

RESEARCH ARTICLE

Do “Newly Born” orphan proteins resemble “Never Born” proteins? A study using three deep learning algorithms

Jing Liu^{1,2} | Rongqing Yuan³ | Wei Shao⁴ | Jitong Wang³ | Israel Silman⁵  | Joel L. Sussman⁶ 

¹Department of Biotechnology and Food Engineering, Guangdong Technion-Israel Institute of Technology, Shantou, China

²Faculty of Biotechnology and Food Engineering, Technion-Israel Institute of Technology, Haifa, Israel

³Department of Chemistry, Tsinghua University, Beijing, China

⁴School of Chemistry and Chemical Engineering, Shanghai Jiao Tong University, Shanghai, China

⁵Department of Brain Sciences, The Weizmann Institute of Science, Rehovot, Israel

⁶Department of Chemical and Structural Biology, The Weizmann Institute of Science, Rehovot, Israel

Correspondence

Israel Silman, Department of Brain Sciences, The Weizmann Institute of Science, Rehovot 7610001, Israel.

Email: israel.silman@weizmann.ac.il

Joel L. Sussman, Department of Chemical and Structural Biology, The Weizmann Institute of Science, Rehovot 7610001, Israel.

Email: joel.sussman@weizmann.ac.il

Present address

Jing Liu, Feinberg Graduate School, Weizmann Institute of Science, Rehovot, Israel.

Funding information

Center for Scientific Excellence at the Weizmann Institute of Science

Abstract

“Newly Born” proteins, devoid of detectable homology to any other proteins, known as orphan proteins, occur in a single species or within a taxonomically restricted gene family. They are generated by the expression of novel open reading frames, and appear throughout evolution. We were curious if three recently developed programs for predicting protein structures, namely, AlphaFold2, RoseTTAFold, and ESMFold, might be of value for comparison of such “Newly Born” proteins to random polypeptides with amino acid content similar to that of native proteins, which have been called “Never Born” proteins. The programs were used to compare the structures of two sets of “Never Born” proteins that had been expressed—Group 1, which had been shown experimentally to possess substantial secondary structure, and Group 3, which had been shown to be intrinsically disordered. Overall, although the models generated were scored as being of low quality, they nevertheless revealed some general principles. Specifically, all four members of Group 1 were predicted to be compact by all three algorithms, in agreement with the experimental data, whereas the members of Group 3 were predicted to be very extended, as would be expected for intrinsically disordered proteins, again consistent with the experimental data. These predicted differences were shown to be statistically significant by comparing their accessible surface areas. The three programs were then used to predict the structures of three orphan proteins whose crystal structures had been solved, two of which display novel folds. Surprisingly, only for the protein which did not have a novel fold, and was taxonomically restricted, rather than being a true orphan, did all three algorithms predict very similar, high-quality structures, closely resembling the crystal structure. Finally, they were used

Jing Liu and Rongqing Yuan contributed equally to this work.

This is an open access article under the terms of the [Creative Commons Attribution](https://creativecommons.org/licenses/by/4.0/) License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

© 2023 The Authors. *Proteins: Structure, Function, and Bioinformatics* published by Wiley Periodicals LLC.

to predict the structures of seven orphan proteins with well-identified biological functions, whose 3D structures are not known. Two proteins, which were predicted to be disordered based on their sequences, are predicted by all three structure algorithms to be extended structures. The other five were predicted to be compact structures with only two exceptions in the case of AlphaFold2. All three prediction algorithms make remarkably similar and high-quality predictions for one large protein, HCO_11565, from a nematode. It is conjectured that this is due to many homologs in the taxonomically restricted family of which it is a member, and to the fact that the Dali server revealed several nonrelated proteins with similar folds. An animated Interactive 3D Complement (I3DC) is available in Proteopedia at <http://proteopedia.org/w/Journal:Proteins:3>

KEYWORDS

intrinsically disordered protein, molten globule, orphan protein, protein structure prediction, taxonomically restricted

1 | INTRODUCTION

The accepted view, until quite recently, has been that protein sequences have evolved so as to incorporate the features required for optimal folding and function.¹ Specific amino acid or oligopeptide patterns appear to yield insights into phylogenetic differences between the three kingdoms: prokaryotes, archaea, and eukaryotes.² Surprisingly, however, evidence has been presented that “from a sequence similarity perspective, *real unrelated* proteins are indistinguishable from random amino acid sequences.”³ This, at first sight, seems counterintuitive, because it might be anticipated that natural sequences would differ from random sequences in their folding characteristics. Indeed, this conclusion has been challenged by using an ad hoc Evolutionary Neural Network Algorithm (ENNA) to assess whether and to what extent natural proteins can be distinguished from random sequences.⁴ The ENNA approach could correctly distinguish natural proteins from random sequences with >94% accuracy. Very recently, a distilled protein language model was used to distinguish natural from random sequences with >92% accuracy.⁵

Two sets of studies have shown that random sequences with native-like amino acid composition can be expressed and that in many cases, the expressed polypeptide chains of these “Never Born” proteins⁶ fold in aqueous solution into compact structures that display resistance to proteolysis,⁷ or have substantial secondary structure elements.⁸ For some earlier studies using randomized sequences see references.^{9–11}

Recently, 2000 random sequences of 100 amino acids each were used to generate 3D models with RoseTTAFold (RTF).¹² These initial models were optimized by Monte Carlo sampling in amino acid sequence space to yield “novel proteins spanning a wide range of sequences and predicted structures”.¹³ 129 of these sequences were then expressed in *Escherichia coli*. Of those expressed, 27 yielded monodisperse species with circular dichroism (CD) spectra consistent with a native structure. Three of the 3D structures were determined, and all three displayed novel folds.

The folding propensity of random sequences is of relevance in the context of the issue of orphan genes and of the proteins that they

express, namely, “Newly Born” proteins, devoid of detectable homology to any other proteins, which occur in a single species, or proteins that occur within a taxonomically restricted gene (TRG) family, TRGPs.^{14–20} Such a possibility was considered to be impossible even by such a distinguished figure as François Jacob.²¹ However, to quote a recent review:

The origin of novel protein-coding genes was once considered so improbable as to be impossible. In the last decade, and especially in the last 5 years, this view has been overturned by extensive evidence from diverse eukaryotic lineages.¹⁷

Both the term orphan gene and the term orphan protein are often loosely used in more than one context. In this study, the term orphan gene refers to a gene for which evidence has been presented that it has arisen from what was previously a noncoding DNA sequence, and is expressed as an open reading frame (ORF). Thus, the orphan protein for which it codes is seen only in a single species or in one that is closely related taxonomically, that is, a protein generated by a TRG, a TRGP. The general contention is thus that new genes may appear out of previously noncoding genomic regions, a process also known as exonization,²² and code for novel protein sequences.²³ The question that then arises is how new functional protein domains might evolve out of such random sequences?¹⁴ It is fair to say that this is still an open question and that much more experimental data will be required. Of particular interest are studies of higher primates in which novel genes were identified that are shared by chimpanzees, gorillas, and humans, whereas other *de novo* genes may, for example, be restricted to humans.^{24–29} It should also be mentioned that it has been suggested that novel protein sequences may also be generated by ORFs present in long noncoding RNAs.^{15,30}

Very recently, the field of structural biology has undergone a revolution due to the development of the deep-learning-based protein structure prediction programs, AlphaFold2 (AF2),³¹ RTF,¹² and most recently, Evolutionary Scale Modeling (ESM-2).³² All three algorithms have been shown to predict 3D structures for many natural

sequences closely resembling the experimental structures deposited in the PDB.^{33–35} We were curious as to whether these powerful new protein structure prediction tools might be valuable for distinguishing natural from random sequences. Here, we use all three AI/Deep Learning programs just mentioned to predict the structures of a natural sequence and of “Newly Born” orphan proteins. We also used them to predict the structures of the random sequences of “Never Born” proteins expressed by Tretyachenko et al.,⁸ some of which these authors had shown to fold into compact structures, others to belong to the category of intrinsically disordered proteins (IDPs).³⁶

2 | MATERIALS AND METHODS

2.1 | Protein sequences

Protein sequences for the crystal structures were retrieved from the PDB (<https://www.rcsb.org>), for the IDPs and for the “Newly Born” orphan proteins from UniProt (<https://www.uniprot.org>), and for the 10 “Never Born” proteins from the supplementary information associated with the study of Tretyachenko et al.⁸ (https://static-content.springer.com/esm/art%3A10.1038%2Fs41598-017-15635-8/MediaObjects/41598_2017_15635_MOESM1_ESM.pdf). All these sequences are listed in Table 1.

2.2 | AF2 predictions

The AF2 predictions were performed by the AlphaFold2_advanced Colab⁵¹ at https://colab.research.google.com/github/sokrypton/ColabFold/blob/main/beta/AlphaFold2_advanced.ipynb. The defaults used were:

- Multisequence alignment, mmseqs2
- Template protein structures were not used.

2.3 | Evolutionary scale modeling

The ESM-2 predictions were performed via the *esmfold.py* plug-in to PyMol; [<https://github.com/JinyuanSun/PymolFold>] for sequences containing up to 400 amino acids, and for longer sequences by ESM-Fold.ipynb Colab [<https://colab.research.google.com/github/sokrypton/ColabFold/blob/main/ESMFold.ipynb>]. The defaults for this server were employed.

2.4 | RoseTTAFold predictions

The RTF predictions were performed by the Robetta server using the RTF option¹² at: <https://rosetta.bakerlab.org>. The defaults for this server were employed.

2.5 | Natural protein structure

As a control, we selected for this study one native globular protein whose crystal structure has been experimentally determined to high resolution, 1.36 Å, human carbonic anhydrase (UniProt–P00918, PDB–6pea). The sequence information for this protein was obtained from the UniProt database (<https://www.uniprot.org>) and the crystal structure from the PDB [<https://www.rcsb.org>].

2.6 | Structure prediction and comparison

For each of the sequences submitted to the AF2 and RTF servers, the five most probable models were generated, whereas ESM-2 generated only a single model. In the cases in which the 3D structures were available, the predicted structures were aligned with the experimental structures using PyMol [PyMOL Molecular Graphics System, Version 2.1 ATI-4.8.101, Schrödinger, LLC].

AF2 and ESM-2 produce an estimate of the confidence, on a scale of 0–100, for each residue. This confidence measure is called predicted Local Distance Difference Test (pLDDT), as defined on the IDDT-Cα metric.⁵² It is stored in the B-factor fields of the corresponding PDB files. pLDDT is also used to color-code the residues of the model in the 3D structure viewer (Figure 1).

For models predicted by RTF, the Root Mean Square Deviation (RMSD) values are inserted in place of B-factors in the PDB files generated. They are then converted into pLDDT values using a Python program that we wrote, based on the formulae described by,⁵³ as seen at https://phenix-online.org/version_docs/dev-4380/reference/process_predicted_model.html.

$$\text{RMSD} = 1.5 \exp(4(0.7 - \text{LDDT})) \quad (1)$$

$$\text{LDDT} = 100((0.7 - (\ln(\text{RMSD}) - \ln(1.5))/4)) \quad (2)$$

To render the coloring schemes consistent for all structures shown, the B-factors in the original PDB files for experimentally determined 3D structures were converted to pLDDT using the following equations:

$$B = (\text{rmsd}^2)((8(\pi^2))/3.0) \quad (3)$$

$$\text{rmsd} = \sqrt{(3B/(8\pi^2))} \quad (4)$$

$$\text{LDDT} = 100((0.7 - (\ln(\sqrt{(3B/(8\pi^2))}) - \ln(1.5))/4)) \quad (5)$$

For all structures, we then applied the color scheme used in the PDB AlphaFold2 database. [<https://alphafold.ebi.ac.uk>] shown in Figure 1, via the *color_plddt.py* plugin for PyMol [<https://github.com/JinyuanSun/PymolFold>].

TABLE 1 Sequences of proteins and polypeptides utilized.

Category/protein	AAs	Amino acid sequence
Crystal structure		
Human carbonic anhydrase ³⁷	260	>6PEA Carbonic anhydrase 2 <i>Homo sapiens</i> MSHHWGYGKHNGPEHWHKDFPIAKGERQSPVDIDTHTAKYDPSLKPLSVSYDQATSLRIL NNGHAFNVEFDDSDQKAVLKGGPLDGTYRLIQFHHFWGSLDQGGSEHTVDKKKYAAELHL VHWNTKYGDFGKAVQPPDGLAVLGIFLKVGSAPGLQKVVDVLDSIKTKGKSADFTNFDP RGLLPESLDYWTYPGSLTTPPLECVTWIVLKEPISVSSEQVLKFRKLNFNNGEGEPEELM VDNWRPAQPLKNRQIKASFK
IDPs		
<i>Drosophila</i> gliotactin cytoplasmic domain ³⁸	207	>Q7KT70 Gli-Cyt C-Term RNAKRQSDRFYDEDVFINGEGLEPEQDTRGVDNAHMTNHHALRSRDNIYEYRDSPTKT LASKAHTDTSLSRSPSLAMTQKSSSQASLKSGISLKETNGHLVKQSERAAATPRSQNGS IAKVASPPVEEKRLQLPSSTPVTQLQAEPKRVPTAASVSGSSRSTTPVPSARSTTTHT TTATLSSQPAAPRRTHLVEGVPQTSV
Human CDN1C-Cyclin-dependent kinase inhibitor ³⁹	316	>spP49918 CDN1C_Human Cyclin-dependent kinase inhibitor 1 MSDASLRSTSTMERLVARGTFPVLVRTSACRSIFGPDHEELSRELQARLAELNAEDQNR WDYDFQDMPLRGPGRQLQWTEVDSVPAFYRETVQVGRCLLLAPRPVAVAVAVSPPLE PAAESLDGLEEAPEQLPSVPVPAPASTPPVPVLAPAPAPAPAPVAAPVAAPVAVAVLAP APAPAPAPAPAPVAAPAPAPAPAPAPAPAPAPDAAPQESAEGANQGQGRGQEPLAD QLHSGISGRPAAGTAAASANGAAIKKLSGPLISDFFAKRKRSAPEKSSGDVPAPCPSPSA APGVGSVEQTPRKRLR
Human osteopontin ⁴⁰	314	>spP10451 OSTP_Human Osteopontin MRIAVICFLLGITCAIPVKQADSGSSEKQLYNKYPDVATWLNPDPSQKQNLAPQNA VSSEETNDFKQETLPSKSNESHDMDDMDDEDDDDHVDSDSDSDSDSDVDVDDTDDSHQS DESHHSDESDELVTDFPTDLATEVFTPVVPTVDTYDGRGDSVYGLRSKSKKFRRPDIQ YPDATDEITSHMESEELNGAYKAIPVAQDLNAPSDWDSRGKDSYETSQDQSAETHSH KQSRLYKRKANDESHSDVIDSQELSKVSREFHSHEFHSHEDMLVVDPKSKEEDKHLKF RISHELDSASSEVN
“Never Born” proteins—Group 1⁸		
1856	109	>1856 MYQIEKADFTFDVRRRTAATDIENHAFNMVWLQSWCDVSIKRTLDAYDEAYDAAFQRLK PAEWAIDWVASIQRRRRHYVAYNLSKIKLPVRLEKLSGTTLEHHHHHHH
6387	109	>6387 MIEHCYSKTVYYNLEQEKYDLEVTHIEGWMRAGRKDLADNLLDSGHVFIPEVALQENHY REVHAKIGDAEMRVYKRELFPQIVEVLETPSQLFFAEIELEHHHHHHH
4090	109	>4090 MVERDKPPNIWVYDAELQGGIVVWHLAALYCANVDDYAPQDHLDTMYGFDHQKTNIL SFEDESVAQSYWQYGIIFVKSHWGEDLQGAIVESWRDSRLEHHHHHHH
2298	109	>2298 MKWYGRGREDFGSPDVDVEKNCEGEVIYGTSQLYSNVVFDWWAGISEQPTIFIGSLTTP NTKDDMLWYRNDAKNPGHSILYNLINDYWEATEVSGIGNVVLEHHHHHHH
“Never Born” proteins—Group 3⁸		
665	109	>665 MATKGADHGLAAPQPHAKWDTQIPAEGADREHRSGGGNERRFYNEGAKHAQATWAIPEP AFHLQPAVGEGATTDQAGSLEDQWVRSNNNDVDPTQADETLEHHHHHHH
3703	109	>3703 MSLYKFGQRRRAVDPLPRQCQRDKDYDAFIGGEQNCNLSKSFIVVMSVFLYDPTYNVD SEAQDNKLDHHGSEPTHTGDTPTTSEDTRPGSDRVMRDVPQTLEHHHHHHH
933	109	>933 MVREIDDKTISDYLRGADEGTTAYSLKIPTDKCLFAPTKKHLHGDKSQEADPPTKSPM VEHQFGHEPDPFSCREPDYPGSPLVELTGLNRLTQEPNEELEHHHHHHH
6851	109	>6851 MLGIVEETEHTGREGERDLDKSLQFSDGLFMEVTVQANKLGATKASTCTKEDGPSRCRS QHTITNLIEFSDSVDKTEKTDKEGRTGPMNSELVGGQEDEDLEHHHHHHH
9927	109	>9927 MHPAEVSFSGGAPNNESKWDNRHYVQAESGEDHETGVHLGDFDEYLSLQAGPRLPMPELS GHWGSQLCNDCKGGKKKANIVSPEDVSKDVKSVEGFADRLLEHHHHHHH

TABLE 1 (Continued)

Category/protein	AAs	Amino acid sequence
9693	109	>9693 MMNGERSLIPDSMKISAVRLIICGLIKPATAGLKEVDMHVWPNPTSLAGHSVSLYSSGK QISNLAFGDEESPNRERETAPADEDVIPDAHDTSDSLDGEHLEHHHHHH
Orphan protein crystal structures		
1o22 ⁴¹	170	>1o22 orphan protein TM0875 <i>Thermotoga maritima</i> MGSDKIHIIHHHMLRMDILEILYKKGKEFGILEKKMKKEIFNETGVSLEPVNSELIGRIF LKISVLEEGEEVPSFAIKALTPKENAVDLPLGDWTDLKNVFEEDIDYDSYGDMLKILSEK NWKYKIYVPYSSVKKNNRNLVEEFMKYFFESKGWNPGEYTFVSQVEIDNLF
1mw5 ⁴²	187	>1mw5 HYPOTHETICAL PROTEIN HI1480 <i>Haemophilus influenzae</i> GSHMSETDLLMKMVRQPVKLYSVATLFHEFSEVITKLEHSVQKEPTSLSEENWHKQFLK FAQALPAHGSASWLNLDALQAVVGNRSRSLHQLIAKLKSRHLQVLELNKIGSEPLDLS NLPAPFYVLLPESFAARITLLVQDKALPYVRVSMYWHALEYKGELENDPAANKARKEAL AAATAEQ
3ut8 ⁴³	130	>3ut8 Uncharacterized protein <i>Clostridium thermocellum</i> MTSLRDLIPKHKFDNSTIDQLCKLIDNEIPIIFDLLKWLQDYNWPIAKDILPVVVLHQS IAMPHILTLQGNIDMWKYWVVKLMIPYLIYPNKQLVKSELERLSLEINEDIREIVNL SKDYLHFYYP
“Newly Born” proteins—well characterized		
HCO_011565 ⁴⁴	632	>HCO_011565 from bHaecon-5 strain of MVWLRASIIFAVATLTSGQSPTECGDPAPIAKKDVLTNGTVKLPDSYKISGVISNWSNT THAFTEAANAIEVSTLFSRNDLSQWLAMKNDASQFFYNRTSGICEKSTSLPLAPFELS AISSNLSFSTLLSGLVEFSNKEPGELVDEKVVAGVEGVRVWVSCVNGTNGTNNFQIEVVF AGVWSLKPPSAFNNPLVQSVQISEYGNFSDKSLKSLQSVFEDRYDTVAADESQLFSTP SGTICSGWKEANIPLTNASDPFNVFIEMTDQNKETYKATVYSAKEELVIVSGSKKDGS IFTNESGTPDGAHSVHDFGKGYEYALGHTRCLDLSPNNSADVVLSGTVSMRPLAYI LVAPELKFGNYGQLKTDNRTVNVFRTFDNKTGDVIELHFDGNWLEKYMFTKLTGDRPSL ASYSRYSQSTPMRASQYNELIRACFAKSSKVHNDNNTFILDVKSRSVENVYVSGVETVSS ALAKALSQIAPINPHRVRFYSSGPDLSRVFVSDEKTDKEPSIVPKYNFSAEVSTDEF MQKLNETISKGDWKFVSVADEKTEDWIVAARSLRYAPSSPPSTKYAGYGGAMFVLGV FSLLGVAIGAGGVFFVTQRQRISTLAYQVFE
PBOV1 ⁴⁵	135	>spQ9GZY1 PBOV1_HUMAN Prostate and breast cancer overexpressed MRAFLRNQKYEDMHNIIHILQIRKLRLHLSNFPRLPGILAPETVLLPFCYKVFRRKKEVK RSQKATEFIDYSIEQSHHAILTPLQTHLTMKGSSMKCSSLSSEAILFTLTQLTQLGLE CCLLYLSKTIHPQII
FLJ33706 ⁴⁶	194	>Q8NBC4 CT203_human uncharacterized protein C20orf203 MFPRPVLNSRAQAILLPQPPNMLDHRQWPPRLASFPTKTGMLSRATSVLAGLTAHLWDL GGGAGRRTSKAQRVHPQPSHQRPQPPQHPGYPYQERIWVGEGWGEVGLRLSKVGRDR EVGRGLRAPAGRGRAMGMPRMGTVDGFGQALSSLAWTSTCFQDFCLPSLPGLPAPLIS KQQLSNSRSRLFN
NCYM ⁴⁷	109	>spP40205 NCYM_human cis-antisense gene of the MCYN oncogen MQHPPCEPGNCLSLKEKKITEGSGGVCWGGETDASNPAPALTACCAAEANVEQGLAGR LLLCNYERRVYRRCKIAGRGRAPLGRPLDVSSFKLKEEGRPPCLKINK
Gm13030 ⁴⁸	143	>trA2APQ6 A2APQ6_mouse regulating the pregnancy cycle MCRFHLLQAIKPPEKQMEQSSALGSIMKLSQSHATETTWWLPSQGLRDYLLHPACFHIF RKEGRPDRCRANMIYGFDKTHPRRCCTDLLFQPRLLMLSRVLGPEQLQELLQIPDDLTSP SLSYGSNQNLSQLNFPKHVHTG
TaFROG ⁴⁹	130	>trA0A0K1YY56 A0A0K1YY56_Wheat Fusarium resistance orphan MVWSTSKQGGEREESQHKMVKVKTPIFTQSLFHSPLNKNVKNIEVDRLRLSFTTPK NSTLVPVDSGSDEESDEDRGCSIDSNKPMDEGLDHICSLHAIPRKNKARSACKRSHKI SSRKFKYKIFS
Newtic1 ⁵⁰	375	>Newtic1 MAGSLEPAMATTSYSAFLWALMSTASMVSTVALLLCLCRRRLRKQSGHSISARNSGLTR AQTTLTQVTKNEELPMNGTKTVEGIENGAFSASEPVDPEQVQLPVDVQTSRRSPDSK RASQNKTLQGPQESLKHRLPSIPHIMLNQPDNERNSGVPMQRPPTFPPEPPTYEEVHGG KAAVESGPPERAVGGGLSQRVGPMEHESWIGDDEEWPAAPPLTFDQPFENVYTDLNESI ATLKDQPIIVPTVDSAGTQKSSPQATASVGVEPHEENSTSIAPLSWRFHLALASNTEDN TAVQDPTLTGLYSKVKKSTKRPFLPSSTLEDPLKVEDDVPPPVPVKQFDIEEDLSTQNQ DIEDPLPPPILLIN

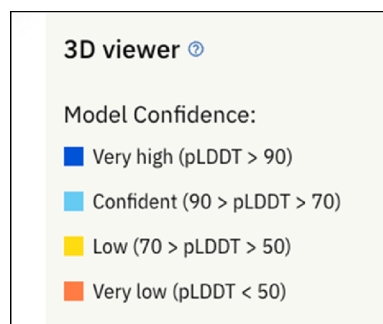


FIGURE 1 Colorcoding scheme based on the PDB AlphaFold database [<https://alphafold.ebi.ac.uk/>].

2.7 | Detection of novel and unique folds

The Dali server (<http://www2.ebi.ac.uk/dali>) was used to check if the protein folds were novel.^{54,55}

2.8 | Calculation of accessible surface area

The accessible surface area (ASA) for 3D protein structures was calculated using the PyMOL Molecular Graphics System, Version 2.1 ATI-4.8.101 Schrödinger, LLC.

2.9 | Morphs for the five top models

To facilitate comparison of each set of the five top models generated by AF2 and RTF, a morph was generated via PyMOL after first aligning the five top model structures on top of each other using PyMOL. These morphs are displayed under Supplementary information S1.

2.10 | Prediction of intrinsically disordered regions

The prediction of intrinsically disordered regions was performed using FoldIndex⁵⁶ (<https://fold.proteopedia.org/cgi-bin/findex>) and fIDPnn⁵⁷ (<http://biomine.cs.vcu.edu/servers/fIDPnn>).

2.11 | Amino acid compositions, pIs, and charge calculations

Amino acid compositions were calculated using the ExPASy ProtParam tool [<https://web.expasy.org/protparam>]. The pI values, and the protein charges at pH 7.4, were calculated using the Prot pi | Protein Tool [<https://www.protpi.ch/Calculator/ProteinTool>].

2.12 | BLASTP sequence searches

BLASTP sequence searches were performed using the NIH-NLM site [<https://blast.ncbi.nlm.nih.gov/Blast.cgi?PAGE=Proteins>].

3 | RESULTS

The AI/Deep Learning tools that have recently been developed have revolutionized the prediction of 3D protein structures.^{31,32,53} This is exemplified by Figure 2, which displays the high-resolution (1.63 Å) crystal structure of human carbonic anhydrase (PDB 6pea), together with the structures predicted by RTF, ESM-2, and AF2. It is immediately apparent that all three algorithms predict this structure very well, displaying high pLDDT scores (Table 2). The RMSD values, relative to the crystal structure, are 0.76 Å, 0.42 Å, and 0.32 Å for RTF, ESM-2, and AF2, respectively. It should be noted that this protein possesses a commonly occurring fold, and has many sequence homologs.

It is now well established that many native proteins are either partially or completely disordered, and are known as IDPs.³⁶ In a review by Uversky⁵⁸ it is stated that "...eukaryotes typically have a higher disorder score than either archaea or prokaryotes, since 52%–67% of eukaryotic proteins have long IDP regions (IDPRs; ≥30 residues) as compared to 26%–51% and 16%–45% proteins with such long IDPRs in archaea and bacteria, respectively.^{59,60}

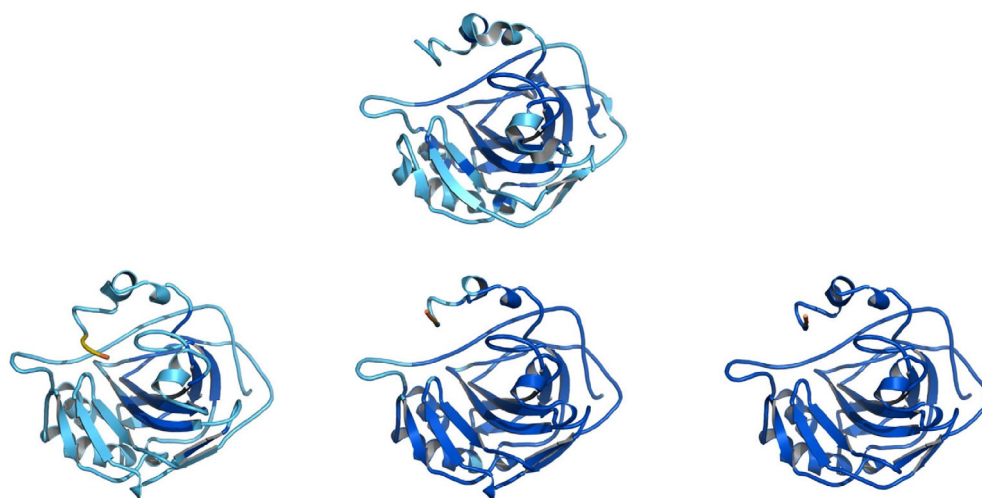
Several groups have examined how well AF2 can be used to predict IDPs.^{61–63} We were also curious as to how well RTF and ESM-2 would predict their structures. Figure 3 shows the predictions for three such proteins, the cytoplasmic domain of gliotactin, the ChE-like adhesion molecule from *Drosophila*,³⁸ human CDN1C-cyclin-dependent kinase inhibitor,³⁹ and human osteopontin.⁴⁰ RTF, ESM-2, and AF2 predict highly unfolded structures for all three proteins. For each individual protein, the models generated by the three algorithms are very different. This is, perhaps, not surprising, since IDPs can adopt large repertoires of conformations. In several of the models, substantial α-helical stretches are predicted, with their percentage in the RTF models being significantly higher than in those generated by the other two methods. They do not appear to correspond to real helical stretches since physicochemical data for the three proteins do not support their presence.

As mentioned in the Introduction, Tretyachenko et al.⁸ predicted, by use of bioinformatic tools, that some of the sequences that they subsequently expressed would be ordered/folded, with a high content of secondary structure elements, whereas others would be disordered/unfolded (IDPs), with a low content of secondary structure elements. They found that only low-significant matches were found by BLAST for the whole set of random sequences.⁸

On the basis of their bioinformatic analysis, they selected groups of 15 proteins each from the random sequence library. Members of Group 1 were predicted to have high secondary structure, low disorder, and high solubility. Members of Group 3 were predicted to have low secondary structure, high disorder, and high solubility (Figure 4). They showed that 31% of the proteins in Group 1 expressed in soluble form, and used CD to demonstrate that they had substantial secondary structure. All the proteins in Group 3 expressed in soluble form, and had low secondary structure.

We used the three algorithms to model the structures of several members of Group 1 (Figure 5) and Group 3 (Figure 6). For all four members of Group 1, RTF predicts compact structures with a high percentage of secondary structure. The structures predicted by

FIGURE 2 Crystal structure of human carbonic anhydrase, PDB 6pea (top), and structures predicted by RoseTTAFold (left), Evolutionary Scale Modeling (center), and AlphaFold2 (right).



ESM-2 and AF2 are more open, with a higher percentage of disordered stretches. In contrast, for Group 3, AF2 and ESM-2 predict, with one exception, very open structures, characteristic of IDPs. RTF predicts some of the structures to be open and others to be more compact. For each individual member of Group 3, the models generated by the three algorithms are very different, as was the case for the well-studied native IDPs referred to in the previous section.

Figure 7 displays the accessible surface areas (ASAs) for the predicted models of proteins in Group 1 (red) and Group 3 (blue), using RTF, ESM-2, and AF2. A student's t-test showed that all three predicted statistically significant differences in ASA between Group 1 and Group 3, with p values being 0.046, 0.002, and 0.007 for RTF, ESM-2, and AF2, respectively. The larger ASA values for the Group 3 proteins relative to those in Group 1 are consistent with their having more open structures. This quantitatively confirms the clear differences in the models predicted by the algorithms that are displayed in Figures 5 and 6.

Despite the fact that research on orphan proteins is a hot topic, largely due to its evolutionary implications,⁶⁴ we were able to find only three crystal structures of orphan proteins in the PDB. Figure 8A shows the crystal structure of orphan protein TM0875 from *Thermotoga maritima* (PDB 1o22),⁴¹ alongside predictions of its structure by RTF, ESM-2 and AF2. The RMSD values, relative to the crystal structure, are 1.37 Å, 17.72 Å, and 1.62 Å for RTF, ESM-2, and AF2, respectively. The crystal structure showed only 149 amino acids, although the sequence that was used for crystallization consisted of 170 amino acids. The pLDDT scores are well correlated with the RMSD values, with ESM-2 showing the lowest agreement with the x-ray structure. The authors pointed out that this was a novel fold, and application of the Dali^{54,65} server reveals that it still maintains this status based on a much large number of experimental structures in the PDB. However, a BLAST search revealed many *Thermotoga* homologs, with the 15th displaying 57% identity for 90% of the sequence. Thus, this protein is a TRGP rather than a true orphan.

Figure 8B shows the crystal structure of a second orphan protein deposited in the PDB, that of the hypothetical protein HI1480 from

Haemophilus influenzae (PDB 1mw5),⁴² alongside the structures predicted by RTF, ESM-2 and AF2. The RMSD values, relative to the crystal structure, are 1.31 Å, 17.6 Å, and 7.3 Å for RTF, ESM-2, and AF2, respectively. It is important to point out that the crystal structure was able to discern only 162 amino acids, although the sequence used for crystallization consisted of 187 amino acids. The pLDDT scores are well correlated with the RMSD values, with RTF being the only method showing good agreement with the x-ray structure, while AF2 showed some similarity. It is worth noting that, even though the 25-residue disordered region is not seen in the crystal structure, it is predicted to be largely helical by RTF. This is in keeping with our observation that in both authentic IDPs, and in members of Group 3 of the “Never Born” proteins, helical stretches were quite frequently predicted by one or other of the three algorithms. In this case, too, the authors pointed out that this is a novel and unique fold, and application of the Dali server again reveals that it still maintains this status based on the much large number of experimental structures now available, as well as the entire AlphaFold Database. A BLAST search confirmed its orphan status.

Cthe_2751, whose crystal structure has been determined [PDB_ID 3ut8], is a protein from *Clostridium thermocellum* with unknown function, which had been reported to be a singleton.⁴³ However, a BLAST search that we performed revealed many homologs from the genus *Clostridium*. So, in fact, it is not an orphan, but rather a TRGP. Examination of its crystal structure revealed an all α -helical topology similar to those observed for nucleic acid processing proteins⁴³ (Figure 8C).

We then went on to use the three algorithms on orphan proteins and TRGPs for which no experimental structures were available. We did this both in order to see how the predictions of the three algorithms would compare, and whether they would predict novel folds. Although many ORFs have been identified which code for putative orphan proteins, only in a limited number of cases has their association with a well-defined biological activity been established. We have identified seven such proteins for which the necessary sequence data are also available. These proteins are:

TABLE 2 pLDDT scores for the structural models predicted by the three algorithms.

Category/protein	Number of AAs	ID ^a	pLDDT			
			Xtal ^b	RTF ^c	ESM-2	AF2 ^d
Crystal structures						
Human carbonic anhydrase	260	6pea/P00918	86.2	86.7	92.7	97.3
IDPs						
<i>Drosophila</i> gliotactin cytoplasmic domain ^e	207	Q7KT70		11.0	49.7	52.3
Human CDN1C-cyclin-dependent kinase inhibitor	316	P49918		11.0	63.2	57.9
Human osteopontin	314	P10451		5.8	40.2	50.2
“Never Born” proteins—Group 1						
#1856	109			39.2	37.8	48.6
#6387	109			24.8	35.5	36.5
#4090	109			25.5	30.55	42.8
#2298	109			45.3	30.9	48.5
“Never Born” proteins—Group 3						
#665	109			37.4	40.7	58.1
#3703	109			24.9	44.2	53.4
#933	109			25.3	41.1	52.9
#6851				39.4	39.7	45.4
#9927				23.8	40.8	72.8
#9693				37.6	44.9	51.6
Orphan protein crystal structures						
1o22 <i>Thermatoga maritima</i> TM0875—unknown function	170	1o22/Q9WZX8	80.7	68.4	34.4	72.8
1mw5 <i>Hemophilus influenzae</i> hypothetical protein HI1480	187	1mw5/P44209	74.9	67.9	38.9	44.6
3ut8 <i>Clostridium thermocellum</i> —unknown function	130	3ut8/Cthe_2751	77.8	88.2	92.6	97.5
“Newly Born” protein orphans and TRGPs						
HCO_011565 ⁴⁴ from the bHaecon-5 strain of <i>H. contortus</i>	632	HCO_011565		59.8	86.2	83.2
PBOV1 ⁴⁵ —Human tumor-specific gene	135	Q9GZY1		36.4	36.4	46.1
FLJ33706 ⁴⁶ —Human gene expressed in neurons (alt ID C20orf203)	194	Q8NBC4		28.8	31.2	50.5
NCYM ⁴⁷ —Human DNA binding transcriptional activator homolog	109	P40205		26.2	30.7	45.2
Gm13030 ⁴⁸ —Mouse involved in regulating the pregnancy cycle	143	A2APQ6		26.3	33.4	36.0
TaFROG— <i>Triticum aestivum</i> (common wheat)	130	A0A0K1YY56		31.4	40.0	60
Newtic1 ⁴⁹ — <i>Cynops pyrrhogaster</i> (Japanese fire-bellied newt)	375	Newtic1		10.9	53.3	57.4

Abbreviations: AAs, amino acids; AF2, AlphaFold2; ESM-2, Evolutionary Scale Modeling; RTF, RoseTTAFold; TRGPs, taxonomically restricted gene proteins.

^aIDs: a four letter code is a PDB ID (<https://www.rcsb.org>), while the longer ID corresponds to a UniProt accession ID (<https://www.uniprot.org>).

^bCalculated from the PDB entry (<https://www.rcsb.org>).

^cCalculated from the RTF highest-ranked model.

^dCalculated from the AF2-Colab highest ranked model.

^eC-terminal 207 residues of the *Drosophila* gliotactin cytoplasmic domain.

- HCO_011565 from the Haecon-5 strain of the nematode, *Haemonchus contortus*. This is a 632-residue protein that has been expressed and characterized by Taki et al.,⁴⁴ who showed that it was the target of a nematocidal small molecule. They also

modeled it using AF2, which predicted that it has a well-defined 3D structure, with a transmembrane domain at its C-terminus. When we performed BLAST on this protein, we found that the first 74 homologs retrieved were all from nematodes. The 75th was

FIGURE 3 3D structure predictions for three intrinsically disordered proteins using RoseTTAFold (RTF), Evolutionary Scale Modeling (ESM-2), and AlphaFold2 (AF2). (A). Gli-Cyt, RTF (left), ESM (center), AF2 (right); (B). CDN1C, RTF (left), ESM (center), AF2 (right); (C). Osteopontin, RTF (left), ESM (center), AF2 (right).

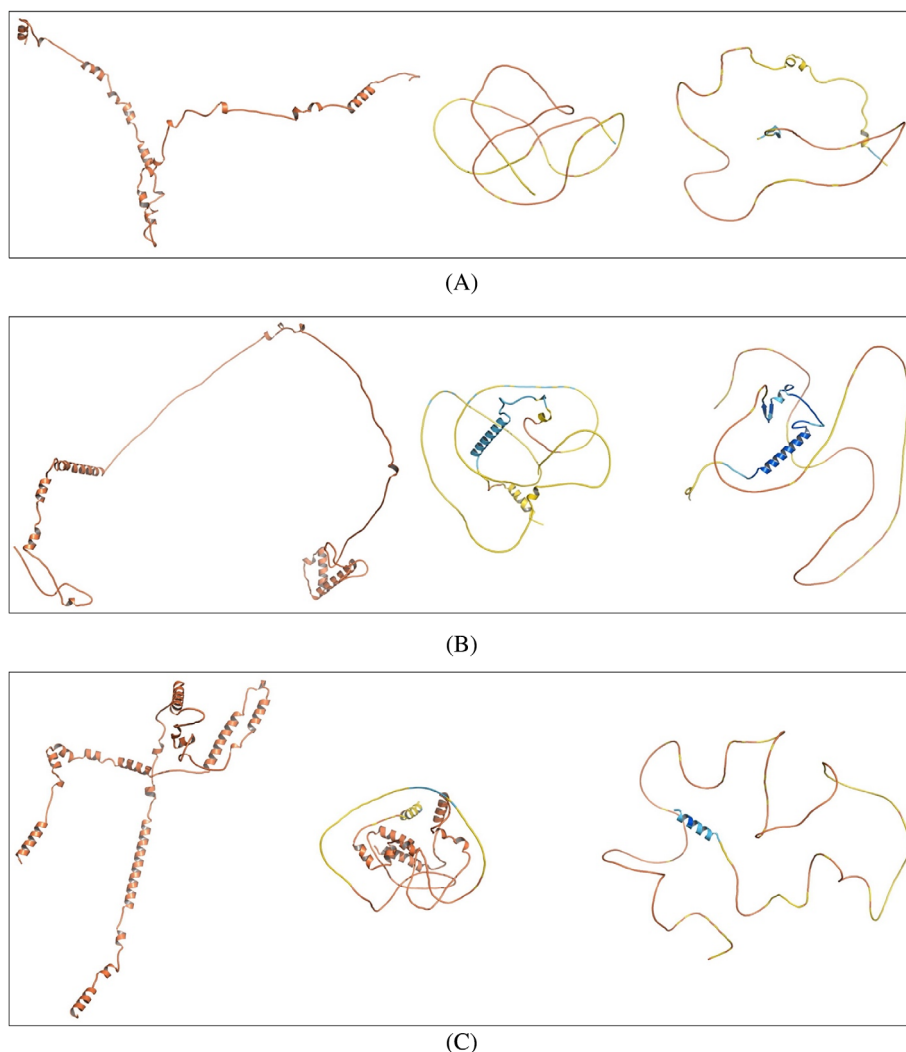
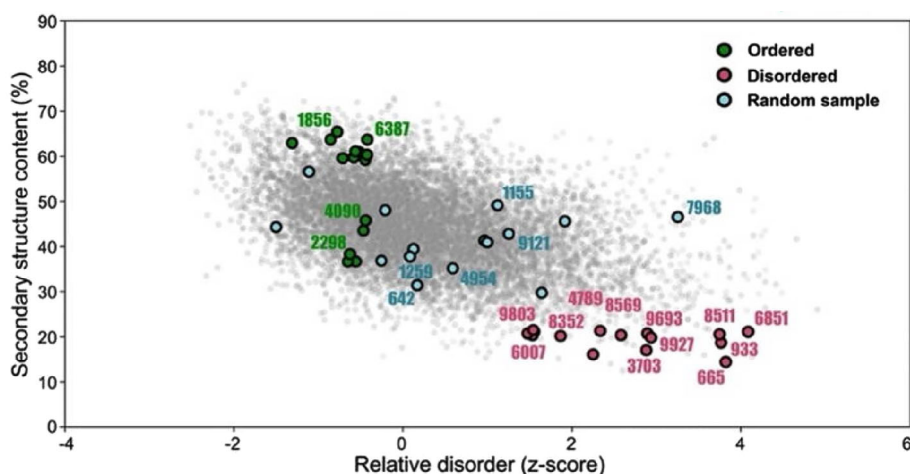


FIGURE 4 Selection of sequences from the set of “Never Born” proteins taken for experimental characterization. Secondary structure is plotted on the y-axis vs. relative disorder on the x-axis. Members of Group 1 (green circles) fall into the category of ordered/folded proteins, and members of Group 3 (red circles) fall into the category of disordered/unfolded proteins. This figure was previously published in Tretyachenko et al.⁸ licensed under a Creative Commons Attribution 4.0 International License <http://creativecommons.org/licenses/by/4.0/> and reproduced here with no change.



from a tick, with the E value increasing from 1×10^{-9} to 3×10^{-6} between the 74th and 75th homologs. Thus, it is clearly a TRGP rather than a true orphan.

- PBOV1 is a human de novo gene with tumor-specific expression that is associated with a positive clinical outcome of cancer.⁴⁵ It is

highly expressed in primary gliomas and breast tumors. PBOV1 codes for a protein that contains 135 residues, whose function is unknown. The authors reported that they were unable to find any orthologs, whether in humans or in any other species. Our BLAST search produced many hits, but all were described as “Low-Quality

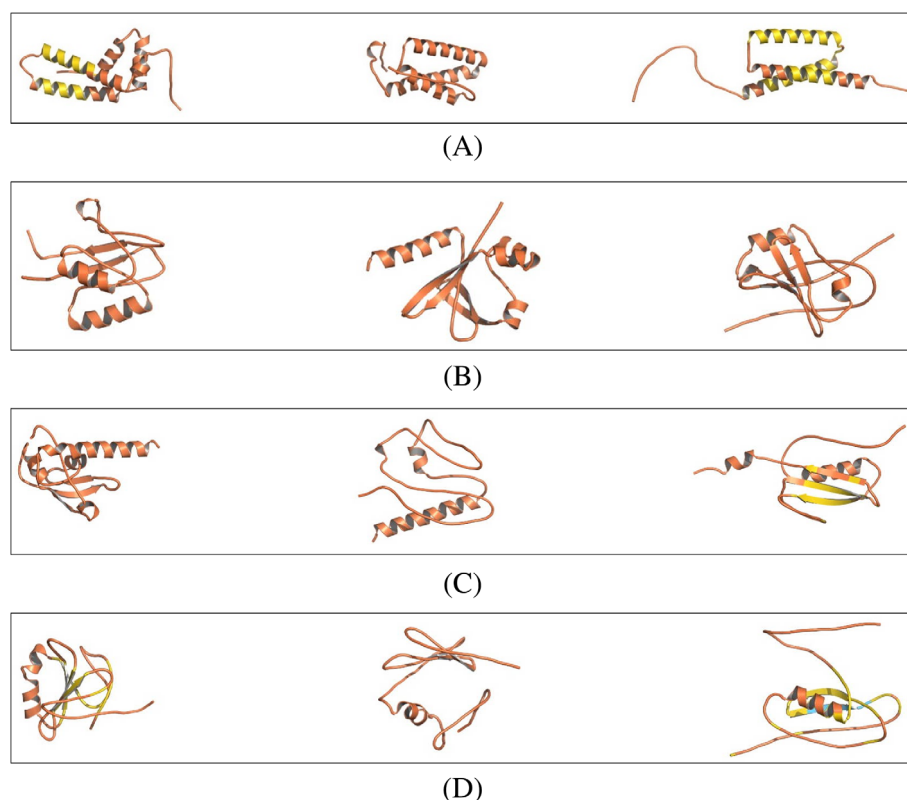


FIGURE 5 3D structure predictions for 4 members of Group 1 of “Never Born” proteins (experimentally shown to be compact, with substantial secondary structure), with RoseTTAFold (left), Evolutionary Scale Modeling (center), and AlphaFold2 (right). (A) #1856; (B) #6387; (C) #4090; (D) #2298.

Proteins”, meaning that it is uncertain if they are, in fact, real proteins. So, the PBOV1 gene product appears to be a true human orphan protein.

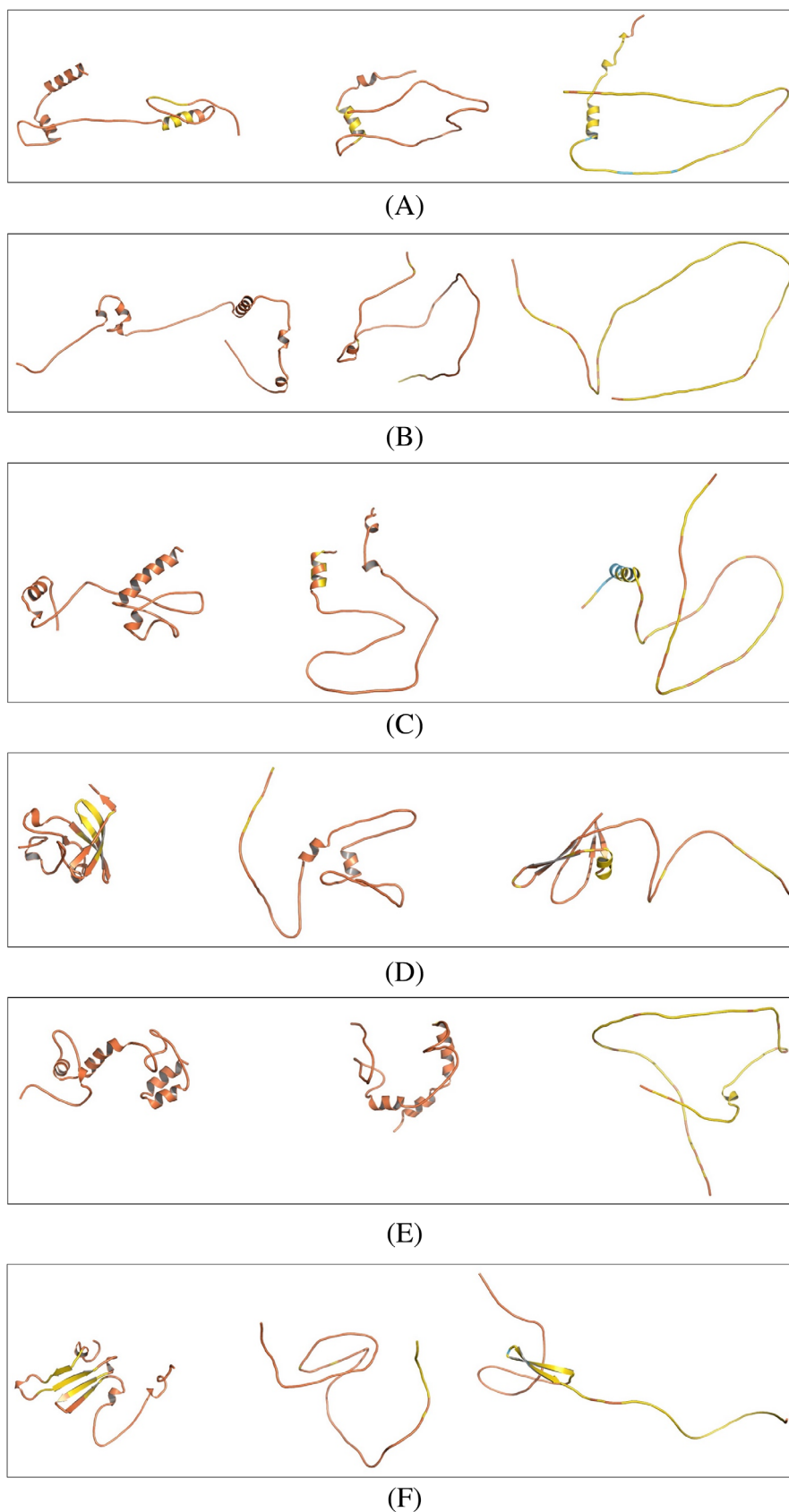
- FLJ33706 (alternative gene ID C20orf203) was identified as a de novo human gene in studies concerned with nicotine addiction.⁴⁶ It is abundantly expressed within neurons in several areas of the human brain, including cortex, cerebellum, and mid-brain, and elevated expression was observed in Alzheimer's disease brain samples. It codes for a protein that contains 194 residues, and the authors identified it as human-specific. Our own BLAST search also revealed no orthologs. As was the case for PBOV1, our BLAST search for FLJ33706 produced many hits, but all were described as “Low-Quality Proteins”, again meaning that it is uncertain if they are, in fact, real proteins. So, the FLJ33706 gene product, too, appears to be a true human orphan protein.
- NCYM, a *cis*-antisense gene of the MYCN oncogene, encodes a de novo evolved protein that regulates the pathogenesis of human cancers, especially neuroblastomas.⁴⁷ The NCYM protein, which contains 109 residues, inhibits the gsk3 β kinase, which, in turn, promotes degradation of the MYCN gene product. So, NCYM is the first de novo evolved protein to act as an oncopromoter in human cancer. For NCYM, our BLAST search revealed many “Low-Quality Homologs” and a few genuine hits, all from monkeys. So, the gene product appears to be a TRGP.
- Gm13030, a protein containing 143 residues, which is encoded by a young protein-coding gene, is specifically expressed in the oviduct of the female mouse.⁴⁸ If the gene expressing it is knocked out, the pregnancy cycle is shortened, and the infanticide rate

increased.⁴⁸ In this case, too, our BLAST search revealed no orthologs. So, the gene product appears to be a true orphan protein.

- TaFROG, which stands for *Triticum aestivum Fusarium* resistance orphan gene, expresses a protein that confers resistance on wheat (*T. aestivum*) to the mycotoxigenic fungus, *Fusarium graminearum*.⁴⁹ The TaFROG protein, which contains 130 residues, and was found to be an IDP, is localized to the nucleus, and acts by interacting with the α subunit of the sucrose non-fermenting-related kinase 1. Our BLAST search revealed two uncharacterized proteins from related wheat, and a few hypothetical proteins. Thus, the TaFROG protein appears not to be a true orphan, but rather a TRGP with a small number of identified homologs.
- Newtic1 is a protein expressed in the regenerating limbs of the adult Japanese fire-bellied newt, *Cynops pyrrhogaster*.⁵⁰ It is found in only one other closely related newt. It is specifically expressed in a subset of erythrocytes that form clumps that accumulate in the distal portion of the regenerating limb, and may be involved in clump formation. It contains 375 amino acids, with a transmembrane sequence near its N-terminus. Our BLAST search did not reveal any additional homologs than the one referred to already. So, the gene product appears to be a true orphan protein.

The number of amino acids for these seven orphans/TMGPs ranges from 109 for NCYM to 632 for HCO_011565, and their pI values range from 4.99 for Newtic1, which contains 7.7% Glu and 5.1% Asp residues, with a net charge at pH 7.4 of -15.67 , to a pI value of 12.23 for FLJ33706, which contains 10.3% Arg and 2.6% Lys residues, with a net charge at pH 7.4 of $+15.25$.

FIGURE 6 3D structure predictions for 6 members of Group 3 of “Never Born” proteins, with RoseTTAFold (left), Evolutionary Scale Modeling (center), and AlphaFold2 (right). (A) #665; (B) No. 3703; (C) No. 933; (D) No. 6851; (E) No. 9927; (F) #9693.



As an initial step in characterizing these seven proteins, we utilized FoldIndex⁵⁶ and fIDPnn⁵⁷ to investigate whether they were predicted to be intrinsically disordered or folded. Although these

algorithms make similar predictions, since fIDPnn was selected as the best disorder predictor in the first Critical Assessment of Protein Intrinsic Disorder Prediction (CAID),⁶⁶ the data displayed in Figure 9

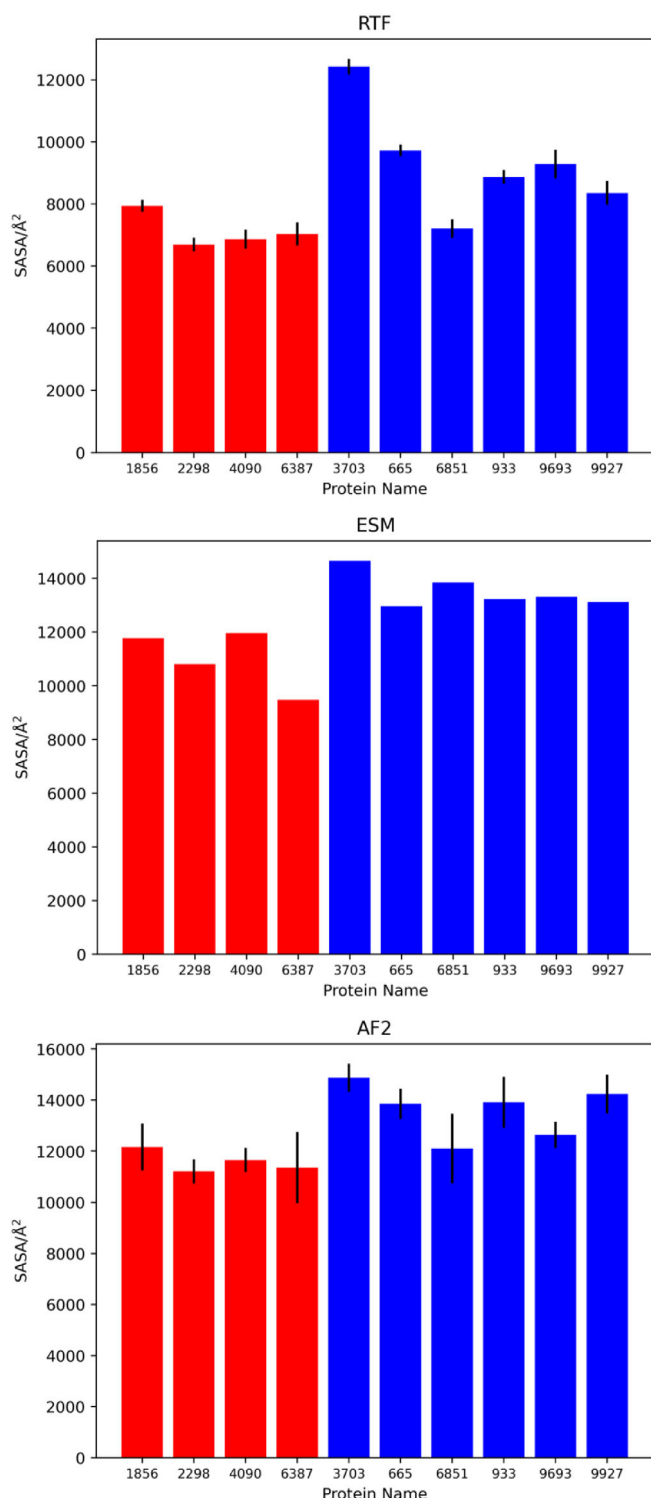


FIGURE 7 Comparison of the accessible surface areas (ASAs) for “Never Born” proteins. ASA values for Group 1 are displayed in red, and those for Group 3 in blue. For RoseTTAFold (RTF) and AlphaFold2 (AF2), the black vertical bars indicate the maximum and minimum values of the five predicted models. For Evolutionary Scale Modeling (ESM-2) no bars are shown, since it predicts only one model. Top panel, RTF; middle panel, ESM-2; bottom panel, AF2.

were generated using it. Five of the proteins are predicted to be almost completely folded, although FLJ33706 has a short, disordered

stretch in the middle of its sequence. The other two, TaFROG and Newtic1, are classified as IDPs, since they are predicted to be disordered throughout almost their entire sequences.

Of the seven proteins studied using the three structure prediction algorithms, only the first, that of HCO_011565, shows fully folded and almost identical structures (Table 3 and Figure 10), with high pLDDT scores (Table 2). Most likely, this is for two reasons. Firstly, rather than being a *true* orphan, HCO_011565 is the product of a TRG,⁴⁴ with the BLAST search having revealed that the first 74 homologous sequences, with the lowest E values, were all from nematodes. Secondly, the DALI server revealed a number of hits for the entire predicted structure, as well as for the three sub-domains predicted by all three algorithms. This is very similar to what was observed when the algorithms were applied to Cthe_2751, whose crystal structure has been solved. It, too, is a TRGP with many homologs, and does not possess a novel fold (Figure 8C). It is interesting that AF2 and ESM-2 give essentially identical structures, although the two algorithms are quite different.

The fIDPnn disorder algorithm predicted that TaFROG⁴⁹ and Newtic1⁵⁰ are mostly unfolded; indeed, all three structure prediction algorithms produce very open structures for both of them, just as was seen for the highly disordered Group 3 “Never Born” proteins (Figure 6). It is interesting that in Newtic1 a substantial stretch near its N-terminus is predicted to be ordered. Casco-Robles et al.⁵⁰ predicted that the sequence F₁₇LWALMSTASMVSTLVALLLCGLC₄₀ is a transmembrane sequence; it is therefore likely to be α -helical, and such an assignment is, indeed, made by all three structure prediction algorithms (Figure 10).

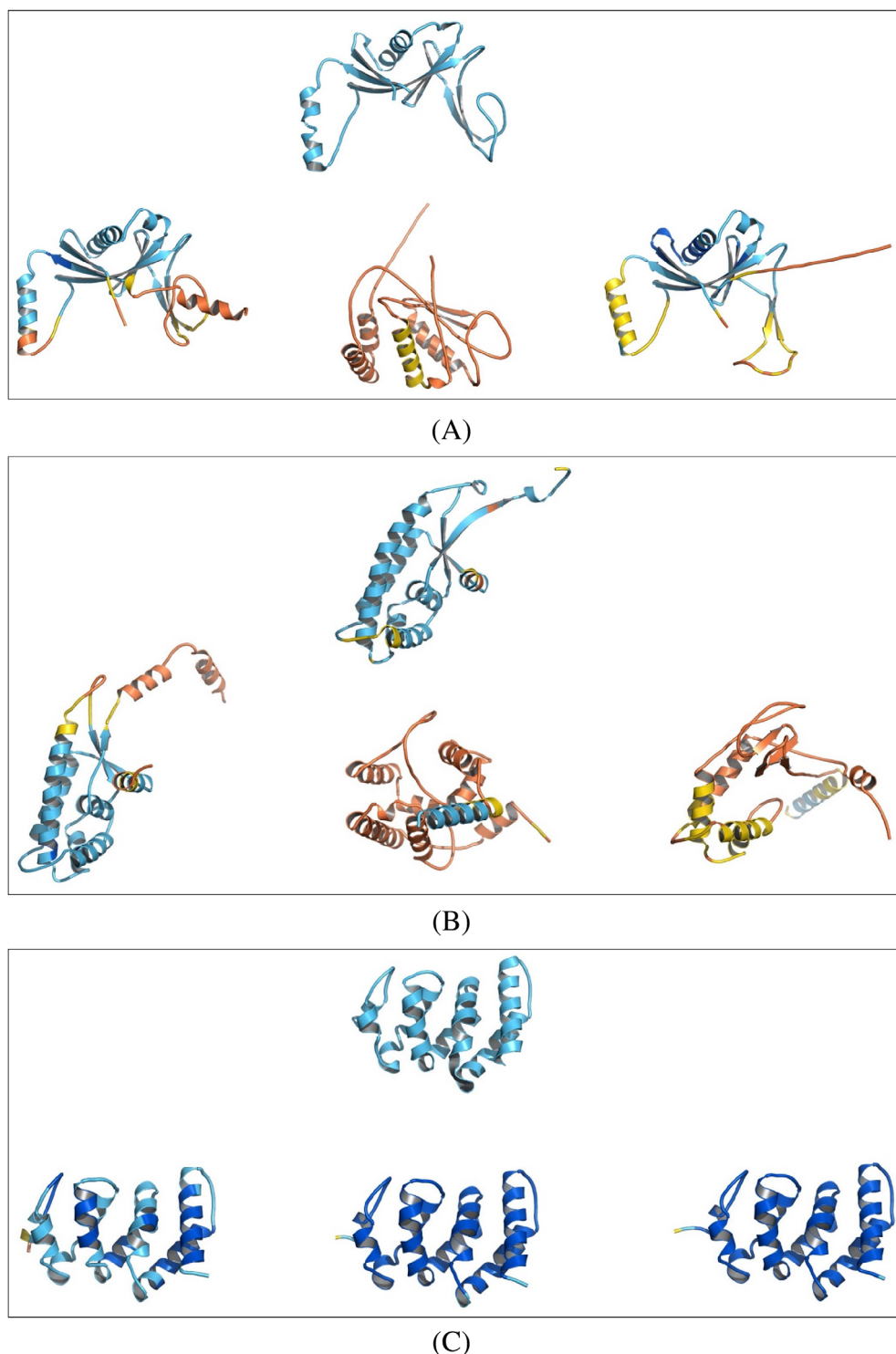
The remaining four proteins^{45–48} are predicted to be predominantly folded by fIDPnn. However, the three structure prediction algorithms yield quite different results for each of them. We should, of course, stress that the pLDDT scores are very low for all these structures, using all three algorithms. However, our present interest is not in their precise 3D structures, but rather in their overall shapes. For all four proteins, RTF yields the most compact structures, with a high percentage of secondary structure (mostly α -helical), while for the structures predicted by ESM-2, there is substantial secondary structure, but the structures are somewhat less compact. In the case of AF2, three of the structures are very open, but that of PBOV1 is quite compact, and is somewhat similar to that predicted by the other two algorithms.

4 | DISCUSSION

The discovery of “Newly Born” orphan proteins,^{23–26} and of proteins coded for by TRGs, namely, TRGPs, which utilize DNA sequences that were previously not expressed, raises cogent questions as to how the polypeptides coded for by such sequences evolve into biologically active proteins.

Are the powerful structure prediction algorithms that have emerged in the past 3 years applicable to the structure analysis of such emergent proteins? In this context, it is worth noting that both

FIGURE 8 Crystal structures of an orphan proteins and of two taxonomically restricted gene proteins (TRGPs), and the models generated by the three structure prediction algorithms. (A) Crystal structure of the TRGP TM0875 from *Thermatoga maritima*, PDB-ID 1o22 (top), and the structures predicted by RoseTTAFold (RTF; left), Evolutionary Scale Modeling (ESM; center), and AlphaFold2 (AF2; right), (B) Crystal structure of the orphan protein hypothetical protein HI1480 from *Haemophilus influenzae*, PDB-ID 1mw5 (top), and structures predicted by RTF (left), ESM (center) and AF2 (right), (C) Crystal structure of the TRGP, Cthe_2751, from *Clostridium thermocellum*, PDB-ID 3ut8 (top), and structures predicted by RTF (left), ESM (center) and AF2 (right).



AF2 and RTF use MSA of homologous proteins as an important element for structure prediction,^{12,31} while ESM-2 does not.³²

True orphan proteins have no sequence homology to any existing protein. We thought, therefore, that the “Never Born” proteins generated and investigated by Tretyachenko *et al.*⁸ would serve as a valuable benchmark for comparison. These authors expressed and purified a substantial number of polypeptides that had compositions similar to those of authentic proteins but whose sequences were random. Based

on CD measurements, they found that a significant number of the random sequences, Group 1, were compact, with substantial secondary structure, while another set, Group 3, appeared to be IDPs. We examined members of these two groups of polypeptides using all three structure prediction algorithms. Although the structures predicted differed significantly in detail, and had low pLDDT scores, by analysis of the ASA data, it was clear that the members of Group 1 were statistically significantly more compact than those of Group

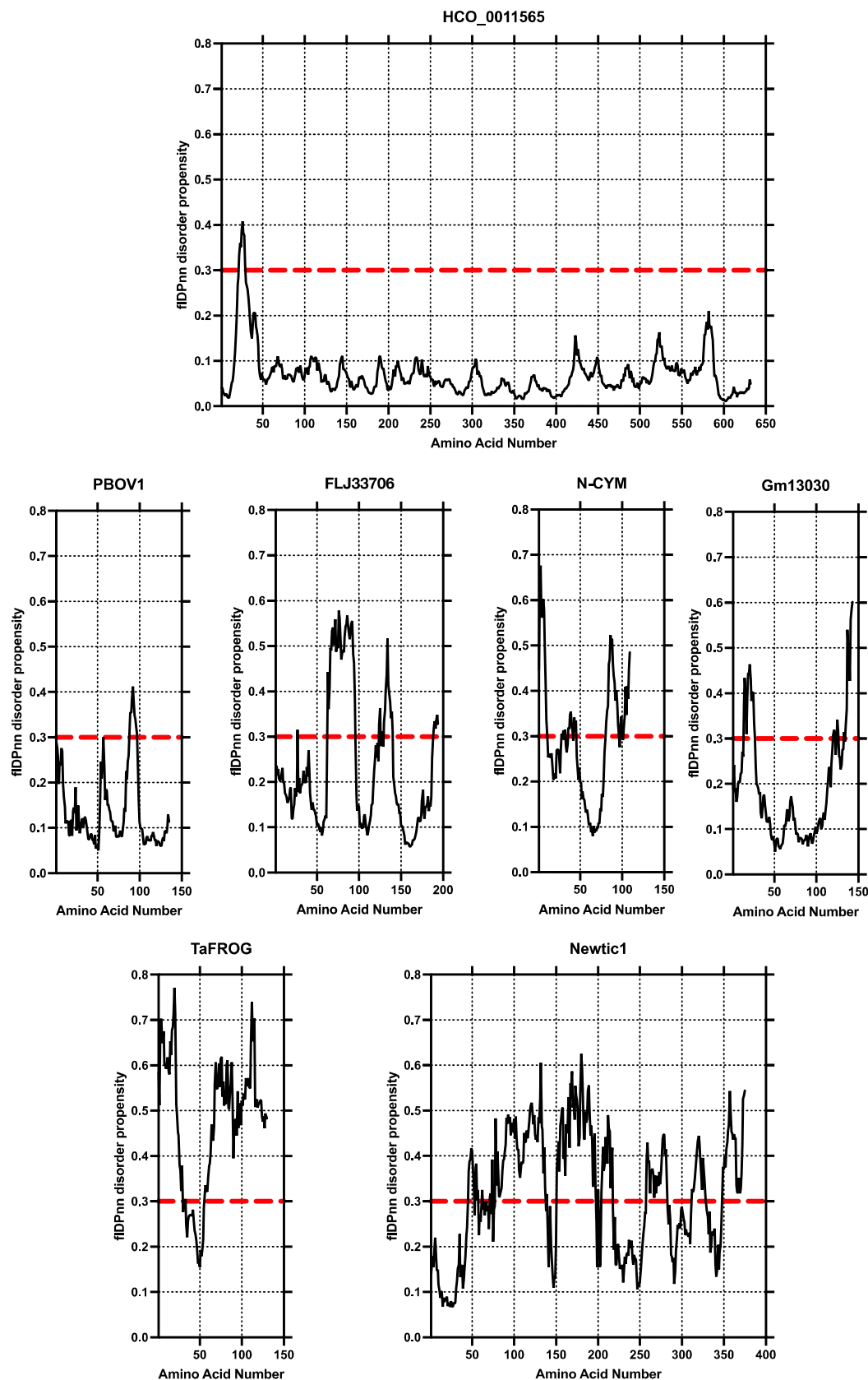


FIGURE 9 . fIDPnn predictions of order/disorder in seven well-studied orphan or taxonomically restricted gene (TRG) proteins. (A) HCO_011565; (B) PBOV1; (C) FLJ33706; (D) NCYM; (E) Gm13030; (F) *TaFROG*; (G) *Newtic1*. Sequences above the red horizontal line are classified as disordered, and below the line as ordered.

TABLE 3 Comparison of the RMSD values for the 3D structures of HCO_011565 predicted by the three algorithms.

Prediction algorithms	RMSD (Å)	Number of atoms used
ESM-2 vs. RTF	3.5	3640
AF2 vs. RTF	3.6	3782
AF2 vs. ESM-2	1.2	3.773

Note: The alignments were performed with PyMOL super command, which automatically determines the number of atoms to be used. Abbreviations: RTF, RoseTTAFold; AF2, AlphaFold2; ESM-2, Evolutionary Scale Modeling.

3 (Figure 7). Moreover, for members of Group 1, both RTF and ESM-2 predict significantly more compact structures than AF2. For Group 3, for which both physicochemical evidence and IDP algorithms indicated that they were intrinsically disordered,⁸ all three algorithms predicted remarkably well that they were disordered, with one exception, #6851 in the case of RTF (Figure 6D). It would be interesting to determine the structures of members of Group 1, to examine directly whether they fold into compact structures.

We similarly analyzed “Newly Born” orphan proteins and TRGPs. Parenthetically, some of the proteins that we examined, which had been defined as orphans in the cited publications, turned out to be TRGPs.

As presented under Results, screening of the PDB revealed only three ‘orphan’ proteins for which crystal structures had been deposited. Only one of these was a true orphan, PDB-ID 1mw5. The other two, PDB-ID 1o22 and PDB-ID 3ut8, were TRGPs. The first two had novel folds, while the third, 3ut8, had a fold that had already been described in proteins that were functionally completely unrelated. It should be noted that the models predicted for this structure by all three algorithms are in excellent agreement with the crystal structure, very much better than the predictions made for the other two orphans/TMGPs displayed and discussed above. This may be attributed to the large number of homologs revealed by BLAST and possibly due to the fact that structures of proteins with similar folds are present in the PDB. Thus, perhaps not surprisingly, not all orphan proteins, or TMGPs, display novel folds.

For the seven orphans, or TRGPs, for which no 3D structure data were available, the disorder predictor, fIDPnn, indicated that five were compact and two were IDPs.

For one of the compact structures, HCO_011565, the three structure algorithms predicted remarkably similar structures with very high pLDDT scores. In contrast, for the other four, all three algorithms predicted structures with very low pLDDT scores. It is plausible that the high quality and similarity of the structures predicted for HCO_011565 is due to the fact that the Dali server identified many structures in the PDB with similar folds for the entire protein and for each of its three domains. In addition, BLAST identified many homologous nematode sequences. The quality of the data is similar to that of the data obtained for 3ut8, for which, in addition to the existence of homologous sequences detected by BLAST, and of proteins from

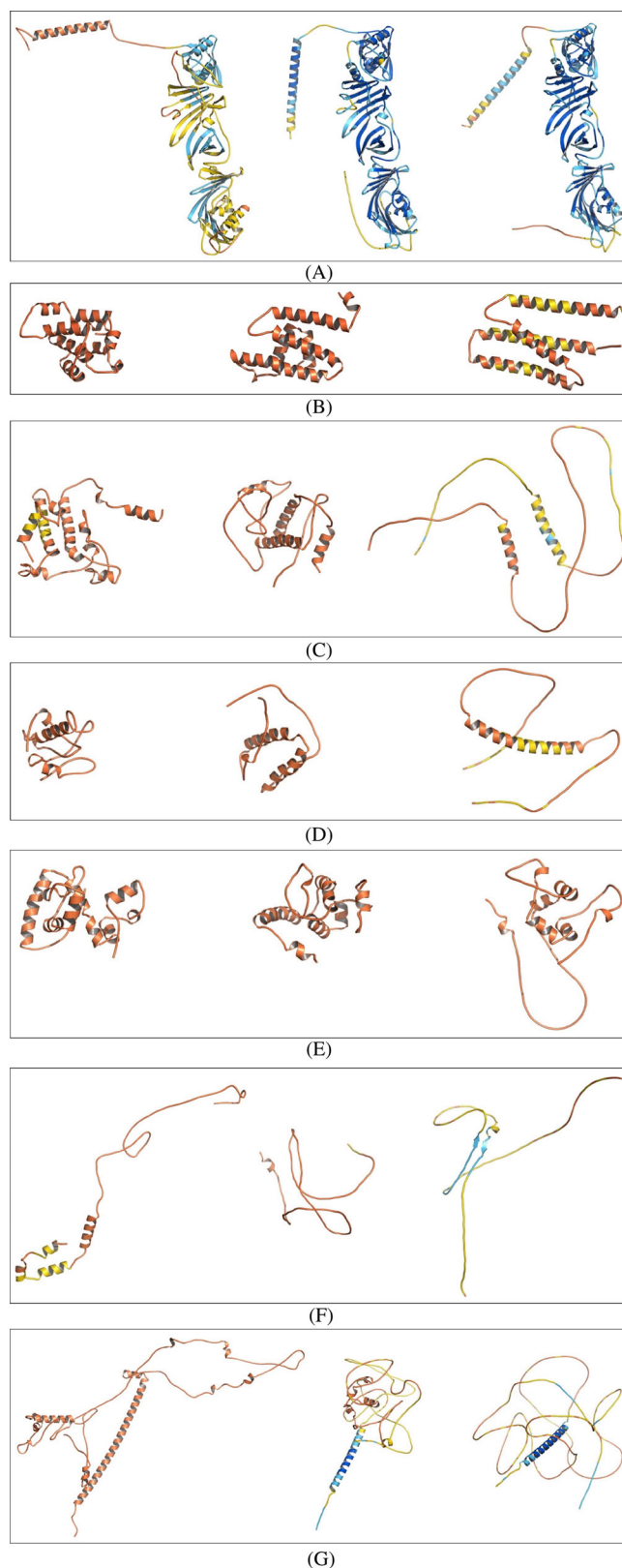


FIGURE 10 3D structure predictions for seven orphan/TMG proteins using RoseTTAFold (RTF; left), Evolutionary Scale Modeling (center), and AlphaFold2 (AF2; right). (A) HCO_011565; (B) PBOV1; (C) FLJ33706; (D) NCYM; (E) Gm13030; (F) TaFROG; (G) Newtic1. Morphs of the five multiple models generated by RTF and AF2 can be visualized in the Supplementary information S1.

other families that displayed the same fold that Dali detected, the crystal structure was available to confirm the predictions.

Monzon et al.⁶⁷ recently examined 250 protein sequence families in the AntiFam resource,⁶⁸ which are thought to be *spurious* proteins. Specifically, they are believed to be ORFs either on the opposite strand or in a different, overlapping reading frame, with respect to the true protein-coding or noncoding RNA gene.⁶⁸ Monzon et al.⁶⁷ conjectured that proteins belonging to these families would not fold into well-folded globular structures. Using AF2, they confirmed this prediction with one exception. To the best of our knowledge, these spurious protein sequences were not examined with RTF or ESM-2. Since some of the data presented above suggest that AF2 may over-predict sequences lacking homologs to be unstructured, it would be worth using these two alternative algorithms to see whether they would predict that some of these sequences might indeed fold into compact structures. Thus, they may not be spurious and possibly have biological functions that could be tested experimentally.

In the Introduction, we referred to a recent study from the Baker lab,¹³ in which 129 random sequences were modeled with RTF. Although these models have low pLDDT scores, and thus most likely differ significantly in detail from the actual structures, they can provide good low-resolution starting points to evolve in the lab into well-folded and biologically active structures.

One conclusion that can be drawn from the data presented above is that some orphan proteins and TRGPs display novel folds that do not overlap with folds already present in the PDB. Indeed, the total number of distinct protein folds has been the topic of heated controversy.^{69,70} In this study, novel folds are observed for two of the three orphan proteins for which crystal structures exist, that is, 1o22 and 1mw5 (Figure 8A,B).

One of the major reservations that have been made with respect to orphan proteins being “Newly Born” proteins, coded for by sequences that were previously noncoding sequences, is that the genes in question might have undergone such rapid evolution that their homology to their predecessors was no longer recognizable.^{71,72} The fact that two of the orphan proteins studied here, whose 3D structures have been experimentally determined, display novel folds substantially weakens this argument.

In retrospect, it is not surprising that many random polypeptide sequences of a suitable amino acid composition, at a first approximation with a high content of hydrophobic residues and a low net charge,⁷³ will yield a compact structure containing substantial secondary structure motifs, as shown by Tretyachenko et al.⁸ This may be due to the fact that the sequences retained the amino acid compositions of the natural proteins from which they were generated, thus not being completely random.

The paradigm change introduced by Kuwajima and by Ptitsyn in the 1980s^{74,75} resulted in the realization that the newly synthesized polypeptide that emerges from the ribosome does not persist as an extended unfolded polypeptide unless it is an IDP, but rather collapses to what is termed a “Molten Globule” (MG), a compact structure somewhat larger than the fully folded native structure. The MG contains substantial secondary structure elements, but lacks the precise

tertiary interactions of the native structure. Small proteins may spontaneously undergo a transition to the native state, whereas larger proteins may require the assistance of molecular chaperones to complete the folding process. The spectroscopic data of Tretyachenko et al.⁸ only tell us that compact structures, with secondary structure elements, have been produced by their random polypeptide sequences. When a native protein unfolds to a MG or some other partially unfolded species, hydrophobic amino acid side chains buried in the hydrophobic core become exposed. The degree of their exposure can be checked by use of the amphiphilic probe, 1-anilinonaphthalene-8-sulfonate (ANS), whose fluorescence is enhanced upon interaction with the hydrophobic residues.⁷⁶ It would, therefore, be interesting to compare the ANS fluorescence of the “Never Born” Group 1 proteins generated by Tretyachenko et al.⁸ to that of typical globular proteins in their native state. It is worth mentioning that in an early study on the folding of polypeptides with random sequences of simplified amino acid composition, NMR data indicated loose packing of the folded state.⁹

In any event, one can speculate that “Newly Born” proteins might, initially, assume an MG-like conformation that would resemble that of the “Never Born” proteins, and that mutations, coupled with natural selection, might convert some of them into “native” orphan proteins with novel biological activities. This may be considered analogous to what occurred in the study in which the hallucinatory proteins were generated.¹³

Why do RTF and ESM-2 do relatively well in predicting plausible compact structures for randomized sequences that have been shown to be compact experimentally, whereas AF2 often makes predictions that are either clearly wrong or implausible? In a recent brief survey of the principles underlying AF2,⁷⁷ it was emphasized that it makes extensive use of MSA for the detection of conserved interactions of residues that are remote from each other in the linear sequence. This approach was earlier proposed, and implemented with a certain degree of success by Marks and Sander.⁷⁸ Obviously, such conserved interactions of distant residues would not exist in randomized sequences. Nor would such conserved interactions be available in orphan proteins, which lack ancestral homologs and, furthermore, in some cases display novel folds. Apparently, even if RTF makes use of such conserved long-distance interactions, it is able to successfully model the overall shape of novel proteins consistently, even in the absence of such information. The fact that ESM-2 is based on natural language, and does not make use of MSA may explain why it, too, does better than AF2.

This study's principal conclusion is that many orphan proteins and TRGPs are predicted to have a compact 3D structure, and, in some cases, a novel fold. It will be interesting to express and purify more of these proteins, so as to determine their experimental structures, in order to find out whether some of them also have novel folds.

Because the sequences of orphan proteins lack homology information, protein structure prediction for them has recently become a hot topic.⁷⁹ The approaches and methodologies that we have implemented in this study may provide a starting point for datasets and protocols to evaluate the performance of structure prediction algorithms on sequences that lack homology to other sequences.

AUTHOR CONTRIBUTIONS

Jing Liu: Investigation; methodology; validation; writing – review and editing; writing – original draft. **Rongqing Yuan:** Investigation; writing – original draft; methodology; validation; writing – review and editing; software. **Wei Shao:** Investigation. **Jitong Wang:** Investigation. **Israel Silman:** Conceptualization; investigation; writing – original draft; methodology; writing – review and editing; formal analysis; supervision; project administration. **Joel L. Sussman:** Conceptualization; investigation; writing – original draft; methodology; validation; visualization; writing – review and editing; software; formal analysis; project administration; supervision; resources.

ACKNOWLEDGMENTS

The Israeli tutors and the Chinese students acknowledge the support of the YutChun-Weizmann Program that enabled this study. The study was also supported by a research grant from the Center for Scientific Excellence at the Weizmann Institute of Science. We are grateful to Dr. Shifra Ben-Dor for valuable discussions, to Prof. Robin Gasser (University of Melbourne) for providing us with the sequence of HCO_011565, to Prof. Keith Dunker (University of Indiana) for recommending the fIDPnn algorithm for disorder prediction, and to Dr. Sergey Ovchinnikov (Harvard University) for valuable advice concerning the use of AF2 Colab.

CONFLICT OF INTEREST STATEMENT

The authors declare no conflicts of interest.

PEER REVIEW

The peer review history for this article is available at <https://www.webofscience.com/api/gateway/wos/peer-review/10.1002/prot.26496>.

DATA AVAILABILITY STATEMENT

All data are presented in the article or in pointers to UniProt (<https://www.uniprot.org>) or the PDB (<https://www.rcsb.org>). The data that support the findings of this study are available from the corresponding author upon reasonable request.

ORCID

Israel Silman  <https://orcid.org/0000-0003-1923-0829>

Joel L. Sussman  <https://orcid.org/0000-0003-0306-3878>

REFERENCES

- Bränden C, Tooze J. *Introduction to Protein Structure*. 2nd ed. Garland Publishing, Inc.; 1999.
- Pe'er I, Felder CE, Man O, Silman I, Sussman JL, Beckmann JS. Proteomic signatures: amino acid and oligopeptide compositions differentiate among phyla. *Proteins*. 2004;54(1):20–40.
- Lavelle DT, Pearson WR. Globally, unrelated protein sequences appear random. *Bioinformatics*. 2010;26(3):310–318.
- De Luca D, Slanzi D, Poli I, Polticelli F, Minervini G. Do natural proteins differ from random sequences polypeptides? Natural vs. random proteins classification using an evolutionary neural network. *PLoS One*. 2012;7(5):e36634.
- Geffen Y, Ofra Y, Unger R. DistilProtBert: a distilled protein language model used to distinguish between real proteins and their randomly shuffled counterparts. *Bioinformatics*. 2022;38(Supplement_2):ii95–ii98.
- Chiarabelli C, Vrijbloed JW, Thomas RM, Luisi PL. Investigation of de novo totally random biosequences, part I: a general method for in vitro selection of folded domains from a random polypeptide library displayed on phage. *Chem Biodivers*. 2006;3(8):827–839.
- Chiarabelli C, Vrijbloed JW, De Luca D, et al. Investigation of de novo totally random biosequences, part II: on the folding frequency in a totally random library of de novo proteins obtained by phage display. *Chem Biodivers*. 2006;3(8):840–859.
- Tretyachenko V, Vymetal J, Bednarova L, et al. Random protein sequences can form defined secondary structures and are well-tolerated in vivo. *Sci Rep*. 2017;7(1):15449.
- Davidson AR, Lumb KJ, Sauer RT. Cooperatively folded proteins in random sequence libraries. *Nat Struct Biol*. 1995;2(10):856–864.
- Keefe AD, Szostak JW. Functional proteins from a random-sequence library. *Nature*. 2001;410(6829):715–718.
- Hecht MH, Das A, Go A, Bradley LH, Wei Y. De novo proteins from designed combinatorial libraries. *Protein Sci*. 2004;13(7):1711–1723.
- Baek M, DiMaio F, Anishchenko I, et al. Accurate prediction of protein structures and interactions using a three-track neural network. *Science*. 2021;373(6557):871–876.
- Anishchenko I, Pellock SJ, Chidyausiku TM, et al. De novo protein design by deep network hallucination. *Nature*. 2021;600(7889):547–552.
- Tautz D, Domazet-Loso T. The evolutionary origin of orphan genes. *Nat Rev Genet*. 2011;12(10):692–702.
- Carvunis AR, Rolland T, Wapinski I, et al. Proto-genes and de novo gene birth. *Nature*. 2012;487(7407):370–374.
- Reinhardt JA, Wanjiu BM, Brant AT, Saelao P, Begun DJ, Jones CD. De novo ORFs in *Drosophila* are important to organismal fitness and evolved rapidly from previously non-coding sequences. *PLoS Genet*. 2013;9(10):e1003860.
- McLysaght A, Guerzoni D. New genes from non-coding sequence: the role of de novo protein-coding genes in eukaryotic evolutionary innovation. *Philos Trans R Soc Lond B Biol Sci*. 2015;370(1678):20140332.
- Schmitz JF, Bornberg-Bauer E. Fact or fiction: updates on how protein-coding genes might emerge de novo from previously non-coding DNA. *F1000Res*. 2017;6:57.
- Vakirlis N, Carvunis AR, McLysaght A. Synteny-based analyses indicate that sequence divergence is not the main source of orphan genes. *Elife*. 2020;9:e53500.
- Li J, Singh U, Bhandary P, et al. Foster thy young: enhanced prediction of orphan genes in assembled genomes. *Nucleic Acids Res*. 2022;50(7):e37.
- Jacob F. Evolution and tinkering. *Science*. 1977;196(4295):1161–1166.
- Schmitz J, Brosius J. Exonization of transposed elements: a challenge and opportunity for evolution. *Biochimie*. 2011;93(11):1928–1934.
- Khalturin K, Hemmrich G, Fraune S, Augustin R, Bosch TC. More than just orphans: are taxonomically-restricted genes important in evolution? *Trends Genet*. 2009;25(9):404–413.
- Knowles DG, McLysaght A. Recent de novo origin of human protein-coding genes. *Genome Res*. 2009;19(10):1752–1759.
- Siepel A. Darwinian alchemy: human genes from noncoding DNA. *Genome Res*. 2009;19(10):1693–1695.
- Toll-Riera M, Bosch N, Bellora N, et al. Origin of primate orphan genes: a comparative genomics approach. *Mol Biol Evol*. 2009;26(3):603–612.
- Wu DD, Irwin DM, Zhang YP. De novo origin of human protein-coding genes. *PLoS Genet*. 2011;7(11):e1002379.

28. Xie C, Zhang YE, Chen JY, et al. Hominoid-specific de novo protein-coding genes originating from long non-coding RNAs. *PLoS Genet*. 2012;8(9):e1002942.
29. Dowling D, Schmitz JF, Bornberg-Bauer E. Stochastic gain and loss of novel transcribed open reading frames in the human lineage. *Genome Biol Evol*. 2020;12(11):2183-2195.
30. Ruiz-Orera J, Messegue X, Subirana JA, Alba MM. Long non-coding RNAs as a source of new peptides. *Elife*. 2014;3:e03523.
31. Jumper J, Evans R, Pritzel A, et al. Highly accurate protein structure prediction with AlphaFold. *Nature*. 2021;596(7873):583-589.
32. Lin Z, Akin H, Rao R, et al. Language models of protein sequences at the scale of evolution enable accurate structure prediction. *bioRxiv*. 2022:2022.2007.2020.500902.
33. Sussman JL, Lin D, Jiang J, et al. Protein data bank (PDB): a database of 3D structural information of biological macromolecules. *Acta Crystallogr D Biol Crystallogr*. 1998;54(Pt 6 Pt 1):1078-1084.
34. Berman HM, Westbrook J, Feng Z, et al. The Protein Data Bank. *Nucleic Acids Res*. 2000;28(1):235-242.
35. Kryshchavych A, Schwede T, Topf M, Fidelis K, Moult J. Critical assessment of methods of protein structure prediction (CASP)-round XIV. *Proteins*. 2021;89(12):1607-1617.
36. Dunker AK, Silman I, Uversky VN, Sussman JL. Function and structure of inherently disordered proteins. *Curr Opin Struct Biol*. 2008;18(6):756-764.
37. Andring JT, Kim CU, McKenna R. Structure and mechanism of copper-carbonic anhydrase II: a nitrite reductase. *IUCr*. 2020;7(Pt 2):287-293.
38. Zeev-Ben-Mordehai T, Rydberg EH, Solomon A, et al. The intracellular domain of the *Drosophila* cholinesterase-like neural adhesion protein, gliotactin, is natively unfolded. *Proteins*. 2003;53(3):758-767.
39. Adkins JN, Lumb KJ. Intrinsic structural disorder and sequence features of the cell cycle inhibitor p57Kip2. *Proteins*. 2002;46(1):1-7.
40. Kurzbach D, Platzer G, Schwarz TC, Henen MA, Konrat R, Hinderberger D. Cooperative unfolding of compact conformations of the intrinsically disordered protein osteopontin. *Biochemistry*. 2013;52(31):5167-5175.
41. Bakolitsa C, Schwarzenbacher R, McMullan D, et al. Crystal structure of an orphan protein (TM0875) from *Thermotoga maritima* at 2.00-Å resolution reveals a new fold. *Proteins*. 2004;56(3):607-610.
42. Lim K, Sarikaya E, Galkin A, et al. Novel structure and nucleotide binding properties of HI1480 from *Haemophilus influenzae*: a protein with no known sequence homologues. *Proteins*. 2004;56(3):564-571.
43. Cheng C, Shaw N, Zhang X, et al. Structural view of a non Pfam singleton and crystal packing analysis. *PLoS One*. 2012;7(2):e31673.
44. Taki AC, Wang T, Nguyen NN, et al. Thermal proteome profiling reveals *Haemonchus* orphan protein HCO_011565 as a target of the nematocidal small molecule UMW-868. *Front Pharmacol*. 2022;13:1014804.
45. Samusik N, Krukovskaya L, Meln I, Shilov E, Kozlov AP. PBOV1 is a human de novo gene with tumor-specific expression that is associated with a positive clinical outcome of cancer. *PLoS One*. 2013;8(2):e56162.
46. Li CY, Zhang Y, Wang Z, et al. A human-specific de novo protein-coding gene associated with human brain functions. *PLoS Comput Biol*. 2010;6(3):e1000734.
47. Suenaga Y, Islam SM, Alagu J, et al. NCYM, a cis-antisense gene of MYCN, encodes a de novo evolved protein that inhibits GSK3beta resulting in the stabilization of MYCN in human neuroblastomas. *PLoS Genet*. 2014;10(1):e1003996.
48. Xie C, Bekpen C, Kunzel S, et al. A de novo evolved gene in the house mouse regulates female pregnancy cycles. *Elife*. 2019;8:e44392.
49. Perochon A, Jianguang J, Kahla A, et al. TaFROG encodes a *Pooideae* orphan protein that interacts with SnRK1 and enhances resistance to the Mycotoxigenic fungus *Fusarium graminearum*. *Plant Physiol*. 2015;169(4):2895-2906.
50. Casco-Robles RM, Watanabe A, Eto K, et al. Novel erythrocyte clumps revealed by an orphan gene *Newt1* in circulating blood and regenerating limbs of the adult newt. *Sci Rep*. 2018;8(1):7455.
51. Mirdita M, Schütze K, Moriawaki Y, Heo L, Ovchinnikov S, Steinegger M. ColabFold: making protein folding accessible to all. *Nat Methods*. 2022;19(6):679-682.
52. Mariani V, Biasini M, Barbato A, Schwede T. IDDT: a local superposition-free score for comparing protein structures and models using distance difference tests. *Bioinformatics*. 2013;29(21):2722-2728.
53. Hiranuma N, Park H, Baek M, Anishchenko I, Dauparas J, Baker D. Improved protein structure refinement guided by deep learning based accuracy estimation. *Nat Commun*. 2021;12(1):1340.
54. Holm L, Sander C. Touring protein fold space with Dali/FSSP. *Nucleic Acids Res*. 1998;26(1):316-319.
55. Holm L. DALI and the persistence of protein shape. *Protein Sci*. 2020;29(1):128-140.
56. Prilusky J, Felder CE, Zeev-Ben-Mordehai T, et al. FoldIndex: a simple tool to predict whether a given protein sequence is intrinsically unfolded. *Bioinformatics*. 2005;21(16):3435-3438.
57. Hu G, Katuwawala A, Wang K, et al. fIDPnn: accurate intrinsic disorder prediction with putative propensities of disorder functions. *Nat Commun*. 2021;12(1):4438.
58. Uversky VN. Intrinsically disordered proteins and their "mysterious" (meta)physics. *Front Phys*. 2019;7:10.
59. Dunker AK, Obradovic Z, Romero P, Garner EC, Brown CJ. Intrinsic protein disorder in complete genomes. *Genome Inform*. 2000;11:161-171.
60. Oldfield CJ, Cheng Y, Cortese MS, Brown CJ, Uversky VN, Dunker AK. Comparing and combining predictors of mostly disordered proteins. *Biochemistry*. 2005;44(6):1989-2000.
61. He J, Turzo SBA, Seffernick JT, Kim SS, Lindert S. Prediction of intrinsic disorder using Rosetta ResidueDisorder and AlphaFold2. *J Chem Phys B*. 2022;126(42):8439-8446.
62. Ruff KM, Pappu RV. AlphaFold and implications for intrinsically disordered proteins. *J Mol Biol*. 2021;433(20):167208.
63. Alderson T, Pritisanac I, Moses A, Forman-Kay J. Systematic identification of conditionally folded intrinsically disordered regions by AlphaFold2. *bioRxiv*. 2022:2022.2002.2018.481080.
64. Vakirlis N, Vance Z, Duggan KM, McLysaght A. De novo birth of functional microproteins in the human lineage. *Cell Rep*. 2022;41(12):111808.
65. Holm L. Dali server: structural unification of protein families. *Nucleic Acids Res*. 2022;50(W1):W210-W215.
66. Necci M, Piovesan D, CAID Predictors, DisProt Curators, Tosatto SCE. Critical assessment of protein intrinsic disorder prediction. *Nat Methods*. 2021;18(5):472-481.
67. Monzon V, Haft DH, Bateman A. Folding the unfoldable: using AlphaFold to explore spurious proteins. *Bioinform Adv*. 2022;2(1):vbab043.
68. Eberhardt RY, Haft DH, Punta M, Martin M, O'Donovan C, Bateman A. AntiFam: a tool to help identify spurious ORFs in protein annotation. *Database*. 2012;2012:bas003.
69. Rost B. Did evolution leap to create the protein universe? *Curr Opin Struct Biol*. 2002;12(3):409-416.
70. Coulson AF, Moult J. A unfold, mesofold, and superfold model of protein fold use. *Proteins*. 2002;46(1):61-71.
71. Light S, Basile W, Elofsson A. Orphans and new gene origination, a structural and evolutionary perspective. *Curr Opin Struct Biol*. 2014;26:73-83.
72. Bornberg-Bauer E, Hlouchova K, Lange A. Structure and function of naturally evolved de novo proteins. *Curr Opin Struct Biol*. 2021;68:175-183.
73. Uversky VN, Gillespie JR, Fink AL. Why are "natively unfolded" proteins unstructured under physiologic conditions? *Proteins*. 2000;41(3):415-427.
74. Ptitsyn O. How molten is the molten globule? *Nat Struct Biol*. 1996;3(6):488-490.

75. Arai M, Kuwajima K. Role of the molten globule state in protein folding. *Adv Protein Chem.* 2000;53:209-282.
76. Dolginova EA, Roth E, Silman I, Weiner LM. Chemical modification of *Torpedo* acetylcholinesterase by disulfides: appearance of a "molten globule" state. *Biochemistry.* 1992;31(48):12248-12254.
77. Jumper J, Hassabis D. Protein structure predictions to atomic accuracy with AlphaFold. *Nat Methods.* 2022;19(1):11-12.
78. Marks DS, Hopf TA, Sander C. Protein structure prediction from sequence variation. *Nat Biotechnol.* 2012;30(11):1072-1080.
79. Chowdhury R, Bouatta N, Biswas S, et al. Single-sequence protein structure prediction using a language model and deep learning. *Nat Biotechnol.* 2022;40(11):1617-1623.

SUPPORTING INFORMATION

Additional supporting information can be found online in the Supporting Information section at the end of this article.

How to cite this article: Liu J, Yuan R, Shao W, Wang J, Silman I, Sussman JL. Do "Newly Born" orphan proteins resemble "Never Born" proteins? A study using three deep learning algorithms. *Proteins.* 2023;91(8):1097-1115. doi:10.1002/prot.26496