

# Distributed prediction of protein structures with alphafold

**Glen Hocky**

ParslFest 2021

Department of Chemistry, NYU

(in collaboration with Parallel Works, Inc)

# ACKNOWLEDGEMENTS

## FUNDING AND SUPPORT



Hocky Research Group, Fall 2021



New York University  
Faculty of Arts and Sciences



Department of Energy  
SBIR DE-SC0019695

## COLLABORATORS

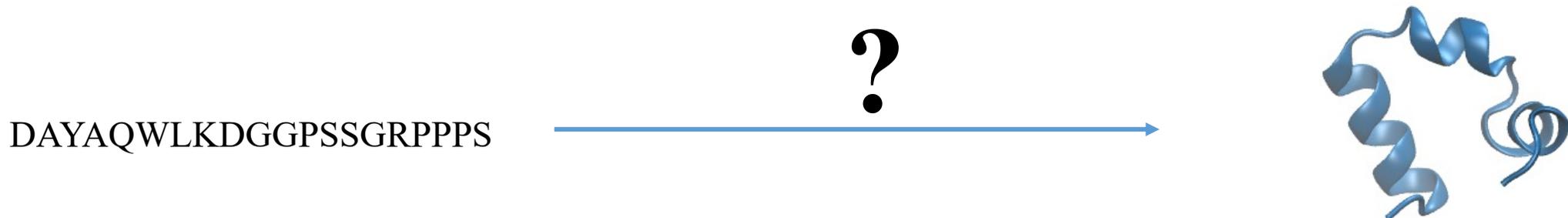
Parallel Works, Inc  
Shenglong Wang, NYU HPC

# BIG PICTURE

## Protein folding problem



## Structure prediction problem



# PAST WORK

## Protein folding with Swift (2007-2009), on BlueGene/P

### Towards petascale *ab initio* protein folding through parallel scripting

Glen Hocky<sup>1</sup>, Michael Wilde<sup>2,3\*</sup>, Joe DeBartolo<sup>4,5</sup>, Mihael Hategan<sup>2</sup>, Ian Foster<sup>2,3,6</sup>, Tobin R. Sosnick<sup>2,4,5\*</sup>, Karl F. Freed<sup>1,2,7\*</sup>

<sup>1</sup>Department of Chemistry, University of Chicago

<sup>2</sup>Computation Institute, University of Chicago & Argonne National Laboratory, USA

<sup>3</sup>Mathematics and Computer Science Division, Argonne National Laboratory, Argonne IL, USA

<sup>4</sup>Department of Biochemistry and Molecular Biology, University of Chicago

<sup>5</sup>Institute for Biophysical Dynamics, University of Chicago

<sup>6</sup>Department of Computer Science, University of Chicago, IL, USA

<sup>7</sup>James Franck Institute, University of Chicago

\*wilde@mcs.anl.gov, trsosnic@uchicago.edu, freed@uchicago.edu

#### Abstract

*Petascale computers allow scientists and engineers not only to address old problems better, but also to consider new methods and new problems. We report here on work that both applies new methods and tackles new problems in the area of structural biology. The project combines an efficient protein structure prediction algorithm implemented in the Open Protein System (OOPS) system with the Swift parallel scripting system to enable the rapid and flexible composition of OOPS components into parallel programs.*

As OOPS becomes more accurate and efficient, a number of related computational challenges emerge in our desire to tackle proteins of increasing size because current prediction methods have limited accuracy even for proteins on the order of 100 residues when homology-based information is minimal. To predict the structures of larger and multi-domain proteins, statistical sampling becomes a limiting factor, and thus we require significantly more computing resources.

## Parallel scripting for applications at the petascale and beyond

Authors Michael Wilde, Ian Foster, Kamil Iskra, Pete Beckman, Zhao Zhang, Allan Espinosa, Mihael Hategan, Ben Clifford, Ioan Raicu

Publication date 2009/11/13

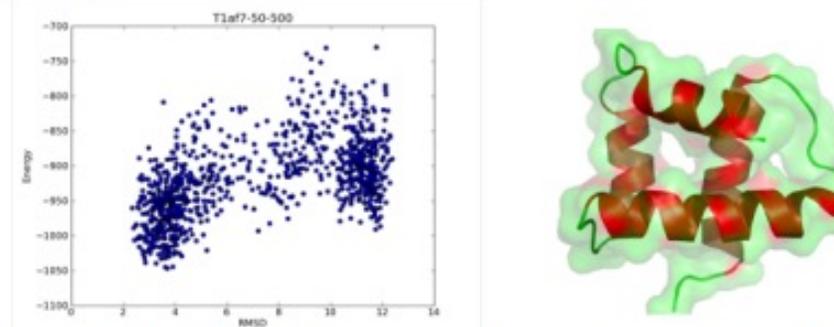
#### Display Swift Output:

Choose run:  
1093auadr1

key	AverageRunTime	LowestRMSD	LowestRMSDEnergy	PredictionEnergy	PredictionRMSD	TotalRuns
T1af7-50-500	50995	2.34647	-960.561	-1047.0	3.70194	985
T1b72-50-500	36777	2.10575	-469.571	-568.818	3.20041	1012
T1dcj-50-500	12617	6.30753	-2584.0	-5591.67	8.47168	995
T1di2-50-500	13484	3.4034	-5281.69	-8442.19	5.78596	1005
T1mkj-50-500	12129	5.61484	-3885.61	-4729.66	8.55371	1004
T1r69-50-500	34001	2.23661	-592.601	-664.585	7.72925	998
T1tif-50-500	9050	4.1065	-4994.31	-6344.37	9.34496	1012
T1ubq-50-500	14857	4.84337	-4645.08	-7223.23	9.56997	991

W3C XHTML 1.0 W3C CSS

This page last updated on Sun, 12 Apr 2009 22:02



**Figure 6** – Results of running eight proteins on 2 racks (8192 CPUS) on Argonne's BG/P, Intrepid. Below are results from this investigation for T1af7. On the left is a scatter plot showing the correlation between our statistical energy potential and accuracy of the protein structures for the 985 simulations that ran to completion. On the right is an image showing the lowest RMSD structure. This table, plot and image were all automatically generated by our scripting mechanism, and the table is presented by a simple CGI script at our web site [2].

# ALPHAFOLD ADVANCE

NEWS | 30 November 2020

## 'It will change everything': DeepMind's AI makes gigantic leap in solving protein structures

Google's deep-learning program for determining the 3D shapes of proteins stands to transform biology, say scientists.

## AlphaFold Is The Most Important Achievement In AI— Ever



**Rob Toews** Contributor ⓘ

AI

*I write about the big picture of artificial intelligence.*

Follow

### Article

## Highly accurate protein structure prediction with AlphaFold

<https://doi.org/10.1038/s41586-021-03819-2>

Received: 11 May 2021

Accepted: 12 July 2021

Published online: 15 July 2021

Open access

Check for updates

John Jumper<sup>1,4</sup>✉, Richard Evans<sup>1,4</sup>, Alexander Pritzel<sup>1,4</sup>, Tim Green<sup>1,4</sup>, Michael Figurnov<sup>1,4</sup>, Olaf Ronneberger<sup>1,4</sup>, Kathryn Tunyasuvunakool<sup>1,4</sup>, Russ Bates<sup>1,4</sup>, Augustin Žídek<sup>1,4</sup>, Anna Potapenko<sup>1,4</sup>, Alex Bridgland<sup>1,4</sup>, Clemens Meyer<sup>1,4</sup>, Simon A. A. Kohl<sup>1,4</sup>, Andrew J. Ballard<sup>1,4</sup>, Andrew Cowie<sup>1,4</sup>, Bernardino Romera-Paredes<sup>1,4</sup>, Stanislav Nikolov<sup>1,4</sup>, Rishub Jain<sup>1,4</sup>, Jonas Adler<sup>1</sup>, Trevor Back<sup>1</sup>, Stig Petersen<sup>1</sup>, David Reiman<sup>1</sup>, Ellen Clancy<sup>1</sup>, Michal Zielinski<sup>1</sup>, Martin Steinegger<sup>2,3</sup>, Michalina Pacholska<sup>1</sup>, Tamas Berghammer<sup>1</sup>, Sebastian Bodenstein<sup>1</sup>, David Silver<sup>1</sup>, Oriol Vinyals<sup>1</sup>, Andrew W. Senior<sup>1</sup>, Koray Kavukcuoglu<sup>1</sup>, Pushmeet Kohli<sup>1</sup> & Demis Hassabis<sup>1,4</sup>✉

### RESEARCH

#### RESEARCH ARTICLE

#### PROTEIN FOLDING

## Accurate prediction of protein structures and interactions using a three-track neural network

Minkyung Baek<sup>1,2</sup>, Frank DiMaio<sup>1,2</sup>, Ivan Anishchenko<sup>1,2</sup>, Justas Dauparas<sup>1,2</sup>, Sergey Ovchinnikov<sup>3,4</sup>, Gyu Rie Lee<sup>1,2</sup>, Jue Wang<sup>1,2</sup>, Qian Cong<sup>5,6</sup>, Lisa N. Kinch<sup>7</sup>, R. Dustin Schaeffer<sup>6</sup>, Claudia Millán<sup>8</sup>, Hahnbeom Park<sup>1,2</sup>, Carson Adams<sup>1,2</sup>, Caleb R. Glassman<sup>9,10,11</sup>, Andy DeGiovanni<sup>12</sup>, Jose H. Pereira<sup>12</sup>, Andria V. Rodrigues<sup>12</sup>, Alberdina A. van Dijk<sup>13</sup>, Ana C. Ebrecht<sup>13</sup>, Diederik J. Opperman<sup>14</sup>, Theo Sagmeister<sup>15</sup>, Christoph Buhlheller<sup>15,16</sup>, Tea Pavkov-Keller<sup>15,17</sup>, Manoj K. Rathinaswamy<sup>18</sup>, Udit Dalwadi<sup>19</sup>, Calvin K. Yip<sup>19</sup>, John E. Burke<sup>18</sup>, K. Christopher Garcia<sup>9,10,11,20</sup>, Nick V. Grishin<sup>6,7,21</sup>, Paul D. Adams<sup>12,22</sup>, Randy J. Read<sup>8</sup>, David Baker<sup>1,2,23\*</sup>

# ALPHAFOLD ADVANCE

NEWS | 30 November 2020

## 'It will change everything'

### DeepMind's AI makes a

### gigantic leap in solving

### protein structures

Google's deep-learning program for determining protein structures stands to transform biology, say scientists.

## AlphaFold Is The Most Important Achievement In AI Ever



**Rob Toews** Contributor ⓘ

AI

*I write about the big picture of artificial intelligence.*

### Article

## High-accuracy protein structure prediction

deepmind / alphafold Public

<> Code Issues 60 Pull requests 15 Actions Projects Security Insights

main 1 branch 2 tags

Go to file

Add file

Code

	Augustin-Zidek and Copybara-Service Accept any ordering given by ListDir i...	1d43aaf on Sep 10	25 commits
alphafold	Skip obsolete PDB templates that don't have a replacement.		2 months ago
docker	Fix a few typos.		3 months ago
imgs	Initial release of AlphaFold.		3 months ago
notebooks	Fix TensorFlow versions in AlphaFold Colab notebook.		2 months ago
scripts	Remove a redundant space.		3 months ago
.dockerignore	Collapse hh-suite install steps into single layer.		3 months ago
CONTRIBUTING.md	Initial release of AlphaFold.		3 months ago
LICENSE	Initial release of AlphaFold.		3 months ago
README.md	Update the bibtex citation with the issue number and pages		2 months ago
requirements.txt	Switch to Tensorflow CPU-only. GPU not needed for data pipeline.		3 months ago
run_alphafold.py	Use pLDDT in the B-factor column of the output PDBs.		2 months ago
run_alphafold_test.py	Accept any ordering given by ListDir in the assert.		2 months ago
setup.py	Use tensorflow-cpu in setup.py as well.		2 months ago

Alexander Pritzel<sup>1,4</sup>, Tim Green<sup>1,4</sup>, Michael Figurnov<sup>1,4</sup>, asuvunakool<sup>1,4</sup>, Russ Bates<sup>1,4</sup>, Augustin Židek<sup>1,4</sup>, Clemens Meyer<sup>1,4</sup>, Simon A. A. Kohl<sup>1,4</sup>, Bernardino Romera-Paredes<sup>1,4</sup>, Stanislav Nikolov<sup>1,4</sup>, Back<sup>1</sup>, Stig Petersen<sup>1</sup>, David Reiman<sup>1</sup>, Ellen Clancy<sup>1</sup>, Michalina Pacholska<sup>1</sup>, Tamas Berghammer<sup>1</sup>, Oriol Vinyals<sup>1</sup>, Andrew W. Senior<sup>1</sup>, Koray Kavukcuoglu<sup>1,4</sup>

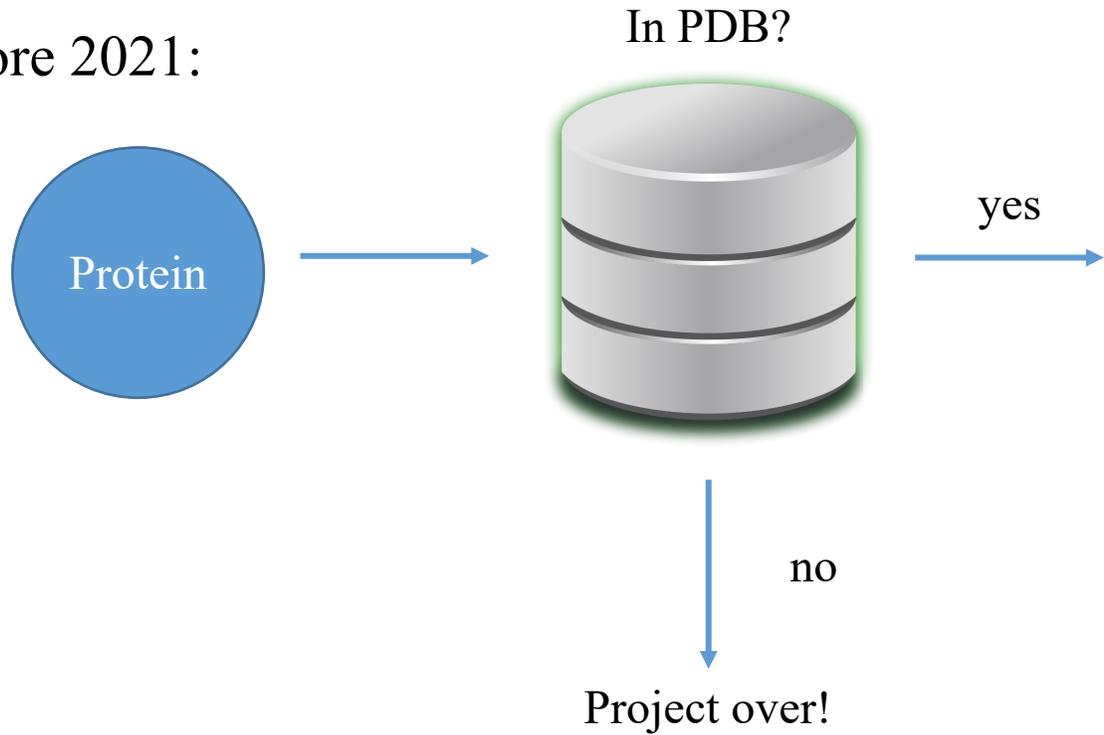
## Protein structures and neural network

Dauparas<sup>1,2</sup>, Sergey Ovchinnikov<sup>3,4</sup>, Justin Schaeffer<sup>6</sup>, Claudia Millán<sup>8</sup>, Andy DeGiovanni<sup>12</sup>, Jose H. Pereira<sup>12</sup>, Diederik J. Opperman<sup>14</sup>,

Theo Sagmeister<sup>15</sup>, Christoph Buhheller<sup>15,16</sup>, Tea Pavkov-Keller<sup>15,17</sup>, Manoj K. Rathinaswamy<sup>18</sup>, Udit Dalwadi<sup>19</sup>, Calvin K. Yip<sup>19</sup>, John E. Burke<sup>18</sup>, K. Christopher Garcia<sup>9,10,11,20</sup>, Nick V. Grishin<sup>6,7,21</sup>, Paul D. Adams<sup>12,22</sup>, Randy J. Read<sup>8</sup>, David Baker<sup>1,2,23\*</sup>

# WHY ARE WE EXCITED?

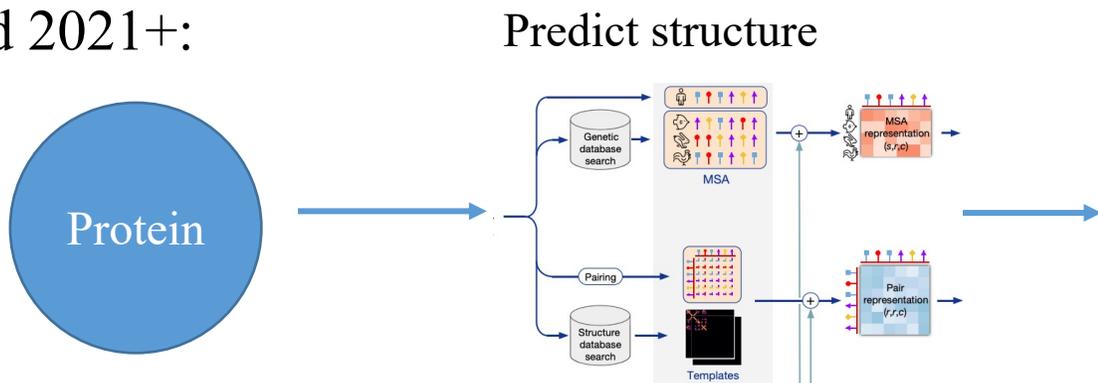
Before 2021:



Do simulation



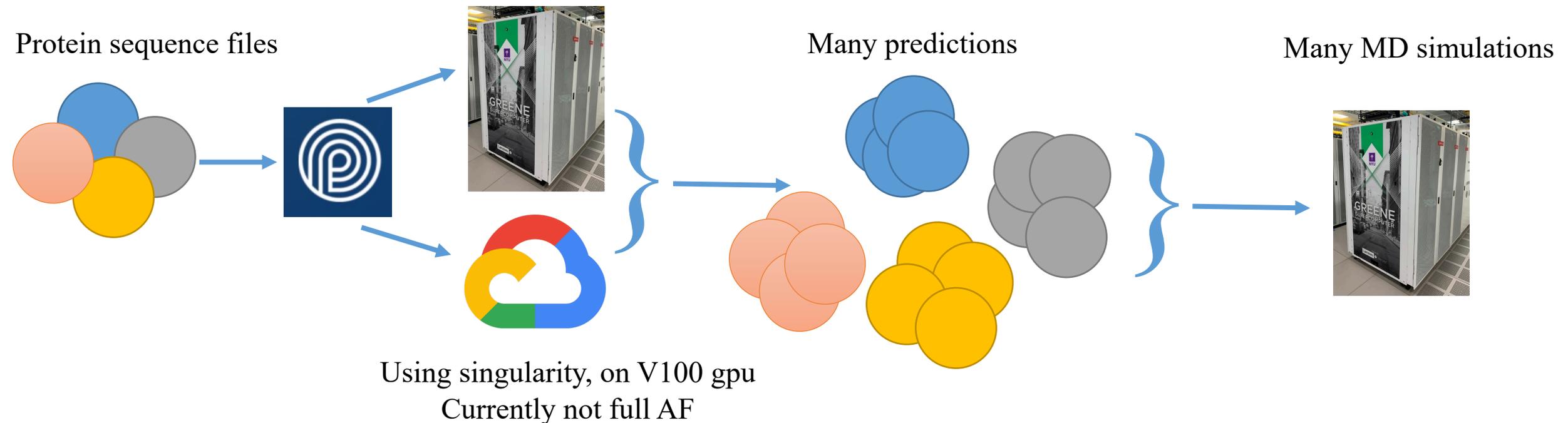
Mid 2021+:



# CHALLENGES

- Alphafold and similar software are computationally expensive
- Hard to install and require a lot of data
- Not high-throughput “out of the box”

## Parallel workflow w/ parsl:



# PW IN AND OUT

ALPHAFOLD

**Case Name**  **Cloud (True) / (False)**  True  False **Number of parallel seeds**

**Fasta files (.fasta)**  
 Select/Unselect all

/storage ▼

**Workflow Resource (GCP\_STRUCTURE\_PREDICT\_GPU) Not Started. [Please Start Selected Resource on Main Compute Page.](#)**

```
95     send_files = [Path(f) for f in glob.glob("*.sh")+glob.glob("*.py")+glob.glob("*.tcl")]
96
97     for run_file_name in run_file_names:
98
99         run_file = Path(run_file_name)
100        for i in range(1,n_seeds+1):
101            r = run_alphafold(
102                runscript=runscript,
103                random_seed=i,
104                inputs = [run_file]+send_files,
105                outputs=[out_dir,Path("af.stdout"),Path("af.stderr")])
106            runs.append(r)
107
108        print("Running", len(runs), "alphafold executions...")
109        [r.result() for r in runs]
110
111    gen_table(pwargs.outcsv,pwargs.outhtml)
```

# PW IN AND OUT

## Parametric Variables

- label
- model
- seed
- is\_relaxed
- rg
- plddt

## Cases

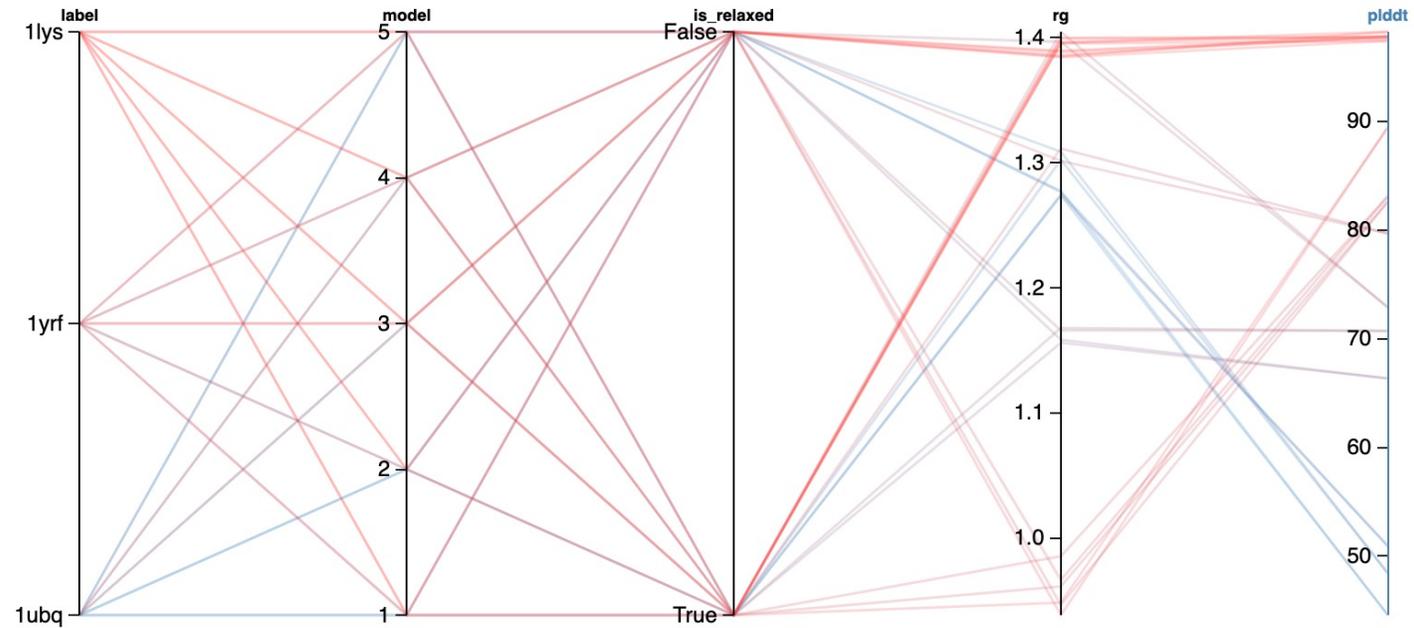
label  
1yrf 1lys

model  
1 5

seed: 1

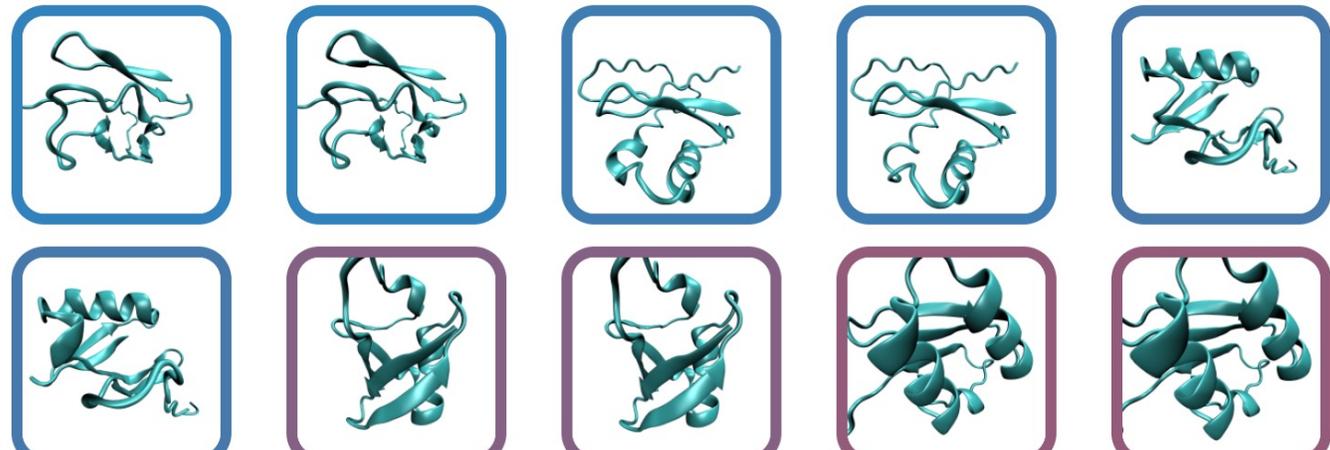
is\_relaxed  
True False

Reset Selection Exclude Selection Zoom to Selection Save Selection to File Thumbnail Size L M S [Open in New Window](#)



Sort by: plddt

Metric: model



# FINAL THOUGHTS

- Use of containers and platform specific bash scripts make parallel script more platform agnostic (but is there a place for app definitions per site?)
- Use of PW platform enables use of alphafold and visualization of results for total novice, using local or cloud resources (but see ongoing work on google colab notebooks for non-high-throughput cases)
- Farming out predictions may be good case for funcX, but composing with other functions like generating MD inputs or viz files still good case for parsl

