



# APACE: AlphaFold2 and advanced computing as a service for accelerated discovery in biophysics

Hyun Park<sup>a,b,c</sup>, Parth Patel<sup>a,d,e</sup>, Roland Haas<sup>e</sup>, and E. A. Huerta<sup>f,f,g,1</sup>

Edited by Herbert Levine, Northeastern University, Boston, MA: received August 14, 2023; accepted December 25, 2023

The prediction of protein 3D structure from amino acid sequence is a computational grand challenge in biophysics and plays a key role in robust protein structure prediction algorithms, from drug discovery to genome interpretation. The advent of AI models, such as AlphaFold, is revolutionizing applications that depend on robust protein structure prediction algorithms. To maximize the impact, and ease the usability, of these AI tools we introduce APACE, AlphaFold2 and advanced computing as a service, a computational framework that effectively handles this AI model and its TB-size database to conduct accelerated protein structure prediction analyses in modern supercomputing environments. We deployed APACE in the Delta and Polaris supercomputers and quantified its performance for accurate protein structure predictions using four exemplar proteins: 6AWO, 6OAN, 7MEZ, and 6D6U. Using up to 300 ensembles, distributed across 200 NVIDIA A100 GPUs, we found that APACE is up to two orders of magnitude faster than off-the-self AlphaFold2 implementations, reducing time-to-solution from weeks to minutes. This computational approach may be readily linked with robotics laboratories to automate and accelerate scientific discovery.

Al for science | biophysics | supercomputing | automation

Innovation at the interface of AI and advanced computing is enabling breakthroughs in science and engineering (1–7). The rise of AI models such as GPT-4 (8), AlphaFold (9), among others, provides new capabilities to accelerate and automate scientific discovery. However, some of these models have not been released to the public, breaking a strong tradition in the AI community. It has been argued that the sheer size of these AI models prohibits their use by a large cross-section of potential users.

To address this shortcoming, we demonstrate how to combine large AI models with high-performance computing platforms to empower a broad cross-section of users to fully exploit the capabilities of AI for scientific discovery. We have selected AlphaFold2 (10) as the science driver for this study, since this AI model is revolutionizing discovery in biophysics, and its use for accurate and rapid protein structure prediction (PSP) demands an optimal use of modern supercomputing environments. Here, we demonstrate how to optimize AlphaFold2 and its database, which exceed 2.6 TB in data storage, to reduce the time needed for accurate PSPs from weeks to minutes.

AlphaFold2's Features. With a deep learning technique—based structure prediction, AlphaFold2 (9) showed unprecedented performance at the 14th Community Wide Experiment on the Critical Assessment of Techniques for Protein Structure Prediction (11), later improved further with multimer prediction (12). Ever since, efforts to make AlphaFold2 faster (13, 14), to predict protein complex, e.g., antibody (10, 15), and of Chicago C to sample diverse protein conformations (16–18) have come to fruition. In this study, we use AlphaFold2 version 2.3.0., as of August, 2023. The pre-trained neural network parameters used include both monomer and multimer v3.

AlphaFold2 utilizes Central Processing Units (CPUs) to compute key input features: multiple sequence alignment (MSA), and structural templates. MSA represents a collection of protein sequence homologues related to the query protein. MSA captures evolutionary relationships between various proteins such as conserved and variation amino acid residues. MSA is computed using CPU-based sequence alignment algorithms The authors declare no competing interest. such as Jackhmmer, which align a query protein sequence with known sequence homologues obtained from databases such as Uniclust. AlphaFold2 can glean key residue Copyright © 2024 the Author(s). Published by PNAS. interactions from MSA (9). Furthermore, structural templates refer to experimentally known protein homologue structures that share significant sequence similarity with the query protein. These template structures are used to improve the accuracy of AlphaFold2 so whom correspondence may be addressed. predictions. CPU-based algorithms, such as HHsearch (for monomer) and Hmmsearch (for multimer), search public protein structure databases such as the Protein Data Bank

# **Significance**

We introduce APACE, AlphaFold2 and advanced computing as a service, a computational framework that optimizes AlphaFold2 to run at scale in high-performance computing platforms, and which effectively handles this TB-size AI model and database. We showcase the use of APACE in the Delta and Polaris supercomputers to accelerate protein structure prediction for a variety of proteins, and demonstrate that using 200 ensembles distributed over 300 NVIDIA A100 GPUs, APACE reduces time-to-insight from days to minutes. This framework may be readily linked with self-driving laboratories to enable automated discovery at scale.

Author affiliations: <sup>a</sup>Data Science and Learning Division, Argonne National Laboratory, Lemont, IL 60439: <sup>b</sup>Theoretical and Computational Biophysics Group, Beckman Institute for Advanced Science and Technology, University of Illinois at Urbana-Champaign, Urbana, IL 61801; <sup>c</sup>Center for Biophysics and Quantitative Biology, University of Illinois at Urbana-Champaign, Urbana, IL 61801; <sup>d</sup>Department of Computer Science, University of Illinois at Urbana-Champaign, Urbana, IL 61801; <sup>e</sup>National Center for Supercomputing Applications, University of Illinois at Urbana-Champaign, Urbana, IL of Chicago, Chicago, IL 60637; and <sup>g</sup>Department of Physics, University of Illinois at Urbana-Champaign, Urbana, IL 61801

Author contributions: E.A.H. designed research: H.P. P.P., and R.H. performed research; H.P., P.P., and E.A.H. contributed new reagents/analytic tools; E.A.H. analyzed data; and H.P., P.P., and E.A.H. wrote the

This article is a PNAS Direct Submission.

This open access article is distributed under Creative Attribution-NonCommercial-NoDerivatives Commons icense 4.0 (CC BY-NC-ND).

elihu@anl.gov

Published June 24, 2024.

(PDB). Then, AlphaFold2 extracts spatial information from the template relevant to the query protein (9).

In the Graphics Processing Unit (GPU) phase, AlphaFold2 utilizes the features generated from MSA and templates, passing them through the evoformer network. Evoformer refines representations for both the MSA and pair interactions while iteratively exchanging information between them in a crisscross fashion to extract amino acid residue relationships. The updated representations then enter a structure module, where predictions for rotations and translations are made to position each residue (9).

The resulting predicted 3D structure undergoes a relaxation process via minimization by Molecular Dynamics (MD) engine to enhance accuracy. Upon generating the final structure, the information cycles back to the beginning of the evoformer blocks compute node consists of 1 AMD EPYC Milan processor, four in a recycle procedure, further refining the structure predictions. Overall, AlphaFold2 is trained end-to-end, leading to remarkable endpoints, and two NVMe SSDs. The system interconnect is accuracy and reliability in predicting protein 3D structures (9).

APACE's Improvement over AlphaFold2. We introduce APACE, AlphaFold2 and advanced computing as a service, a computational framework to accelerate AlphaFold2 through CPU & GPU optimizations, and distributed computing in supercomputing environments. Key features of this approach encompass:

Data management, First, APACE facilitates the usage of AlphaFold2's 2.6 TB AI model and database (9) by hosting it at the Delta and Polaris supercomputers (19). AlphaFold2's neural networks can readily access data by leveraging solid state drive (SSD) data storage, and Infinite Memory Engine (IME) data staging.

CPU optimization. Second, inspired by ref. 13, APACE uses Ray library's (20) CPU optimization to parallelize CPU intensive MSA and template computation calculations. As part of CPU optimization, APACE allocates higher CPU cores to MSA/template search tools, rather than default numbers (4 or 8), which showed heuristic speed improvement in our experiments. In addition, in a similar manner as ref. 13, we also implemented a checkpoint to circumvent redundant MSA/template steps if features.pkl file exists, i.e., an intermediate file storing MSA/template search result.

GPU optimization. Third, APACE uses Ray library's (20) GPU optimization to parallelize GPU intensive neural network protein structure prediction steps. An important key difference from ParaFold (13) in terms of GPU speedup is that, ParaFold predicted one conformation with only model 1 (a templatebased pretrained model) mostly on peptide sequences (e.g., with an average size less than 100 amino acid residues), rather than protein sequences (i.e. ~ 400 amino acid residues and more). In stark contrast, APACE predicts multiple protein conformations for each protein sequence, and peptide sequence if necessary, wit all five pretrained models in parallel, which is computationally demanding.

New functionalities. Fourth, APACE can predict multiple monomer conformations per pretrained neural network model (out of five models), a simple functionality existing only in multimer prediction in the original AlphaFold2 model. APACE includes functionalities such as enabling dropout during structure prediction, changing number of Evoformer (9) recycles, or subsampling MSA options, as provided in ref. 14.

## **Results and Discussions**

tailed comparison between APACE and the original AlphaFold2 we started Ray worker processes for the remaining compute

model. We describe each experiment at a time and then provide the corresponding results.

These results were obtained using the Delta and Polaris supercomputers, housed at the National Center for Supercomputing Applications, and at the Argonne Leadership Computing Facility (ALCF), respectively. Both machines provide highly capable GPU-focused compute environment for GPU and CPU workloads.

Delta offers a mix of standard and reduced precision GPU resources, as well as GPU-dense nodes with both NVIDIA and AMD GPUs. It also provides high-performance node-local SSD scratch file systems, as well as both standard Lustre and relaxed-POSIX parallel file systems spanning the entire resource. On the other hand, the Polaris supercomputer has 560 nodes. Each NVIDIA A100 GPUs, unified memory architecture, two fabric HPE Slingshot 11, and uses a Dragonfly topology with adaptive routing.

We compared APACE and AlphaFold2 performance using both NVIDIA A100 and A40 GPUs in Delta. and NVIDIA A100 GPUs in Polaris. The computational benchmarks we report below in terms of CPU and GPU runtimes were extracted from the generated timings ison file of both APACE and AlphaFold2.

Experiment 1: Predicting Structures for Four Benchmark Proteins. Four proteins were selected as benchmarks to assess the effectiveness and operational proficiency of APACE. To predict protein structures with APACE, we developed scientific software that enables users to provide suitable headers in sbatch scripts and to load the appropriate environment module, that are used to successfully submit and complete simulations in the Delta and Polaris supercomputers. These are the Simple Linux Utility for Resource Management (SLURM) parameters we used: --mem=240g, --nodes=10, --exclusive, --ntasks-per-node=1, --cpus-per-task=64, --gpus-per-task=4, --gpus-per-node=4. The neural network and MSA/templaterelated parameters were the same as AlphaFold2.

Monomers. We used the monomer protein 6AWO (serotonin transporter) as a basic structure to test baseline prediction accuracy and conformational diversity using a total of five models. Thus, we created a Ray cluster consisting of eight NVIDIA A100/A40 GPUs (equivalent to 2 A100/A40 GPU nodes in Delta and Polaris) to facilitate both CPU and GPU parallel execution and relaxation for all five models, i.e., one structure per model, as in the case of AlphaFold2.

Multimers. For multimer proteins, we tested 6OAN, Duffybinding protein bound with single-chain variable fragment antibody (15); 7MEZ, phosphoinositide 3-kinase (10); and 6D6U, a three distinct chain heteropentamer GABA transporter, which represents a more challenging case for multimer prediction. For each of these proteins, we had eight structure predictions per model, yielding a total of 40 predictions (five ensemble modes eight predictions per model). To facilitate concurrent execution and relaxation for the entire array of 40 models, 40 NVIDIA A100 and also 40 A40 GPUs (10 A100/A40 GPU nodes) were harnessed using a Ray cluster.

To initiate a Ray cluster utilizing compute nodes (as described in Methods), we first fetched a list of available compute nodes and their IP addresses. We then launched a head Ray process using We completed three computational experiments to carry out a de- one of these nodes, referred to as the "head node." Subsequently,

Table 1. Performance benchmarks between off-the-shelf framework for four exemplar proteins

# AlphaFold2 and our APACE CPU & GPU optimized

	# of	AlphaFold2-A40 AlphaFold2-A100		APACE-A40	APACE-A100	AlphaFold2-Polaris APACE-A100 Polaris	
Protein	ensembles	CPU/GPU [min]	CPU/GPU [min]	CPU/GPU [min]	CPU/GPU [min]	CPU/GPU [min]	CPU/GPU [min]
1. 6AWO	5	33.0/17.7	33.0/12.8	16.1/4.0	16.1/2.9	189.0/46.8	92.2/9.4
2. 6OAN	40	99.1/268.1	99.4/181.4	56.5/7.9	57.1/5.6	306.2/593.5	175.9/14.8
3. 7MEZ	40	100.7/3756.3	100.7/2339.4	58.0/100.1	58.8/63.0	556.2/3640.9	324.8/91.0
4. 6D6U	40	143.6/1528.7	143.8/786.9	89.3/72.2	89.4/35.1	485.9/1279.1	302.1/32.0

We present results for two types of GPUs available in the Delta supercomputer, NVIDIA A40 and A100 GPUs. We also present results using the Polaris supercomputer housed at the ALCF.

nodes. Each worker is equipped with all four GPUs, and Ray automatically determines the utilization of available GPUs for running and relaxation of models. The workers are then linked to the head node by providing the head node's address.

We utilize srun via message passing interface (MPI) to start the workers on the compute nodes. This is necessary because the sbatch script executes solely on the first compute node. Given the simultaneous launch of all Ray processes using MPI, we incorporate safeguards to prevent race conditions. The race condition safeguards ensure that the head node is started before the worker nodes and the beginning of predictions.

After the underlying Ray cluster was ready, we established a connection to it using ray.init within the run\_alphafold.py code and initiated the prediction of the protein structure using APACE. Ray automatically allocates resources and concurrently executes MSA tools on CPUs, model runs, and model relaxation on a distinct GPU. This operation efficiently harnesses the full computational potential of both CPU cores and GPUs available on the compute node.

### **Experiment 1: Results and Discussion.**

CPU acceleration. Through the implementation of parallel optimization techniques, APACE achieved an 1.8X average CPU speedup in Delta, and 1.78X average CPU speedup in Polaris. These results are independent of the number of compute nodes. GPU acceleration. APACE achieves significant GPU speedups. The following results were obtained using 8 GPUs for 6AWO, and 40 GPUs for 6OAN, 7MEZ, and 6D6U:

- 1. 6AWO. 4.4× speedup on both A40 and A100 GPUs for Delta; and 4.98× speedup on Polaris.
- 2. 6OAN. 34× and 32.4× speedup on A40 and A100 GPUs, respectively, for Delta; and 40.k speedup for Polaris.
- 3. 7MEZ. 37.5× and 37.1× speedup on A40 and A100 GPUs, respectively, for Delta; and 40X speedup for Polaris.
- 4. 6D6U. 21.2× and 22× speedup on A40 and A100 GPUs, respectively, for Delta; and 40x speedup for Polaris.

We summarize these results in Table 1. We also note that prediction times are consistently shorter when using NVIDIA A100 GPUs. In brief, APACE provides remarkable speedups for basic and complex structures, retaining the accuracy and

robustness of the original AlphaFold2 model. Furthermore, APACE can readily be used for analyses at scale using hundreds of GPUs, as shown below.

Experiment 2: Predicting Protein 7MEZ using 100 and 200 NVIDIA A100 GPUs. To quantify the performance and scalability of APACE in the Delta and Polaris supercomputers, we conducted protein 7MEZ predictions utilizing a significant number of compute nodes. Specifically, we utilized 100 NVIDIA A100 GPUs, which correspond to 25 A100 GPU compute nodes to generate predictions (20 predictions per model). Likewise, we leveraged the computational power of 200 NVIDIA A100 GPUs, equivalent to 50 A100 compute nodes to generate a total of 200 predictions (40 predictions per model). To predict the structures, the sbatch script was modified to allocate the correct number of compute nodes. We also modified the srun and singularity run parameters to successfully complete these calculations.

**Experiment 2:** Results and Discussion. APACE delivered remarkable speedups. If we compute 100 ensembles (distributed over 100 GPUs) for protein 7MEZ, APACE completed the required calculations within 67.8 min, as opposed to AlphaFold2's 6068.8 min (101.1 h/4.2 d) in Delta. In Polaris, we observe that APACE reduced time-to-solution from 8793.3 min (146.5 h/6.1 d) to 87.9 min.

Similarly, if we now require 200 ensembles for the same protein, APACE in Delta completed all predictions within 64 min, as opposed to the 12023.3 min (200.4 h/8.3 d) that would be needed using the original AlphaFold2 method. In Polaris, APACE reduced time-so-solution from 12741.2 min (212.4 h/8.8 d) to only 84.9 min.

Finally, using 300 ensembles for protein 7MEZ, APACE in Delta completed all predictions within 68.2 min, as opposed to the 18064.3 min (301.1 h/12.5 d) that would be needed using the original AlphaFold2 method. In Polaris, APACE reduced time-so-solution from 15295.6 min (254.9 h/10.6 d) to only 76.9 min. These results are summarized in Table 2.

**Experiment 3:** Ensemble Diversity of APACE. AlphaFold2's inherent limitations restrict us to generating merely five predictions per monomer, e.g., one prediction per model, thereby

Table 2. Performance benchmarks between off-the-shelf AlphaFold2 and APACE for protein 7MEZ

Nodes/ensembles	AlphaFold2-Delta CPU/GPU [min]	APACE-Delta CPU/GPU [min]	AlphaFold2-Polaris CPU/GPU [min]	APACE-Polaris CPU/GPU [min]
25/100	100.7/6068.8	58.8/67.8	556.2/8793.3	324.8/87.9
50/200	100.7/12023.3	58.8/64.0	556.2/12741.2	324.8/84.9
75/300	100.7/18064.3	58.8/68.2	556.2/15295.6	324.8/76.9

We present results for 25, 50, and 75 nodes in Delta and the Polaris supercomputers. Each node has 4 NVIDIA A100 GPU.

confining the diversity of protein conformation. Moreover, fine-tuning parameters such as dropout remain inaccessible. However, we successfully addressed this constraint by adapting the ColabFold (14) code. For experimental purposes, we generated 100 structures for protein 6AWO using --num\_multimer\_predictions\_per\_model=20 while employing the parameter --use dropout=True. This was accomplished by configuring the sbatch script with the appropriate parameters. APACE enables users to select the following options (14):

- 1. Ensemble of structure module with -num ensemble.
- 2. Control for recycles with --num recycles,
- Subsampling of MSA with --max seq, --max extra seq,
- Evoformer fusion with --use fuse,
- 5. Bfloat16 mixed precision with --use bfloat16.
- 6. Bernoulli-masking based diverse conformational sampling with --use dropout.

#### **Experiment 3: Results and Discussion.**

Protein structure prediction and conformational diversity by APACE. We have modified AlphaFold2 code to mirror ColabFold (14)'s versatile protein structure prediction pipeline parameter customization. With these improvements, we have successfully expanded the spectrum of predictions, thereby enhancing the overall reliability of the predicted structures. Although protein structure prediction is of great significance, we would like to expand APACE to predict conformational diversity since proteins are not static but malleable and flexible structures. Sampling a wide range of conformational ensemble is important for drug discovery (16–18). In the case of 6AWO (~500 amino acid residues), Fig. 1, we used our parameter customization enhancements (with the option --use dropout=True) and predicted 100 structures of serotonin transporter (SERT). We have found that the structure predicted by APACE is comparable to the ground truth structure. When we visualize most variant transmembrane domain alpha helices (cyan in the Right panel), we observe that TM2, TM6, TM10, and TM12 are highlighted. Among these, TM6, TM10, and TM12 are responsible for conformational change or ligand binding from outward-facing to inward-facing structures (22-24). This implies that APACE learned patterns to predict a wide range of

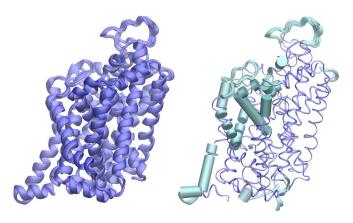


Fig. 1. Protein structure used to test APACE: serotonin transporter (PDB accession: 6AWO; shorthand SERT). The Left panel is 100 SERT predicted conformational ensemble overlaid, which has good agreement with ground truth SERT. The Right panel is high variant transmembrane domains, shown in cyan, and computed with root mean square fluctuations overlayed. Figures are generated with Visual Molecular Dynamics (21).

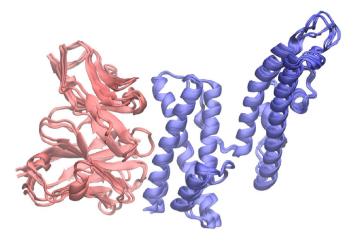


Fig. 2. Protein structure used to test APACE: the antibody-antigen complex Plasmodium vivax Duffy-binding protein (PDB accession: 6OAN). The structure has good agreement with ground truth bound structure conformation. The predicted conformational ensemble of complementary determining region (CDR; loops) of the antibody (red) binding against helical secondary structure epitopes of antibody (blue) are predicted well when compared to

conformational landscape of SERT. APACE makes an accurate prediction for SERT, which is an integral membrane protein. Even without the presence of membranes, APACE manages to predict transmembrane domains with high accuracy, hence demonstrating APACE's promise in drug discovery research. In the case of 6OAN,  $^{7}MEZ$ , and 6D6U ( $^{600}$ , 2,000, and 1,800 amino acid residues, respectively), we have multimer predictions. Both 6OAN and 7MEZ in Figs. 2 and 3 each predict conformational ensemble of heterodimer structures with high accuracy. Especially, the interface binding pose is well predicted and comparable with ground truth structures. Although there may be minor errors in predicted secondary structures not involved in interface binding, correct interface

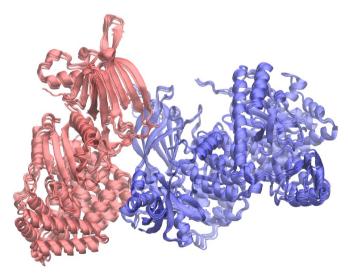


Fig. 3. Protein structure used to test APACE: a phosphoinositide 3-kinase (PĪ3K) consisting of p110  $\,\gamma$  and p101 subunits (PDB accession: structure has good agreement with ground truth bound structure conformation. Although there are mispredictions of loop secondary structures in p101 (red; Top Left helical loop; mispredicted as alpha helix rather than loop) γ (blue) is well subunit, the interface binding pose between p101 and p110 predicted, implying conserved binding interface in evolution. Also, rest of the secondary structures and overall heterodimer structure of the predicted conformational ensemble are comparable with ground truth structure

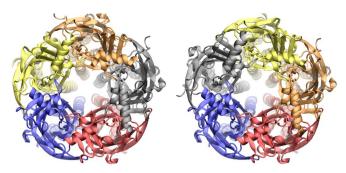


Fig. 4. Protein structure used to test APACE: a pentameric GABA A receptor (PDB accession: 6D6U). We show one predicted heteropentamer structure of neurotransmitter GABA A receptor. The Left panel shows a comparable structure with ground truth predictions. Blue and gray chains form a homodimer while red and orange chains form the other homodimer. Yellow chain is a monomer differing in sequence from other two homodimers. However, the location of transmembrane helices (toward the paper direction) does not exactly reproduce the ground truth structure. This is understandable since APACE does not use membrane as an input to predict the transmembrane domain. However, the overall structure is comparable with the ground truth. On the other hand, in the Right panel, we see an Al-predicted protein whose structure is erroneous, and where blue and gray chains bind to each other. This structure may have high thermodynamic instability and steric hindrance when being crystallized.

binding pose between proteins by nonbonded interactions is of greater importance. Minor misfolding may be addressed with methods such as MD, Monte Carlo, and protein design tools (25–38).

In 6D6U in Fig. 4, we observe comparable structure (*Left* panel) with ground truth and wrong structure (*Right* panel) with wrong homodimer location predictions. Since 6D6U is a membrane pentameric heteromer protein, it lends a challenging case of predicting not only correct structure of each monomer but also alternating chain patterns. The transmembrane helices are therefore mispredicted but overall structure is still comparable Limitations of AlphaFold2 Model. with ground truth.

In short, we have demonstrated APACE's capabilities to predict protein structures, mirroring AlphaFold2's robustness and accuracy, and providing remarkable speedups, reducing time-to-solution from days to minutes. APACE may be limited occasionally when it comes to predicting transmembrane proteins and/or multi-chain multimers, features it has inherited from AlphaFold2.

## **Methods**

Given that Delta and Polaris's containe support is only available or GitHub repositoryhich are intended for Docker containes (AM) we

- to the respective directorly code is available at s://github.com/ hyunp2/alphafold/tree/main.
- sequence representation, as well as for template similarity capturing synfacturesceed normal workstation storage capacity. representation. We also downloaded AlphaFold2 model parameterts/tempoirtyetRequirements: The model's memory footprint can be substa for AlphaFold2's functioning. Every database and model paramete(wistbpaltoger protein requiring higher GPU memory), making it challen date (as of August 2023) and the multimer version is v3.
  - Since Singularity is used on Delta and Polaris for contain we support with limited RAM. built the Singularity image by initially building the Docker image tocalhgle GPU Utilization: The original AlphaFold2 model is designed to Afterwardie pushed this Dockermage to Docker Hubbilizing the a single GPU during inference, limiting its capability to work with mulsingularity putbmmandve converted the Docker image to Singularit@PUs. As a result, it predicts and relaxes one protein structure (save sif format, making it compatible with Delta and Polaris environmentise format) at a time sequentially.

Running the Singularity image built using the default Dockerfile resul an HHsearch (used for template search against PDB database) runti To address this issue, modifications were made to the Dockerfile. Ini the Dockerfile involved cloning and compiling the HHsuite package fi source locally, which posed portability challenges across different ma The compilation process with cmake relied on the processor architecture the user's machine, potentially leading to compatibility issues. For install the user building the Docker image locally had a processor with an Ir SetArchitecture (ISIA) tdiffered from Delta's supported architecture, HHsearch encountered a runtime error with "Illegal Instruction."

To address this issue and ensure cross-machine commandative. modifications to the DockerFile. Instead of compiling HHsuite from so we adopted a differ**app**roach by installing HHsuite using a statically compiled version that supports the AVX2 ISA. This modification elimi the dependency on local processor architecture during the build proc mitigating the potential runtime errors and enhancing the portability of Singularity image.

- 4. To complete setup, we created an output lirectory defaultis /tmp/alphafold) and ensured that it had the necessary permissions to
- 5. Due to Delta and Polarisassence of Dockesupport the standard run\_docker.py scripts notiableInsteadive devised a custom shell script to replicate the essential functionality of run\_docker.py. Employ a singularity run command, effectively bound the necessary mounts and passed the required flags for exemilitioning the procedure of run\_alphafold.py with Docker.
- 6. Upon completion of the deployment process, the output directory cor the predicted structures of the target protein, accurately obtained thro AlphaFold2's advanced prediction capabilities.

In conclusion, deploying AlphaFold2 on Delta and Polaris required a modifications to account for Singularity containerization. Through this a we successfully integrated AlphaFold2's posterful olding prediction capabilities into these supercomputers' environments.

The origina Alpha Fold 2 mode while highly accurate in predicting protein structures bave some limitations in terms of computational efficiency. Some of the key limitations include

- 1. Long Inference Time: The time taken for the model to make prediction considerable, especially for larger and more complex protein structur can hinder its use in time-sensitive applications. Such long time infer including both CPU and GPU computations have been reported and elsewhere in the literature (13).
- 2. Computationallytensive/Limited Real-Time Prediction suriginal AlphaFold2 modelcomputationally demandinguiring significant computational sources and time for accurate predictions ample, MSA and template search have to be performed in CPU while GPU i utilized for each structure predibition, in a sequential anner This

Apptainer/Singularity (39), we modified the instructions provided in Alphatricter's applicability for real-time predictions or on hardware with computational power. The computational demands of the model may describe the steps followed to deploy AlphaFold2 on Delta and Polariseals) me prediction of protein structures, making it less suitable for ti sensitive applications.

1. We began by cloning the AlphaFold2 directory from DeepMind and Resolg ated Intensive: AlphaFold2's inference requires composition that tional resources, including powerful GPUs or tensor processing units may limit its accessibility to researchers or institution accounts to the many limit its accessibility to researchers or institution accounts to the many limit its accessibility to researchers or institution accessibility to research accessibility to research accessibility to research accessibility to research accessibility accessibility to research accessibility access

2. Next, we downloaded the necessary genetic databases for MSA catridutrienged hardware. Also, the storage of database amounts to 2.6 TB

to process multiple protein structures concurrently, particularly on ma

https://doi.org/10.1073/pnas.2311888121

6. Otherpotentialimitations include bate notlimited to little protein within the mode ensemble. When running AlphaFold2/ARASES, can

implemented optimizations in both CPU and GPU constriviation enhanced efficiency and performance.

Uniref90, clustered MGnify and small Bulk File Distribution (BFD) database. On PDB70 in the monomer case and Hmmsearch for template search against in Bapabilities. During GPU utilization, AlphaFold2 performs se Segres in the multimer case. For the multimer only, Jackhmmer basted thre prediction, which is one of the reasons why AlphaFold2 takes database parsing step exists as well.

To process a single quexhphaFold2 limits itstdf8 CPU cores for significantly impacting the overall runtime (see Limitations of Alpha Polarical Process. for bottlenecks for CPU computation).

in parallelAPACE significantly enhances the spaced construction. searches to run concurrently.

Additionallywe allocated 16 CPU cords each MSA search tool and template search tolalckhmmer Hblitsand HHSearch/Hmmsearch (monomer/multimer). By running all three MSA tools in parallel, utilizing a total of 48 CPU cores, we achieved a substantial tube in performance.

After increasing the number of CPU cores, we observed remarkable speed and accuracy of AlphaFold2, and which leverages enhancements in the MSA computation. However, beyond 20 CPUs goesso the uting to reduce time-to-insight from days to minutes. speed-up in MSA calculation plateaued. This observation unveiled the hattleneskinglished this by a) making an efficient use of as being related to input trieval atherthan CPU processirendering it input-bound his inherent nature of input-bound presents hurdles at a staging; b) optimizing CPU and GPU computing; an for straightforward parallelization (i.e., CPU multiprocessing or MPI) methods in scientific software to enable the prediction of Upon the completion of parallel computation for MSA, the template search and conformational ensemble of protein structures. These tools are multimer Jackhmmer-Uniprot, sequentially ensue.

key optimization is retired the entire datases Dyninimizing data retrieval time. Additionally, we leveraged the IME to stage the dataset files into the SSD cache within the /ime file system. This pre-staging allowed jobs to

swiftly access and utilize the required data. IME is a DataDirect Network acalusticand software Availability. within a high-performance computing (HPC) environment. The MSA and structural template search results acquired on CPUs are stored

in features.pland passed to the neunaltwork foprediction on GPUs. protein structure, a similar manner as rest. This optimization ensures efficient processing and avoids redundant computations. GPU optimization. AlphaFold2/APACE employs an ensetintelectural for protein structure prediction. The three models out of five make Geampaign and its Nationahtefor Supercomputing Applications.

the refer to refs. 9 and 18 of independent predictions made by each indivindual hetwork model Institutes of Health.

conformation diversity. Predicting correct, yet diverse protein conformation diversity. Predicting correct, yet diverse protein conformation diversity. is a significant task for drug discovery, partially addressed in refseth? each molties each molties collective predictions from each model in the model ensemble offer diverse final predicted 3D structures, e.g., plasme

Key Optimizations in APACE. To transcend the limitations of AlphaFoldan was partic protease causing malarian (1861), are crucial understand free energy landscaperoftein conformations and to identify important drug discovery targetptic binding pocketscontrast to AlphaFolde2,

CPU optimization. AlphaFold2 utilizes Jackhmmer to conduct MSA searches on APACE a capability to pneutiple monomer structures per the other hand, it employs HHBlits for MSA search on large BFD and the original lphaFold2 modeldesigned to use a single GPU during databases. In addition, AlphaFold2 utilizes HHSearch for template search caylainch does notake fulladvantage offeep learning's parallel

time tillcompletion (13) (Limitation Alpha Fold 2 Mode To expedite the GPU phase, we used the Ray library for GPU parallelization as well for Jackhmmer, four CPU cores for HHblits, and eight CPU cores for Hmmsearcheach ensemble model and its corresponding predictions are Given the vast database sizes (around 2.6 TB) and the considerable distinct GPWs for structure prediction. As a result, APACE can harnest access involved, the MSA search for a single prediction can take several foods; in paradicted by expediting the overall

Following prediction by each model, the corresponding structure und To expedite the CPU stage implemented an approach inspired by relaxation process in a sequenation in AlphaFold an approach the ParaFold (13). By orchestrating the three independent sequential Missing step, once more harnessed the power of the Ray library in APACEhrough this optimizatearch structure predicted by ensemble contrast to AlphaFold2, where UniRef90, MGnify, and BFD datasets workers is signed to an individual cated GPIacilitating parallel sequentially using Jackhmmer and/or HHblits (for large BFD), APA@@xatiplo@accessing. This enhancement by APACE has substantially r the Ray library (20) to simultaneously initiate three processes, enabliocests at timeontributing to the overaticeleration tife relaxation process.

We have introduced APACE, a framework that retains the the Delta and Polaris supercomputer systems' data storage and In the provided researchers with a sequentially ensue.

To further enhance speed in CPU intensive MSA computation, we made two computational framework that may be readily linked with robotic computational framework that may be readily linked with robotic computational framework that may be readily linked with robotic computational framework that may be readily linked with robotic computational framework that may be readily linked with robotic computational framework that may be readily linked with robotic computations.

The data and scientific software designed to facilitate fast data tiering between compute nodes and refibes ytater produce this work are availables at github.com/hyunp2/ alphafold/tree/main (41).

Additionally, we have incorporated a code check in our pipeline to circumvent and Development funding from Argonne National Laboratory, pr CPU-burdening MSA computations features.palready exists (as a by the Directooffice oscience the United States (U.S.) Department of result of storing features.pkl by successfully executing CPU computation at least Contract No. DE-AC02-06CH11357. E.A.H. was partially once for a given protein sequence), the pipeline skips the MSA and structure by NSF award OAC-2209898 research used resourche Afgonne template search and computation steps and proceeds directly to predict the Leadership Computing Facility, which is a Department of Energy (DOE) of Science User Facility supported under Contract DE-AC0**ℤ⊧06**CH1135 research used the Delta advanced computing and data resource which network models to predict the 3D structure of proteins. This ensemble approach the National Science Foundation (award OAC 2005572) a entails using multiple pretrained models with slight variations hyperstate new Delta is a joint efforttoe University to the University tou based on MSA (i.e., models 3 to 5) while the other two models (i.e. authorises acknowledge support from the Intestitionated of General deciding 1 to 2) also rely on templates. For details of how the five models difference settler National Astitutes of lealth under awards P41-GM104601, R24-GM145965, and R01-GM123455. The content is solely the respon-

The "--num\_multimer\_predictions\_per\_model" flag governs the numerathors and does not necessarily represent the official views of the

- Y. LeCun, Y. Bengio, G. Hinton, Deep learning. Nature 521, 436-444 (2015).
- E. A. Huerta et al., Enabling real-time multi-messenger astrophysics discoveries with deeprlatoritiong and transport. Nature 569, 141–145 (2019). Nat. Rev. Phys. 1, 600-608 (2019). 23. M. C. Chan, E. Procko, D. Shukla, Structural rearrangement of the serotonin transporter intr
- M. Krenn et al., On scientific understanding with artificial intelligence. Nat. Rev. Phys. 4, 76gate9nduced by thr276 phosphorylation. ACS Chem. Neurosci. 13, 933–945 (2022). (2022).
- S. Issue, A machine-intelligent world. Science 381, 136-137 (2023).
- S. Bianchini, M. Müller, P. Pelletier, Artificial intelligence in science: An emerging gen Stall @eBradden, R. Neutze, Advances and challenges in time-resolved macromolecular crystallo of invention. Res. Policy 51, 104604 (2022).
- K. Crawford, Atlas of AI: Power, Politics, and the Planetary Costs of Artificial Intelligence 26. M. E. Mäeots, R. I. Enchev, Structural dynamics: Review of time-resolved cryo-EM. Acta Cr (Yale University Press, New Haven, 2021). D Struct, Biol. 78, 927-935 (2022).
- 58-74 (2022)
- OpenAl (2023). GPT-4 technical report.
- J. Jumper et al., Highly accurate protein structure prediction with alphafold. Nature 596, 529,58001 (2019). (2021).
- alphafold2. Nat. Commun. 13, 1265 (2022). 11. J. Moult, K. Fidelis, A. Kryshtafovych, T. Schwede, M. Topf, Critical assessment of techniquaesequorks. Nat. Methods 17, 665-680 (2020).
- protein structure prediction, fourteenth round. CASP 14 Abstract Book (2020).
- 12. R. Evans et al., Protein complex prediction with alphafold-multimer. bioRxiv (2021). https://www.imulations. J. Chem. Inf. Model. 61, 901–912 (2021). biorxiv.org/content/10.1101/2021.10.04.463034v2 (Accessed 30 November 2023).
- 13. B. Zhong et al., "Parafold: paralleling alphafold for large-scale predictions" in International binding. Biochimica et Biophysica Acta (BBA)-Gen. Subj. 1863, 1560–1567 (2019).
- Computing Machine, New York, NY, 2022), pp. 1-9.
- 14. M. Mirdita et al., ColabFold: Making protein folding accessible to all. Nat. Methods 19, 8494682 rofolini et al., Integrative approaches in structural biology: A more complete picture from (2022).
- reveals accuracy determinants. Protein Sci. 31, e4379 (2022).
- alphafold. J. Chem. Theory Comput. 19, 4355–4363 (2023).
- typical human kinase domains. bioRxiv (2023). https://www.biorxiv.org/content/10.1101/2528.676 (2016). 21.550125v1 (Accessed 30 November 2023).
- 18. D. Sala, F. Engelberger, H. Mchaourab, J. Meiler, Modeling conformational states of protein modeling suite can do for you. Biochemistry 49, 2987–2998 (2010) alphafold. Curr. Opin. Struct. Biol. 81, 102645 (2023).
- USA, 2018), pp. 561-577.

22. J. A. Coleman et al., Serotonin transporter-ibogaine complexes illuminate mechanisms of

- 24. M. C. Chan, B. Selvam, H. J. Young, E. Procko, D. Shukla, The substrate import mechanism human serotonin transporter. Biophys. J. 121, 715-730 (2022).
- Science 373, eaba0954 (2021).
- J. Dean, A golden decade of deep learning: Computing systems & applications. Daedalls 3.51 Amann, D. Keihsler, T. Bodrug, N. G. Brown, D. Haselbach, Frozen in time: Analyzing
  - dynamics with time-resolved cryo-EM. Structure 31, 4-19 (2023). 28. M. Schmidt, Time-resolved macromolecular crystallography at pulsed X-ray sources. Int. J.
  - 29. F. Martín-García, E. Papaleo, P. Gomez-Puertas, W. Boomsma, K. Lindorff-Larsen, Compa
- 10. P. Bryant, G. Pozzati, A. Elofsson, Improved prediction of protein-protein interactions using olecular dynamics force fields in the essential subspace. PLoS One 10, e0121114 (2015). 30. J. K. Leman et al., Macromolecular modeling and design in Rosetta: Recent methods and
  - 31. D. Sala, A. Giachetti, A. Rosato, Insights into the dynamics of the human zinc transporter Zr

  - 32. D. Sala, A. Giachetti, A. Rosato, An atomistic view of the YiiP structural changes upon Zinc
  - Conference on High Performance Computing in Asia-Pacific Region Workshops (Association 30 or). Matsunaga, Y. Sugita, Use of single-molecule time-series data for refining conformationa
    - dynamics in molecular simulations. Curr. Opin. Struct. Biol. 61, 153-159 (2020).
- combination of individual techniques. Biomolecules 9, 370 (2019). 15. R. Yin, B. Y. Feng, A. Varshney, B. G. Pierce, Benchmarking alphafold for protein configuration. Computational methods for exploring protein conformations. Biochem. Soc. Train 1707-1724 (2020).
- 16. A. Meller, S. Bhakat, S. Solieva, G. R. Bowman, Accelerating cryptic pocket discovery86si6g Bussi, A. Laio, Using metadynamics to explore complex free-energy landscapes. Nat. Re **2**, 200–212 (2020).
- 17. B. Faezov, R. L. Dunbrack Jr., Alphafold2 models of the active form of all 437 catalytically economic Metable Mt. Sali, Comparative protein structure modeling using modeller. Curr. Protoc. Bioin
  - 38. K. W. Kaufmann, G. H. Lemmon, S. L. DeLuca, J. H. Sheehan, J. Meiler, Practically useful:
  - 39. G. M. Kurtzer, V. Sochat, M. W. Bauer, Singularity: Scientific containers for mobility of comp
- 19. J. Towns et al., XSEDE: Accelerating scientific discovery. Comput. Sci. Eng. 16, 62–74 (20PA)S One 12, e0177459 (2017).
  20. P. Moritz et al., "Ray: A distributed framework for emerging Al applications" in 13th USAOLD. Merkel, Docker: Lightweight Linux containers for consistent development and deployment. Symposium on Operating Systems Design and Implementation (OSDI 18) (USENIX Association, J. 2014, 2 (2014). 41. H. Park, P. Patel, R. Haas, E. A. Huerta, Data and Software from "APACE: AlphaFold2 as a
- 21. W. Humphrey, A. Dalke, K. Schulten, VMD: Visual molecular dynamics. J. Mol. Graph. 14(083438e) lerated discovery in biophysics". GitHub. https://github.com/hyunp2/alphafold/tree/m Deposited 20 April 2023.