]. While it’s usually pretty easy to tell which is which based on their shape, sometimes the protein may be more ambiguous. AlphaFold could help out by labeling every 3D structure it predicts as one of these protein types (or, in the case of multimers, what type each strand within is), which would make the model easier to interpret.
- AlphaFold provides plenty of plots and graphs with its output for bioinformatics purposes, but some of these graphs could do with adding in a key and labels for each axis, as well as a quick explanation on how to interpret the plots and graphs for beginners to more easily use the data, and to make the data more accessible.

Large advancements:

Large advancements are the ideas and concepts that will not just improve AlphaFold but elevate the entire field. There are three of these ideas below.

Improved accessibility to new proteins from the Protein Data Bank using web parsers.

There are thousands of new proteins discovered every year, but not all of the proteins discovered end up on the PDB. These novel proteins are usually from unconventional hosts

(any host besides humans/mice) and are usually still present in research documents across the internet. A web parser (or scraper) could theoretically comb through all of these overlooked proteins and automatically run them through AlphaFold (if enough similar proteins are found) and add them and their structure to PDB. This would significantly improve cooperation between researchers by making these overlooked proteins more accessible and shed more light on the fields they represent [30].

Adding features to AlphaFold to overcome some of its limitations.

Now that the concept of using machine learning to predict molecular structures has been proven, AlphaFold can move beyond some of its limitations by adding features, such as training the model to predict several variations of the protein, molecules binding to the protein, free DNA and RNA strands, and protein-nucleic acid complexes. This would dramatically improve the scope of AlphaFold's reach and would be an important stepping stone for future endeavors.

Predicting Protein Mimicry.

Protein Mimicry is a phenomenon where bacterial (or other foreign organism) proteins mold themselves into a shape similar to that of the host's own proteins, usually to avoid detection. When AlphaFold predicts a protein, it could also show bacterial proteins that have been known to mimic that specific protein. This would give important insights and information to researchers studying protein mimicry and expose more people to the field [31].

References and Additional Resources

- [1] "AlphaFold Protein Structure Database," European Bioinformatics Institute, Accessed: Aug. 31, 2024. [Online]. Available: <https://alphafold.ebi.ac.uk/>
- [2] B. Alberts, A. Johnson, J. Lewis, et al., *Molecular Biology of the Cell*, 4th ed. New York: Garland Science, 2002, The Shape and Structure of Proteins. [Online]. Available: <https://www.ncbi.nlm.nih.gov/books/NBK26830/>.
- [3] J. P. Hughes, S. Rees, S. B. Kalindjian, and K. L. Philpott, "Principles of early drug discovery," *Br. J. Pharmacol.*, vol. 162, no. 6, pp. 1239-1249, Mar. 2011, doi: 10.1111/j.1476-5381.2010.01127.x.
- [4] G. M. Ashraf *et al.*, "Protein misfolding and aggregation in Alzheimer's disease and type 2 diabetes mellitus," *CNS Neurol. Disord. Drug Targets*, vol. 13, no. 7, pp. 1280-1293, 2014, doi: 10.2174/1871527313666140917095514.
- [5] J. Jumper, R. Evans, A. Pritzel, et al., "Highly accurate protein structure prediction with AlphaFold," *Nature*, vol. 596, pp. 583–589, 2021, doi: 10.1038/s41586-021-03819-2.
- [6] S. C. Pakhrin et al., "Deep learning-based advances in protein structure prediction," *Int. J. Mol. Sci.*, vol. 22, no. 11, p. 5553, May 2021, doi: 10.3390/ijms22115553.
- [7] M. A. Haggerty, "Explained: Neural networks and deep learning," *MIT News*, Apr. 14, 2017. [Online]. Available: <https://news.mit.edu/2017/explained-neural-networks-deep-learning-0414>. [Accessed: Aug. 31, 2024].
- [8] J. H. Williams, "Language modeling from scratch: Part 2," *Towards AI*, Jan. 27, 2021. [Online]. Available: <https://towardsai.net/p/data-science/language-modeling-from-scratch-part-2>. [Accessed: Aug. 31, 2024].
- [9] I. H. Sarker, "Deep learning: A comprehensive overview on techniques, taxonomy, applications and research directions," *SN Comput. Sci.*, vol. 2, no. 6, p. 420, 2021, doi: 10.1007/s42979-021-00815-1.
- [10] "Reading Protein Structure," *Nursing Hero*. [Online]. Available: <https://www.nursinghero.com/study-guides/bio1/reading-protein-structure>. [Accessed: Sep. 01, 2024].

- [11] M. Mirdita, K. Schütze, Y. Moriwaki, L. Heo, S. Ovchinnikov, and M. Steinegger, "ColabFold: making protein folding accessible to all," *Nature Methods*, vol. 19, pp. 679–682, 2022, doi: 10.1038/s41592-022-01488-1.
- [12] V. Hornak, R. Abel, A. Okur, B. Strockbine, A. Roitberg, and C. Simmerling, "Comparison of multiple Amber force fields and development of improved protein backbone parameters," *Proteins*, vol. 65, no. 3, pp. 712–725, Nov. 2006, doi: 10.1002/prot.21123.
- [13] ChatGPT. (GPT-4). OpenAI. Accessed: Sep. 1, 2024. [Online]. Available: <https://chat.openai.com/chat>.
- [14] "ColabFold Documentation," *UMass Unity*. [Online]. Available: <https://docs.unity.rc.umass.edu/documentation/tools/colabfold/>. [Accessed: Sep. 01, 2024].
- [15] ColabFold, "AlphaFold2.ipynb," 2021. [Online]. Available: <https://colab.research.google.com/github/sokrypton/ColabFold/blob/main/AlphaFold2.ipynb>. [Accessed: Sep. 01, 2024].
- [16] ColabFold, GitHub repository, 2021. [Online]. Available: <https://github.com/sokrypton/ColabFold>. [Accessed: Sep. 01, 2024].
- [17] "torch.nn.Dropout," *PyTorch Documentation*. [Online]. Available: <https://pytorch.org/docs/stable/generated/torch.nn.Dropout.html>. [Accessed: Sep. 01, 2024].
- [18] "Table 2.1," *Western Oregon University*, Apr. 2020. [Online]. Available: <https://wou.edu/chemistry/files/2020/04/Table-2.1.jpg>. [Accessed: Sep. 01, 2024].
- [19] Soeding Lab, *hhdatabase_cif70*, GitHub repository, 2021. [Online]. Available: https://github.com/soedinglab/hhdatabase_cif70. [Accessed: Sep. 01, 2024].
- [20] ColabFold, "batch.py," *GitHub repository*, 2021. [Online]. Available: <https://github.com/sokrypton/ColabFold/blob/main/colabfold/batch.py>. [Accessed: Sep. 01, 2024].
- [21] S. Saranyan, "Investigations into AlphaFold Modelling Types," *Google Docs*, [Online]. Available: <https://docs.google.com/document/d/1qSxyBQB7Vwt77xJsdGWz8se-t6XDk7AVuZNLEb4vKSs/edit>. [Accessed: Sep. 01, 2024].

- [22] P. Bryant, G. Pozzati, and A. Elofsson, "Improved prediction of protein-protein interactions using AlphaFold2," *Nat. Commun.*, vol. 13, p. 1265, 2022, doi: 10.1038/s41467-022-28865-w.
- [23] J. Liu et al., "Enhancing AlphaFold-Multimer-based protein complex structure prediction with MULTICOM in CASP15," *Commun. Biol.*, vol. 6, no. 1, p. 1140, Nov. 2023, doi: 10.1038/s42003-023-05525-3.
- [24] "FASTA Format Overview," *Zhang Lab*, [Online]. Available: <https://zhanggroup.org/FASTA/>. [Accessed: Sep. 01, 2024].
- [25] "Strengths and Limitations of AlphaFold," *EMBL-EBI*, [Online]. Available: <https://www.ebi.ac.uk/training/online/courses/alphafold/an-introductory-guide-to-its-strengths-and-limitations/strengths-and-limitations-of-alphafold/>. [Accessed: Sep. 01, 2024].
- [26] "The Complement System," *NCBI Bookshelf*, [Online]. Available: <https://www.ncbi.nlm.nih.gov/books/NBK27100/>. [Accessed: Sep. 01, 2024].
- [27] "National Center for Biotechnology Information (NCBI)," *NCBI*, [Online]. Available: <https://www.ncbi.nlm.nih.gov/>. [Accessed: Sep. 01, 2024].
- [28] "Modified Copy of AlphaFold2.ipynb," *Google Colab*, [Online]. Available: https://colab.research.google.com/drive/1aZmUSm1k3XBO6W18gfu_InFhUJWPx4H-. [Accessed: Sep. 01, 2024].
- [29] "Protein Classes," *Rose-Hulman Institute of Technology*, [Online]. Available: https://www.rose-hulman.edu/~brandt/Chem330/Protein_classes.pdf. [Accessed: Sep. 01, 2024].
- [30] "Research Collaboratory for Structural Bioinformatics (RCSB) Protein Data Bank," *RCSB*, [Online]. Available: <https://www.rcsb.org/>. [Accessed: Sep. 01, 2024].
- [31] R. Maoz-Segal and P. Andrade, "Molecular mimicry and autoimmunity," in *Infection and Autoimmunity*, 2015, pp. 27–44, doi: 10.1016/B978-0-444-63269-2.00054-4.

