

Large language models generate functional protein sequences across diverse families

Received: 12 July 2022

Accepted: 17 November 2022

Published online: 26 January 2023

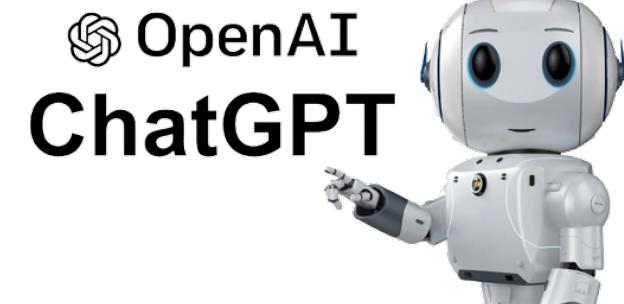
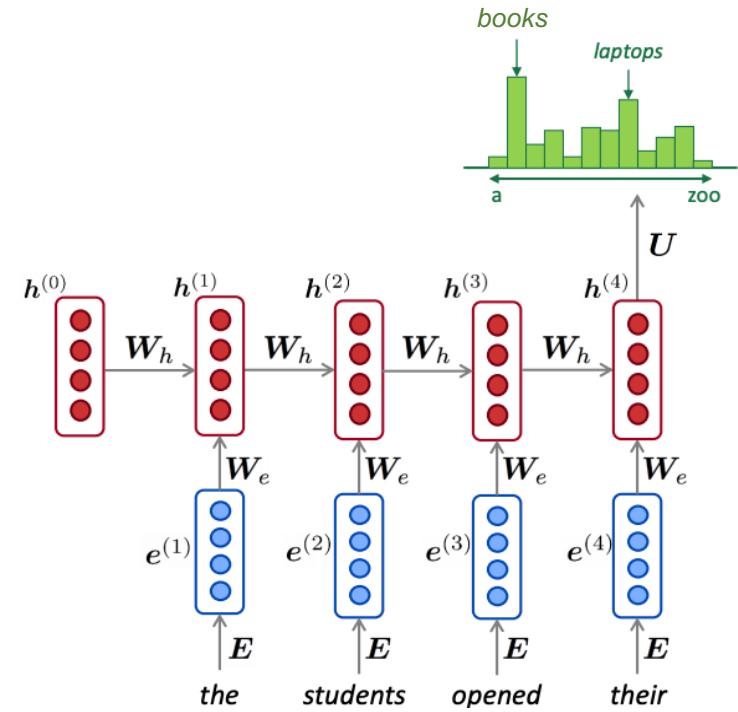
 Check for updates

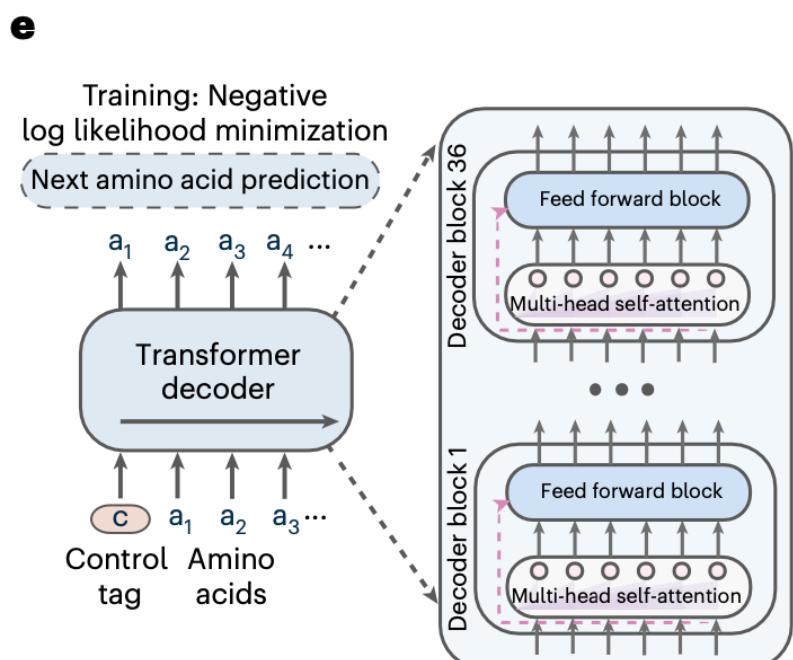
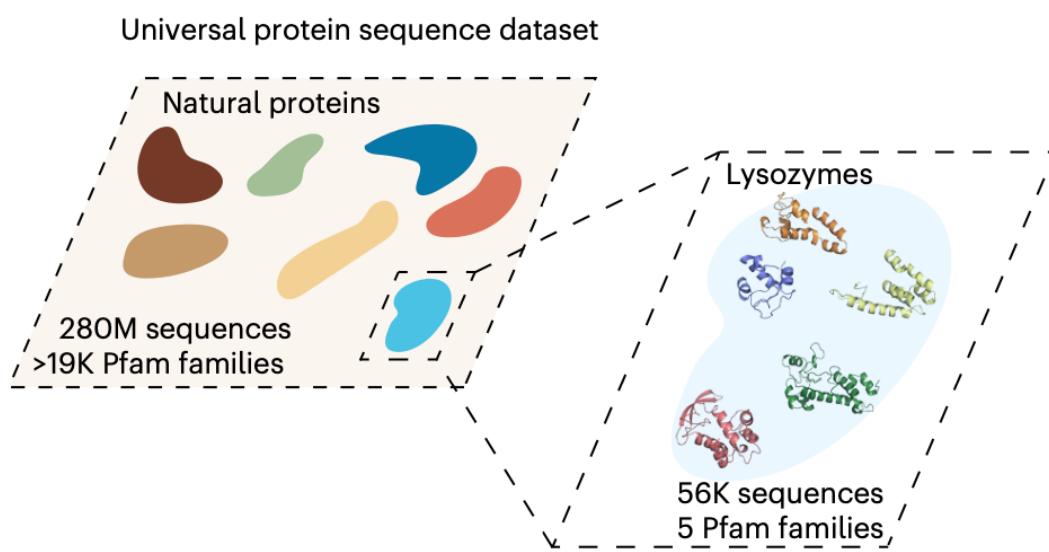
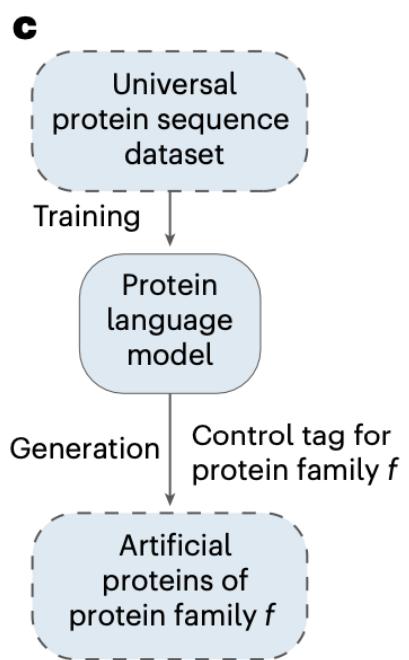
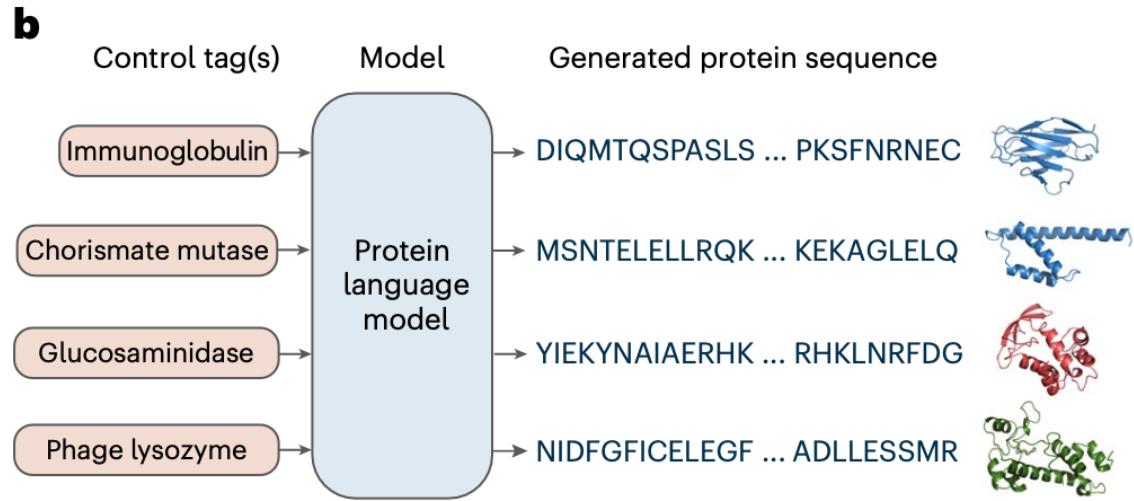
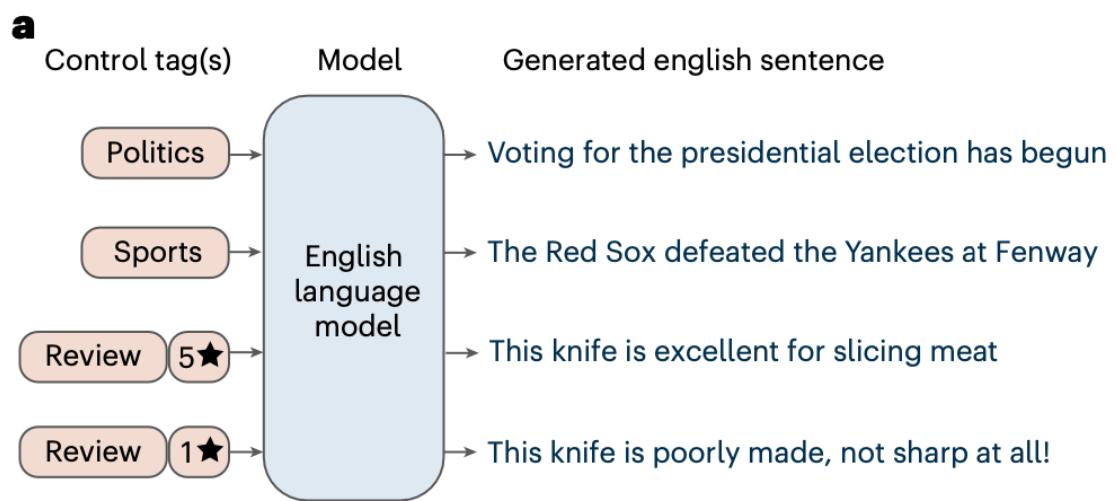
Ali Madani^{1,2}✉, Ben Krause^{1,10}, Eric R. Greene^{3,10}, Subu Subramanian^{4,5}, Benjamin P. Mohr⁶, James M. Holton¹✉, Jose Luis Olmos Jr.³, Caiming Xiong¹, Zachary Z. Sun⁶, Richard Socher¹, James S. Fraser³ & Nikhil Naik¹✉

Deep-learning language models have shown promise in various biotechnological applications, including protein design and engineering. Here we describe ProGen, a language model that can generate protein sequences with a predictable function across large protein families, akin to generating grammatically and semantically correct natural language sentences on diverse topics. The model was trained on 280 million protein sequences from >19,000 families and is augmented with control tags specifying protein properties. ProGen can be further fine-tuned to curated sequences and tags to improve controllable generation performance of proteins from families with sufficient homologous samples. Artificial proteins fine-tuned to five distinct lysozyme families showed similar catalytic efficiencies as natural lysozymes, with sequence identity to natural proteins as low as 31.4%. ProGen is readily adapted to diverse protein families, as we demonstrate with chorismate mutase and malate dehydrogenase.

What's large language models?

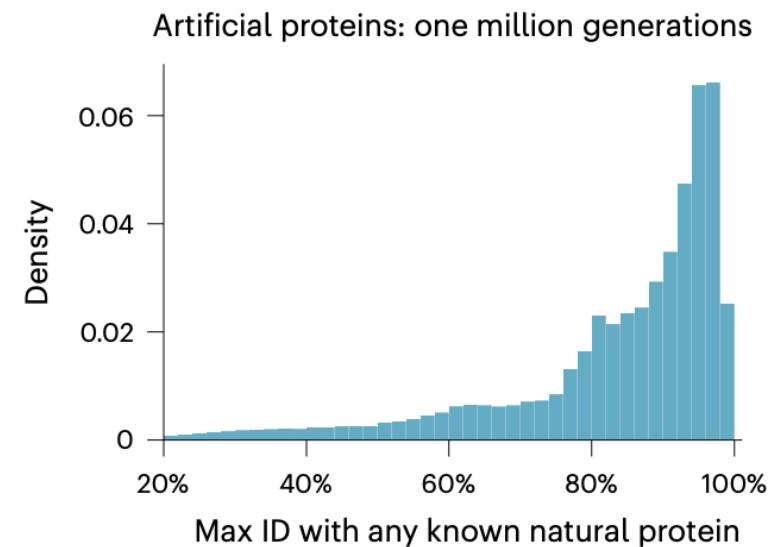
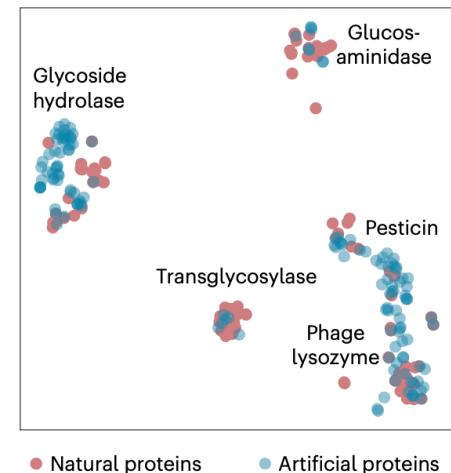
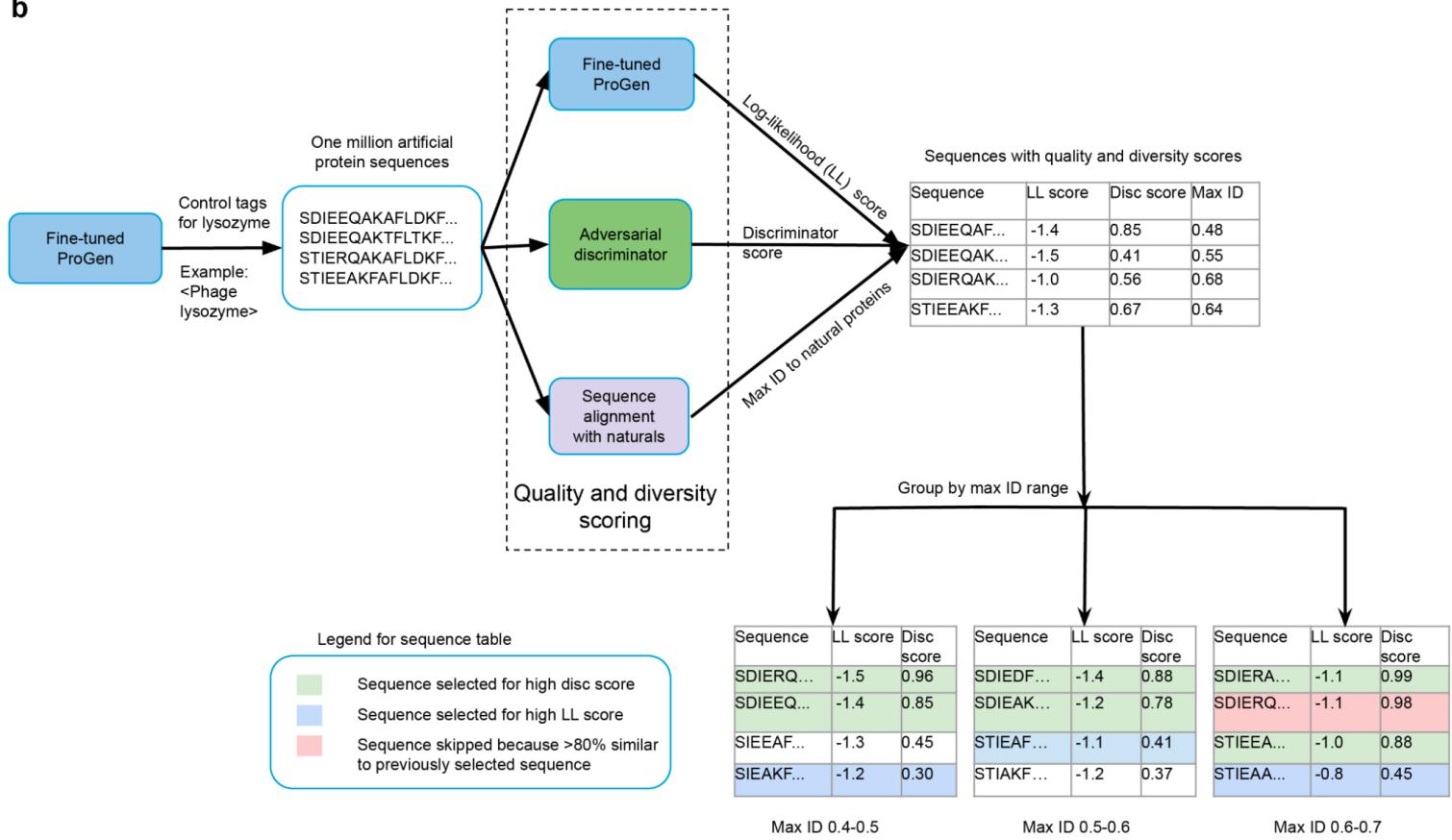
- Large language models are a type of artificial intelligence (AI) system that are trained on vast amounts of text data and can generate natural language responses to a wide range of inputs.
- A deep learning model trained on a large corpus of text data, which enables the model to learn the patterns and structures of language at a very high level of abstraction, allowing it to generate new text that is grammatically correct and semantically meaningful.
- One of the most well-known large language models is GPT-3, which was developed by OpenAI and contains over 175 billion parameters. GPT-3 has been trained on a vast amount of text data, including books, articles, and websites, and can generate responses to a wide range of prompts, such as questions, writing prompts, or even computer code.

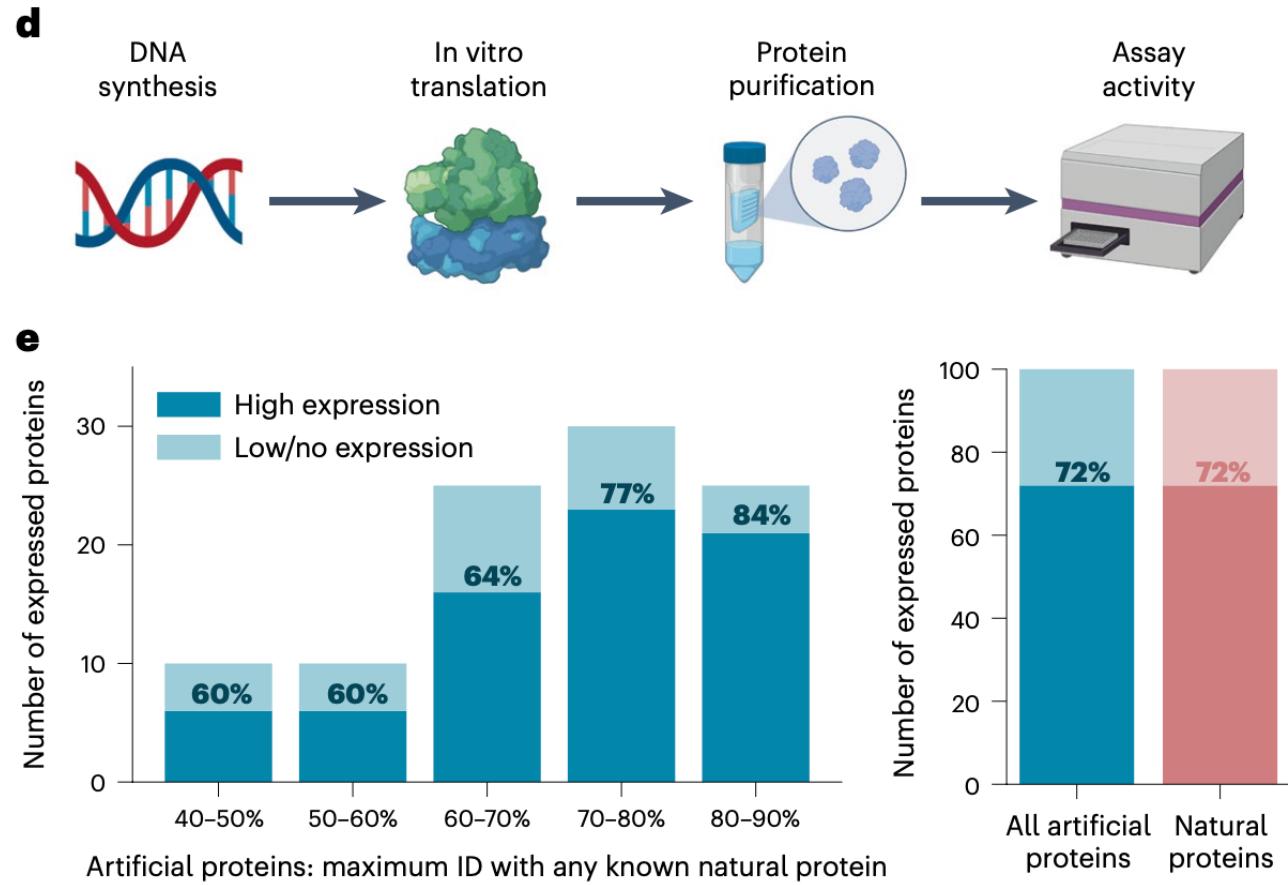
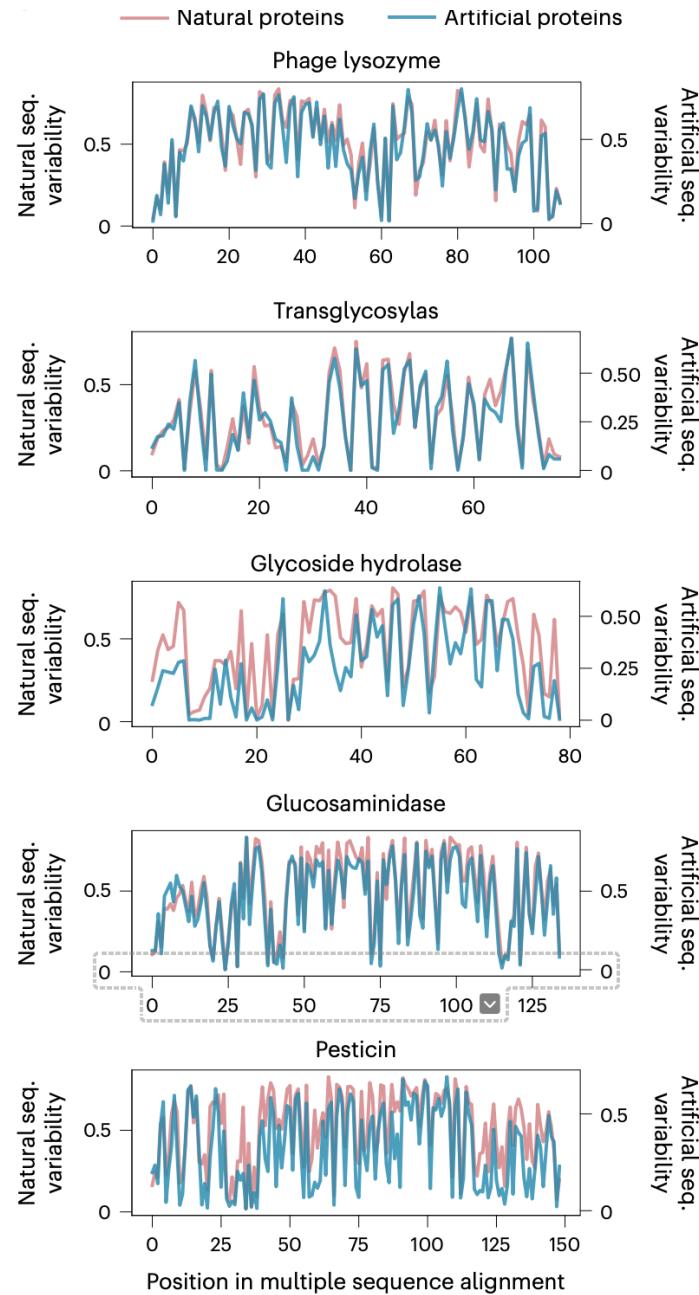


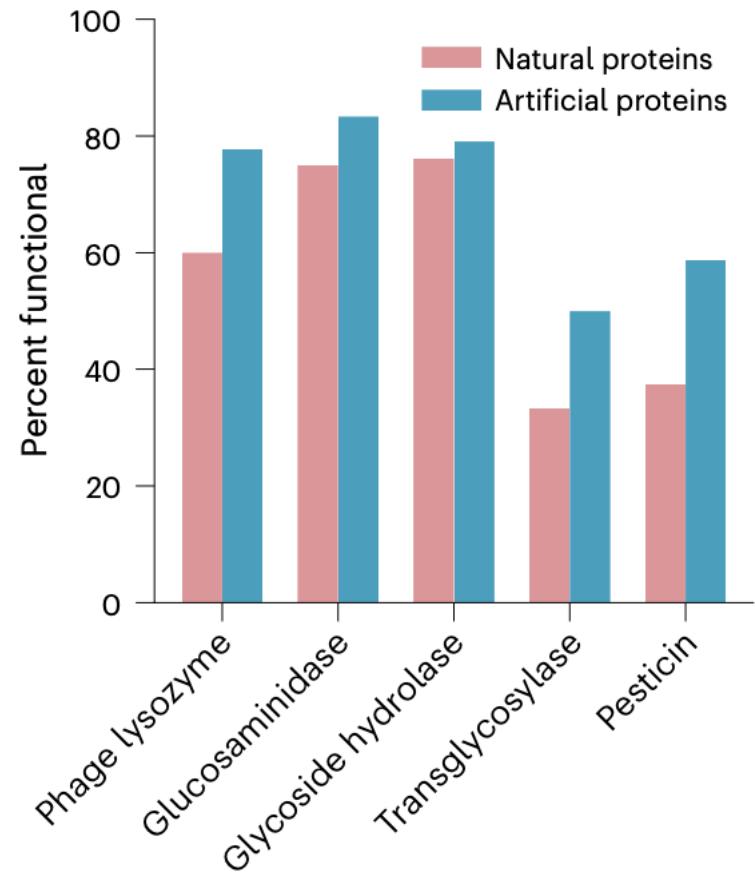
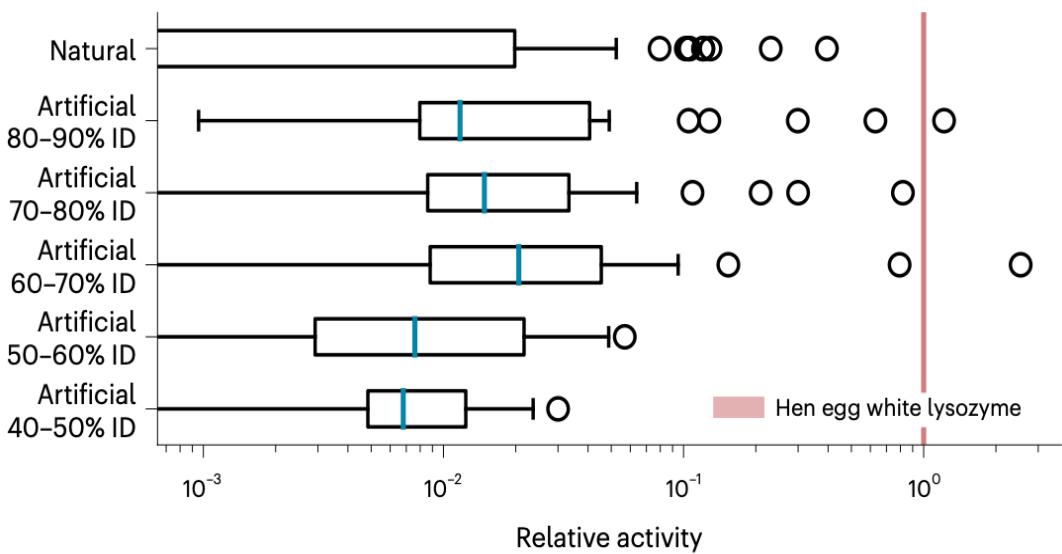
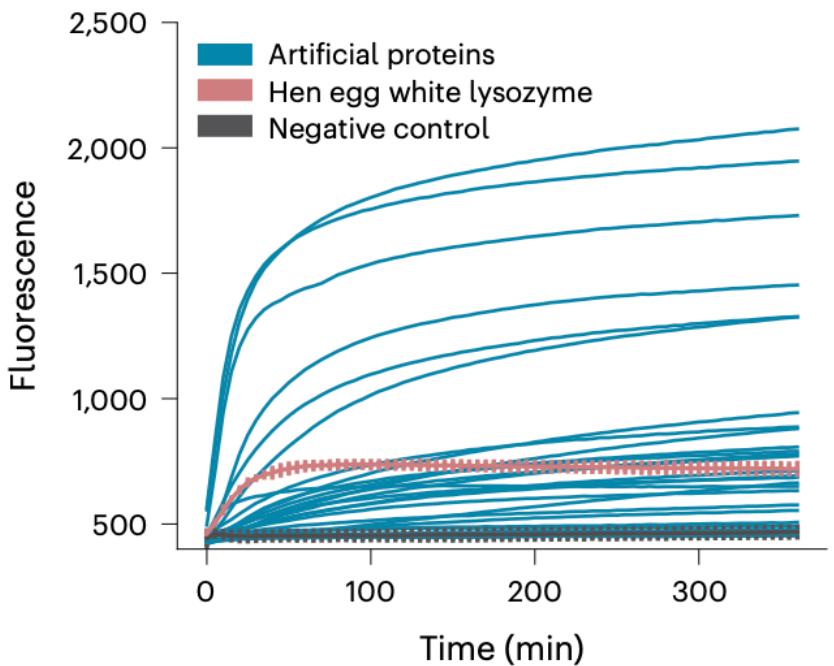


Generate artificial antibacterial proteins using ProGen

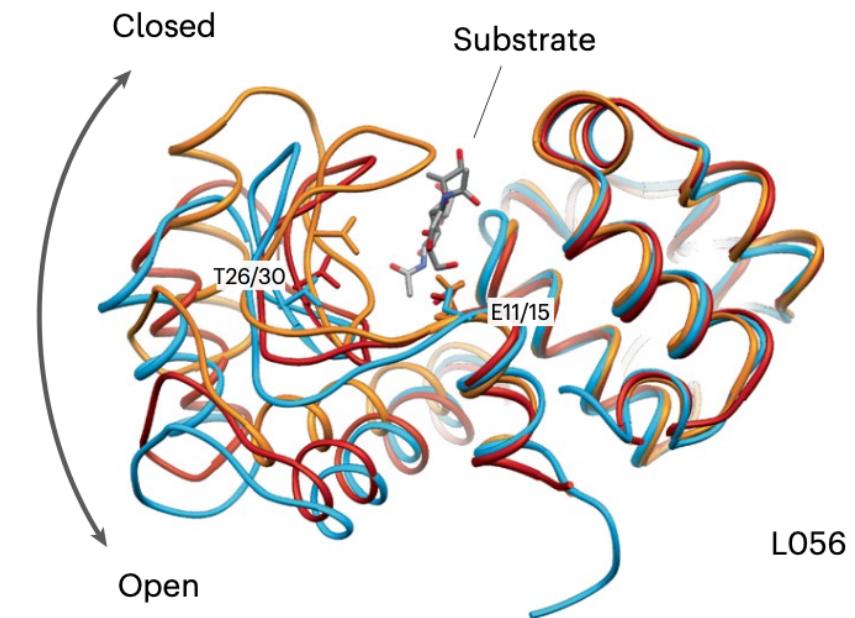
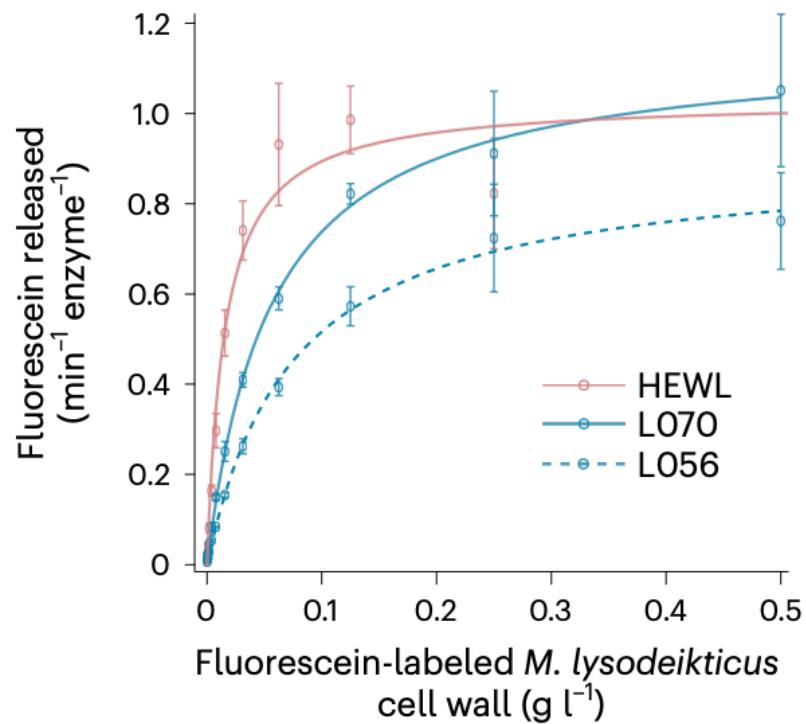
b





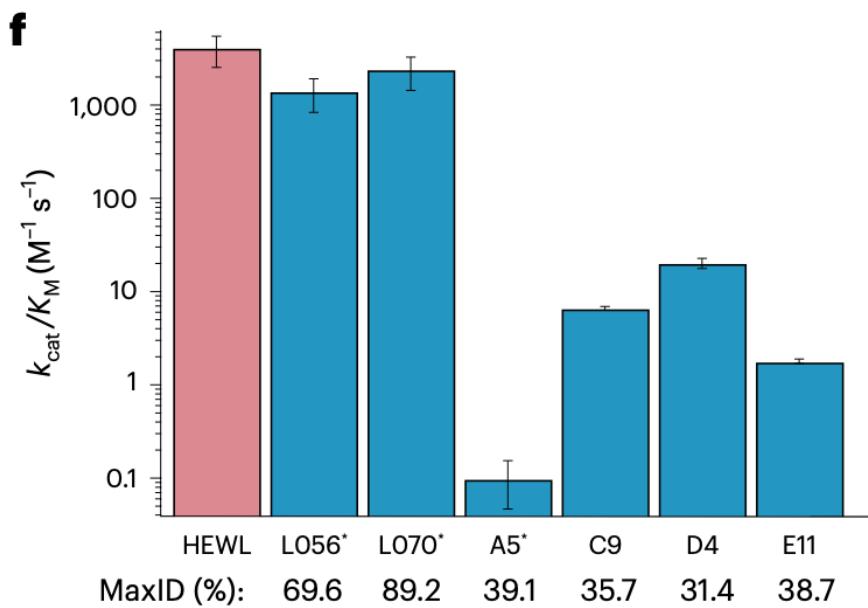


Two synthesized protein L070 and L056 (<40% max ID) gain competent enzyme activity to HEWL with 53 and 18 amino acid different from the nearest natural protein. L070 and L056 only share 17.9% sequence identity.

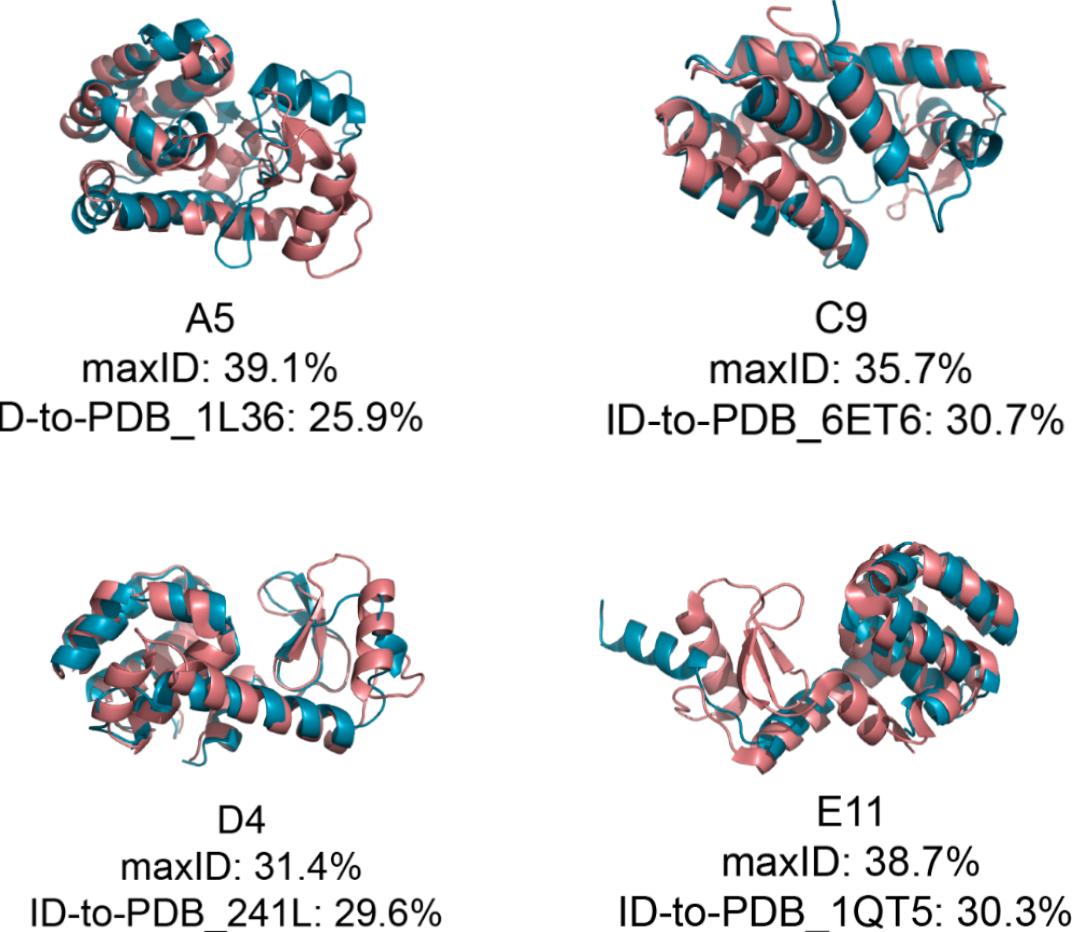


A global overlay of L056 crystal structure with two representative T4 lysozyme conformations is shown with L056 presented in sky blue, 'open' conformation of M6I T4 lysozyme (PDB accession [150L](#)) in dark red, 'closed' conformation of wild-type T4 lysozyme (PDB accession [3FA0](#)) in orange,

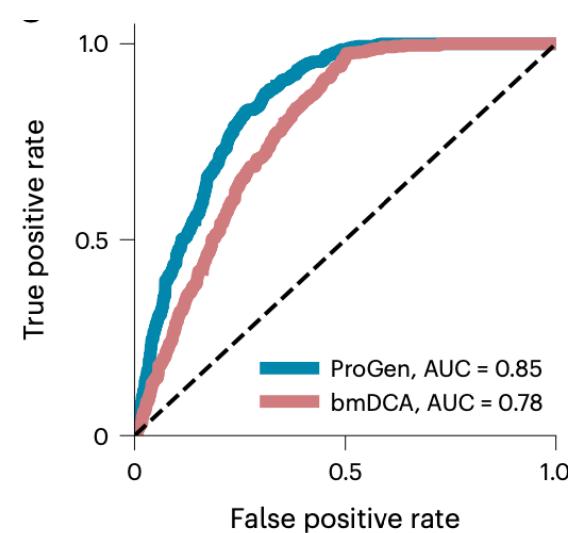
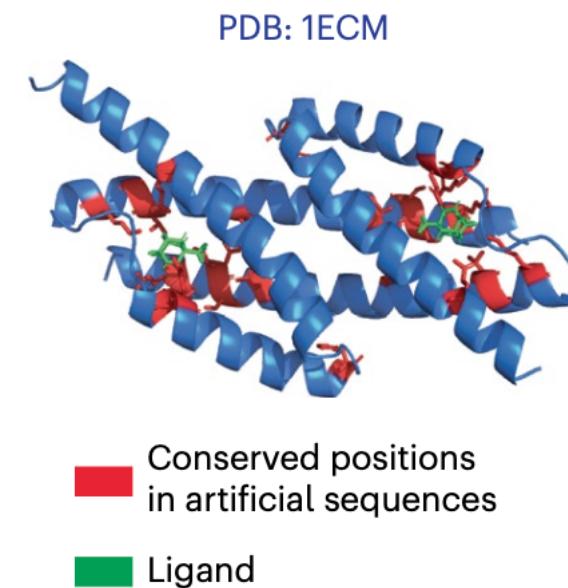
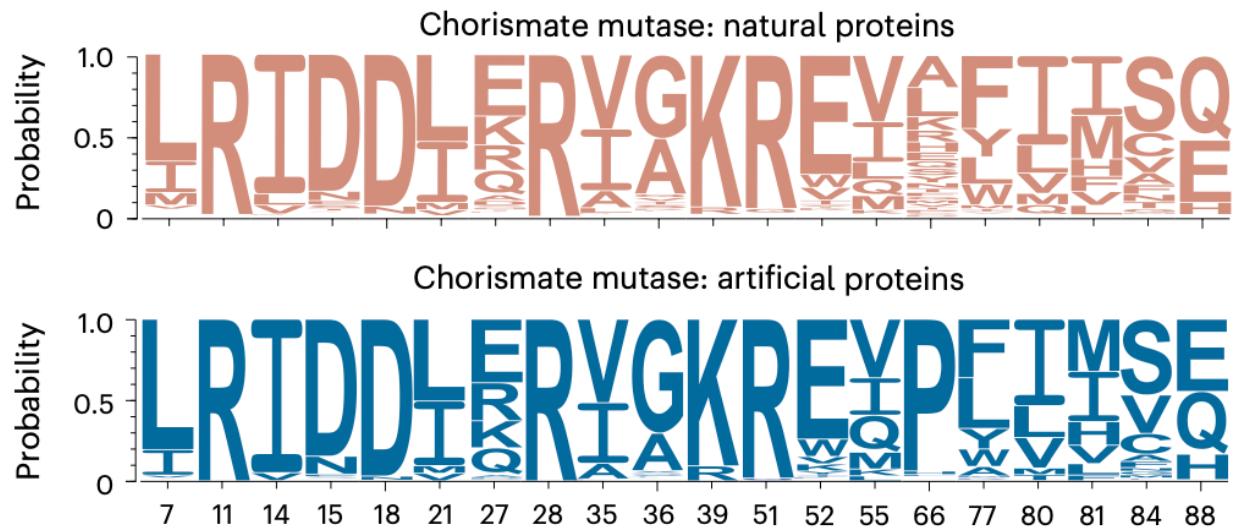
6 artificial sequences (<40% max ID)
with gained lysosome activity



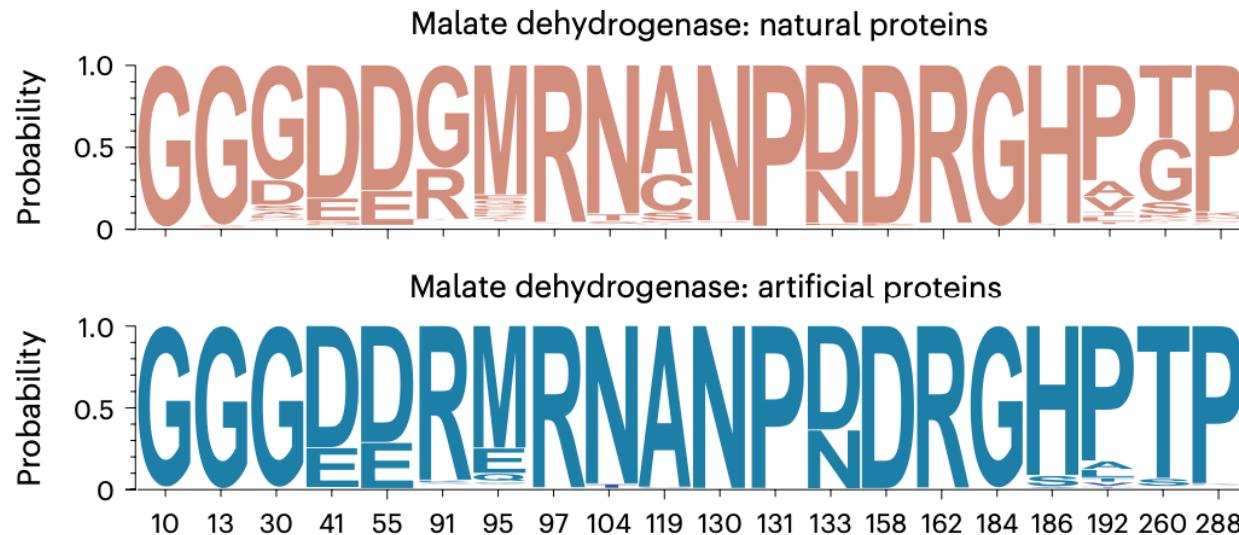
The 3D structure roughly match with nearest
natural protein despite low 2D similarity



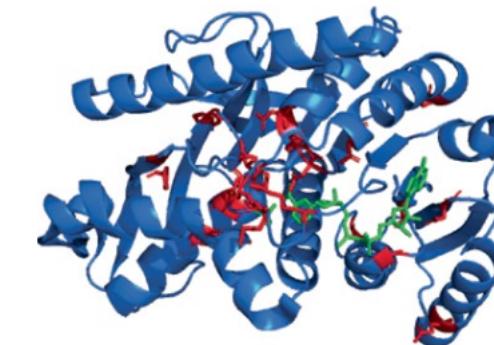
Generate protein sequences for chorismate mutase protein family using ProGen



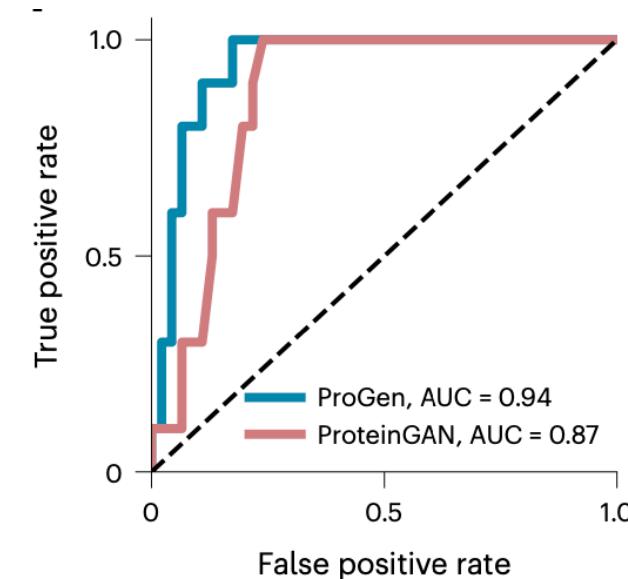
Generate protein sequences for malate dehydrogenase protein family using ProGen



PDB: 4MDH



- Conserved positions in artificial sequences
- Ligand





Resource

<https://doi.org/10.1038/s41592-022-01685-y>

AlphaFill: enriching AlphaFold models with ligands and cofactors

Received: 10 December 2021

Accepted: 18 October 2022

Published online: 24 November 2022

 Check for updates

Maarten L. Hekkelman ^{1,2}, Ida de Vries ^{1,2}, Robbie P. Joosten ^{1,3} & Anastassis Perrakis ^{1,3}

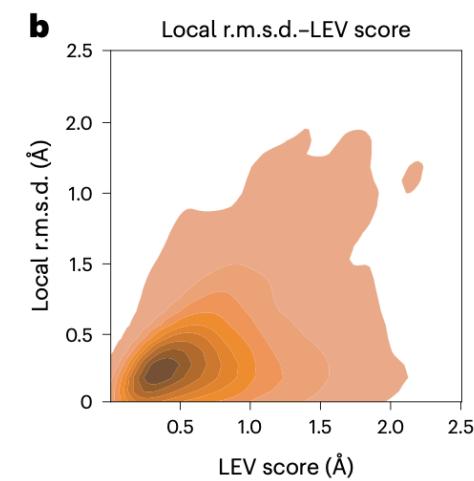
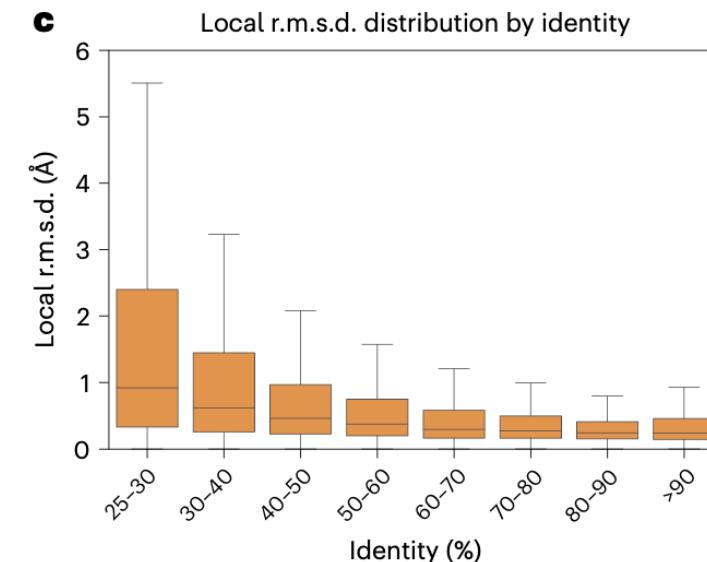
Artificial intelligence-based protein structure prediction approaches have had a transformative effect on biomolecular sciences. The predicted protein models in the AlphaFold protein structure database, however, all lack coordinates for small molecules, essential for molecular structure or function: hemoglobin lacks bound heme; zinc-finger motifs lack zinc ions essential for structural integrity and metalloproteases lack metal ions needed for catalysis. Ligands important for biological function are absent too; no ADP or ATP is bound to any of the ATPases or kinases. Here we present AlphaFill, an algorithm that uses sequence and structure similarity to ‘transplant’ such ‘missing’ small molecules and ions from experimentally determined structures to predicted protein models. The algorithm was successfully validated against experimental structures. A total of 12,029,789 transplants were performed on 995,411 AlphaFold models and are available together with associated validation metrics in the alphafill.eu databank, a resource to help scientists make new hypotheses and design targeted experiments.

Process for ligand co-factor transplant

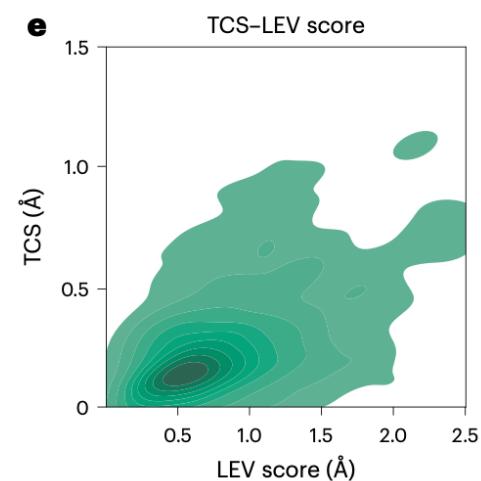
First, we search for sequence homologs for each structure in the AlphaFold database in the PDB-REDO databank. We consider structures with identity higher than 25% over an aligned sequence of at least 85 residues as hits.

Next, the selection of structures with compounds of interest are structurally aligned on the C α -atoms of the AlphaFold model, and the root-mean-square deviation (r.m.s.d.) is calculated (global r.m.s.d.).

Starting from the closest homolog, all backbone atoms within 6 Å from the atoms of each compound that will be considered for ‘transplantation’ are selected and used for a local structural alignment to the AlphaFold model; the r.m.s.d. from this alignment is also calculated (local r.m.s.d.). Compounds are then transplanted into the AlphaFold model to make the AlphaFill model.



Local environment validation (LEV) score as the all-atom r.m.s.d. of any ligand atom and all proteins' atoms within 6.0 Å from the ligand, between the AlphaFill and experimental complexes.



Transplant clash score (TCS) as a function of the van der Waals overlaps between a transplanted ligand and its binding site

P29373

Cellular retinoic acid-binding protein 2

Structure file

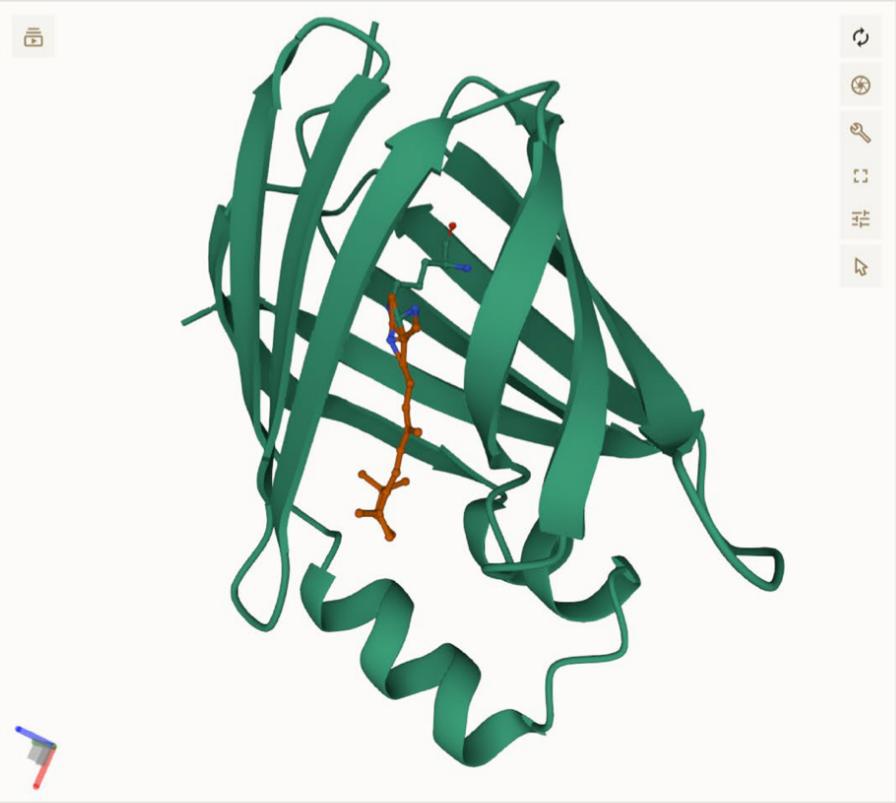
<https://alphafill.eu/v1/aff/P29373-F1>

Metadata

<https://alphafill.eu/v1/aff/P29373-F1/json>

Original AlphaFold model

<https://alphafold.ebi.ac.uk/entry/P29373>

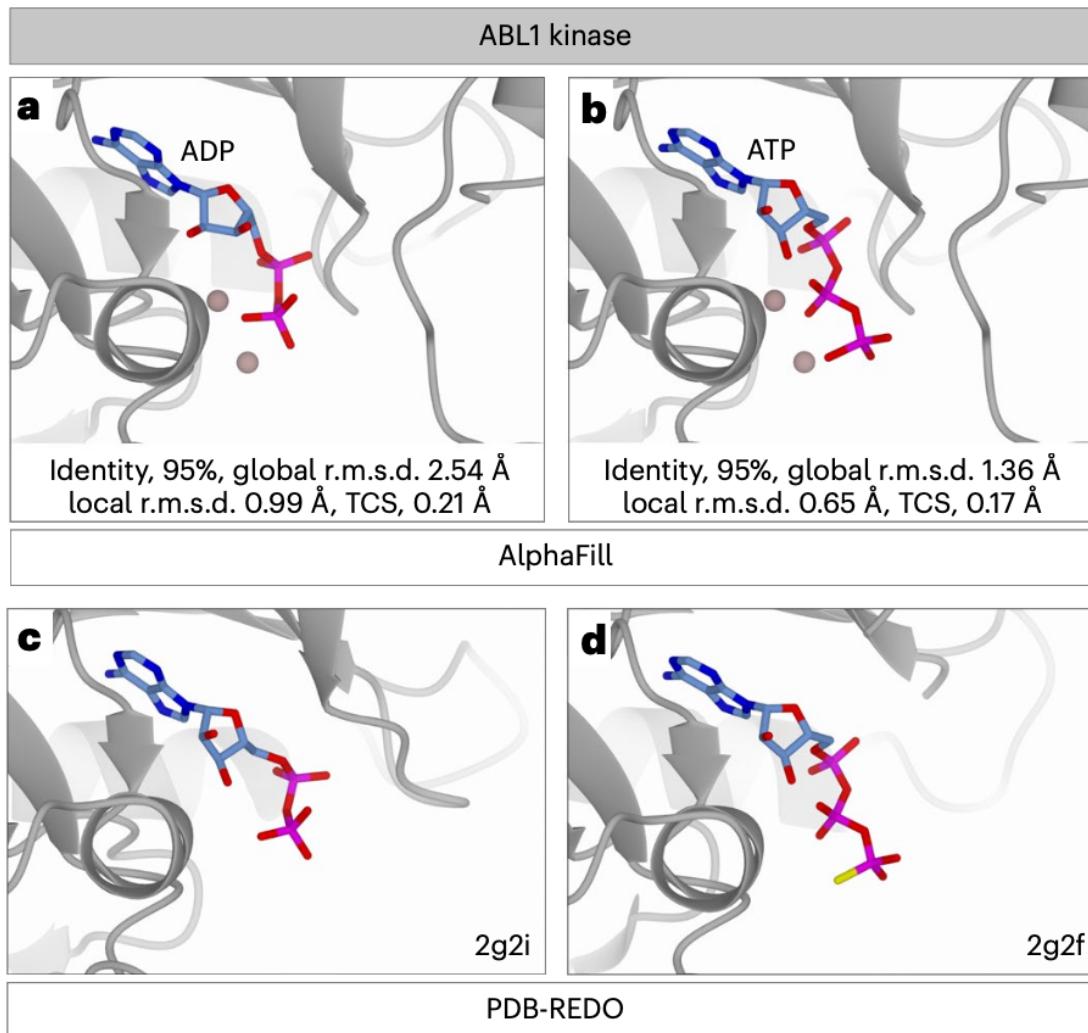


Download PDB

View in AlphaFold

25% identity30% identity40% identity50% identity60% identity70% identity

Compound	PDBID	Global RMSd	Asym	Local RMSd	TCS	Show	
B3P	3fel.A	0.56	M	0.20	0.09	<input type="checkbox"/>	
			L	0.21	0.30	<input type="checkbox"/>	
NA	2g78.A	0.36	G	0.19	0.00	<input type="checkbox"/>	
	2fs6.A	0.37	C	0.17	0.32	<input type="checkbox"/>	
	2g79.A	0.46	I	0.04	0.00	<input type="checkbox"/>	
			J	0.08	0.00	<input type="checkbox"/>	
			H	0.09	0.00	<input type="checkbox"/>	
		2g7b.A	0.52	O	0.13	0.00	<input type="checkbox"/>
			N	0.23	0.34	<input type="checkbox"/>	
		2frs.A	0.69	E	0.22	0.14	<input type="checkbox"/>
			D	0.25	0.09	<input type="checkbox"/>	
			F	0.96	0.32	<input type="checkbox"/>	
REA	3d97.B	0.94	P	0.04	0.00	<input type="checkbox"/>	
			Q	0.05	0.48	<input type="checkbox"/>	
RET	1cbs.A	0.36	B	0.52	0.21	<input type="checkbox"/>	
	1cbr.A	0.62	T	0.04	0.00	<input type="checkbox"/>	
RET	2g79.A	0.46	K	0.59	0.14	<input type="checkbox"/>	
	4i9s.A	0.76	R	0.78	0.99?	<input checked="" type="checkbox"/> optimise	



For the human tyrosine-protein kinase ABL1 ([AF-P00519](#)) the AlphaFill model shows an ADP molecule and an ATP molecule (Fig. 5a,b) allowing different hypotheses for the functional state of this model. The global r.m.s.d. for the ADP source is 2.54 and for ATP 1.36 Å, while the local r.m.s.d. for ADP is 0.99 Å and for ATP 0.65 Å. This suggests that the structure is more representative of the ATP-bound state.

How do you think the roles of AI for scientific research in the future and what we can do to take advantage of it?

We need more creative thinking, which is beyond the existing knowledge and raise questions. Ask a question is more important than solving it.