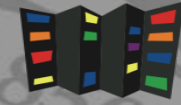


~~Protein~~ Biomolecular structure prediction with AI Alphafold and friends

Martin Čuma
Center for High Performance Computing
University of Utah
m.cuma@utah.edu



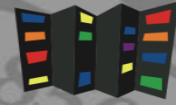
- Why AI structure prediction
- How it works?
- Some history (incl. Nobel Prize)
- Tools we have available
- Efficient use of the tools at CHPC



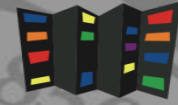
Biomolecular structure prediction



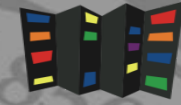
- Proteins are one of the base building blocks of life
- They form 3D structure which is affected by the amino acid sequence and surroundings
- Protein structure prediction
 - experimental - X-ray crystallography, NMR spectroscopy - labor intensive, expensive
 - computational
 - comparative - assembly from known smaller structures
 - ab-initio - physical principles structure minimization (molecular dynamics, ...)



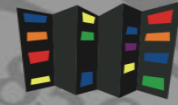
- requires Multiple Sequence Aligned (MSA) input
 - identifies relationships between sequences
 - first step, usually runs only on CPUs and uses large databases (is I/O intensive)
- use deep neural network trained on known (protein) structures
 - runs on GPUs, or CPUs (very slowly)
 - larger sequences need GPUs with more memory
- AI predicted structure may follow up with MD minimization



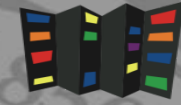
- AI revolutionized structure prediction by
 - much better accuracy than previous methods (70% for AlphaFold2)
 - easier to use and quicker results
- 2024 Nobel Prize for Chemistry
 - D. Baker (U Wash) - RosettaFold - "for computational protein design"
 - D. Hassabis, J. Jumper (Google DeepMind) - AlphaFold - "for protein structure prediction"



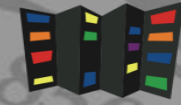
AI structure prediction programs available



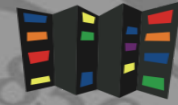
- Alphafold 2, 3, <https://github.com/google-deepmind/alphafold3>
 - be careful about the MSA performance
- Colabfold(batch), <https://github.com/sokrypton/ColabFold>
 - more efficient MSA, Alphafold for the inference
- Boltz, <https://github.com/jwohlwend/boltz>
 - fully open source AF alternative, uses Colabfold server for MSA
 - [Imi4boltz](#) - lower GPU memory alternative
 - [Boltzgen](#) - binder design tool



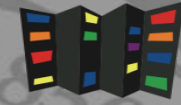
- RFAntibody, <https://github.com/RosettaCommons/RFAntibody>
 - workflow must be run in a container
- RFDiffusion, <https://github.com/RosettaCommons/RFdiffusion>
 - trickier to set up due to user space requirements
- Google Colab notebooks
 - RFDiffusion, Alphafold, ProteinMPNN
 - can run on CHPC resources
- RosettaFold All-Atom, <https://github.com/baker-laboratory/RoseTTAFold-All-Atom>
 - open source alternative to Alphafold 3, requires setup in user space



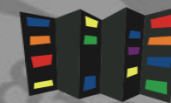
Performance considerations



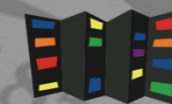
- AI enabled prediction runs in 2-3 steps
- Step 1 - MSA
 - Mostly on CPU (except for mmseqs-gpu), very I/O intensive
 - it's worth to have the databases or at least their indices in RAM
 - even better - use MSA server
- Step 2 - AI structure inference
 - Much more efficient on GPUs
 - Can be faster than the MSA search
- Step 3 (optional, AF2) - MD structure refinement
 - GPU or CPU, GPU faster for larger structures



- Alphafold uses more accurate but much slower MSA program
 - AF2 uses indexed databases - indices in the RAM disk (~30 GB)
 - AF3 does not = all databases on VAST file system
 - ~ 30 min for reference protein, AF3 a bit slower
 - for that reason 2 jobs, one CPU job for MSA, one GPU job for inference
- Colabfold, Boltz
 - use MSA server which runs mmseqs
 - CHPC servers fit the important databases in RAM (700 GB)
 - servers store past MSAs so it can reuse them and return result quicker
 - ~ 5 min for reference protein
 - just a single job on a GPU is usually OK, but split for large sequences



- CHPC has GPUs from many different generations
 - Their performance and capabilities vary widely
 - https://www.chpc.utah.edu/documentation/guides/gpus-accelerators.php#gpu_types
- (Nvidia) GPU classification:
 - Generation (code name - Maxwell, Pascal, Volta, Turing, Ampere, Hopper, Blackwell)
 - Compute Capability (5.2, 6.0, 6.1, 7.0, 7.5, 8.0, 8.6, 8.9, 9.0, 12.0)
<https://developer.nvidia.com/cuda-gpus#compute>
 - Theoretical compute throughput (single, double precision, tensor)
https://en.wikipedia.org/wiki/Nvidia_Tesla
<https://en.wikipedia.org/wiki/Quadro>
<https://en.wikipedia.org/wiki/GeForce>
https://en.wikipedia.org/wiki/GeForce_40_series
 - Amount of memory (~10-80 GB)



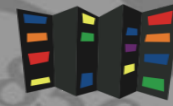
- Example: A40, A100, H100
 - Data center GPUs, high double precision performance, large memory
 - Expensive (H100 ~\$20000 w/ edu discount), hard to get
 - Good for simulations that need high numerical precision (engineering), or high memory (AI)

Model	Micro-architecture	Launch	Core	Core clock (MHz)	Shaders			Memory				Processing power (GFLOPS) ^[a]			CUDA compute capability ^[b]	TDP (W)	
					CUDA cores (total)	Base clock (MHz)	Max boost clock (MHz) ^[c]	Bus type	Bus width (bit)	Size (GB)	Clock (MT/s)	Bandwidth (GB/s)	Half precision Tensor Core FP32 Accumulate	Single precision (MAD or FMA)			Double precision (FMA)
A40 GPU accelerator (PCIe card) ^[43]		October 5, 2020	1× GA102	—	10,752	1,305	1,740	GDDR6	384	48	7,248	695.8	149,680	37,420	1,168	8.6	300
A100 GPU accelerator (PCIe card) ^{[44][45]}		May 14, 2020 ^[46]	1× GA100-883AA-A1	—	6,912	765	1410	HBM2	5,120	40 or 80	1,215	1,555	312,000	19,500	9,700	8.0	250
H100 GPU accelerator (PCIe card) ^[47]	Hopper	March 22, 2022 ^[48]	1× GH100 ^[49]	—	14,592	1,065	1,755 CUDA 1620 TC	HBM2E	5120	80	1,000	2,039	756,449	51,200	25,600	9.0	350
H100 GPU accelerator (SXM card)				—	16,896	1,065	1,980 CUDA 1,830 TC	HBM3	5,120	80	1,500	3,352	989,430	66,900	33,500	9.0	700



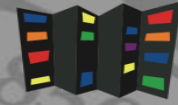
- Example: RTX6000, A6000
 - Mid size models (graphical workstations), more memory than the gaming graphics cards, price about 1/2 to 1/3 of the data center GPUs
 - Very low double precision performance, good single and tensor performance
 - Good for lower precision numerical calculations, AI with small to medium size models

Quadro GPU	Launch	Core	Core clock	Memory clock	Memory size (GB)	Memory type	Memory bandwidth	CUDA cores	Tensor cores	RT cores	Half precision	Single precision	Double precision	CUDA Compute Capability
Units			MHz	MHz	GB		GiB/s				TFLOPS	TFLOPS	GFLOPS	
RTX 6000 Ada Generation ^[202]	2022-12-03	AD102-870	915–2505	2500	48	384-bit GDDR6	960	18176	568	142	91.06 ^[203]	91.06	1423	8.9
RTX A6000 ^{[185][186]}	2020-10-05	GA102-875	1410–1800	2000	48 (96 with NVLink 3.0)	384-bit GDDR6	768	10752	336	84	38.709 ^[187]	38.709	1209.677	8.6



- Example: RTX3090
 - Consumer grade graphics cards, recent don't fit to compute servers
 - Very low double precision performance, good single and tensor performance
 - More affordable (a few thousand \$)
 - Good for lower precision numerical calculations, AI with small to medium size models

Model	Launch	Launch MSRP (USD)	Code name(s) ^[b]	Transistors (billion)	Die size (mm ²)	Core config ^[c]	SM count ^[d]	L2 cache (MB)	Clock speeds ^[e]		Fillrate ^{[f][g]}		Memory				Processing power (TFLOPs) ^[h]				TDP (watts)
									Core (MHz)	Memory (GT/s) ^[i]	Pixel (Gpx/s)	Texture (Gtex/s)	Size (GB)	Bandwidth (GB/s)	Type	Bus width (bit)	Half (boost)	Single (boost)	Double (boost)	Tensor compute [sparse]	
GeForce RTX 3090 Ti ^{[40][65]}	Mar 29, 2022	\$1,999	GA102-350	28.3	628.4	10752 336:112:84:336	84	6	1560 (1860)	10.5	174.7 (208.3)	524.1 (625)	24	1008	GDDR6X	384	33.55 (39.99)	33.55 (39.99)	0.524 (0.625)	160 ^[66] [320]	450
GeForce RTX 3090 ^{[40][63]}	Sep 24, 2020	\$1,499	GA102-250 ^[64] GA102-300	28.3	628.4	10496 328:112:82:328	82	6	1395 (1695)	9.75	156.2 (189.8)	457.6 (556)	24	936	GDDR6X	384	29.28 (35.58)	29.28 (35.58)	0.458 (0.556)	142 ^[34] [284]	350



Variables

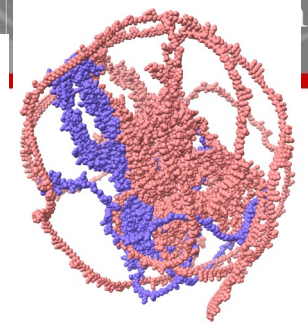
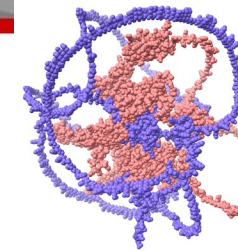
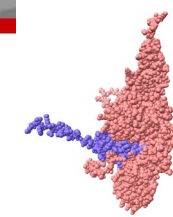
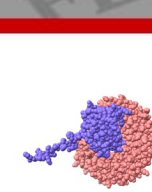
- Protein size
- GPU type
- GPU availability
- Job parallelization

Calculation times

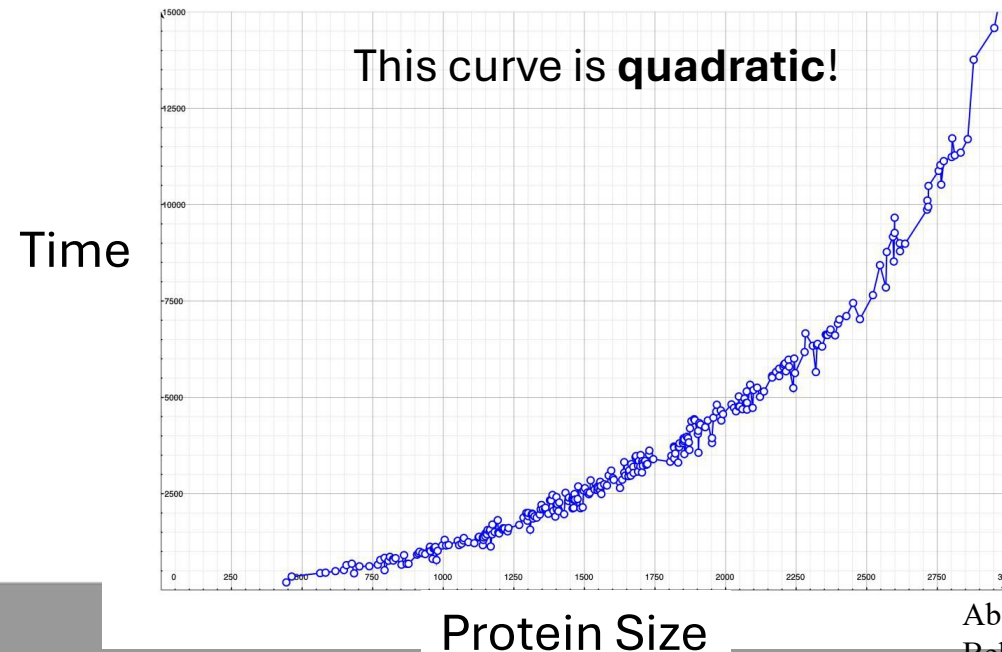


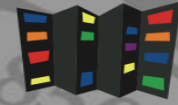
Variables

- Protein size
- GPU type
- GPU availability
- Job parallelization



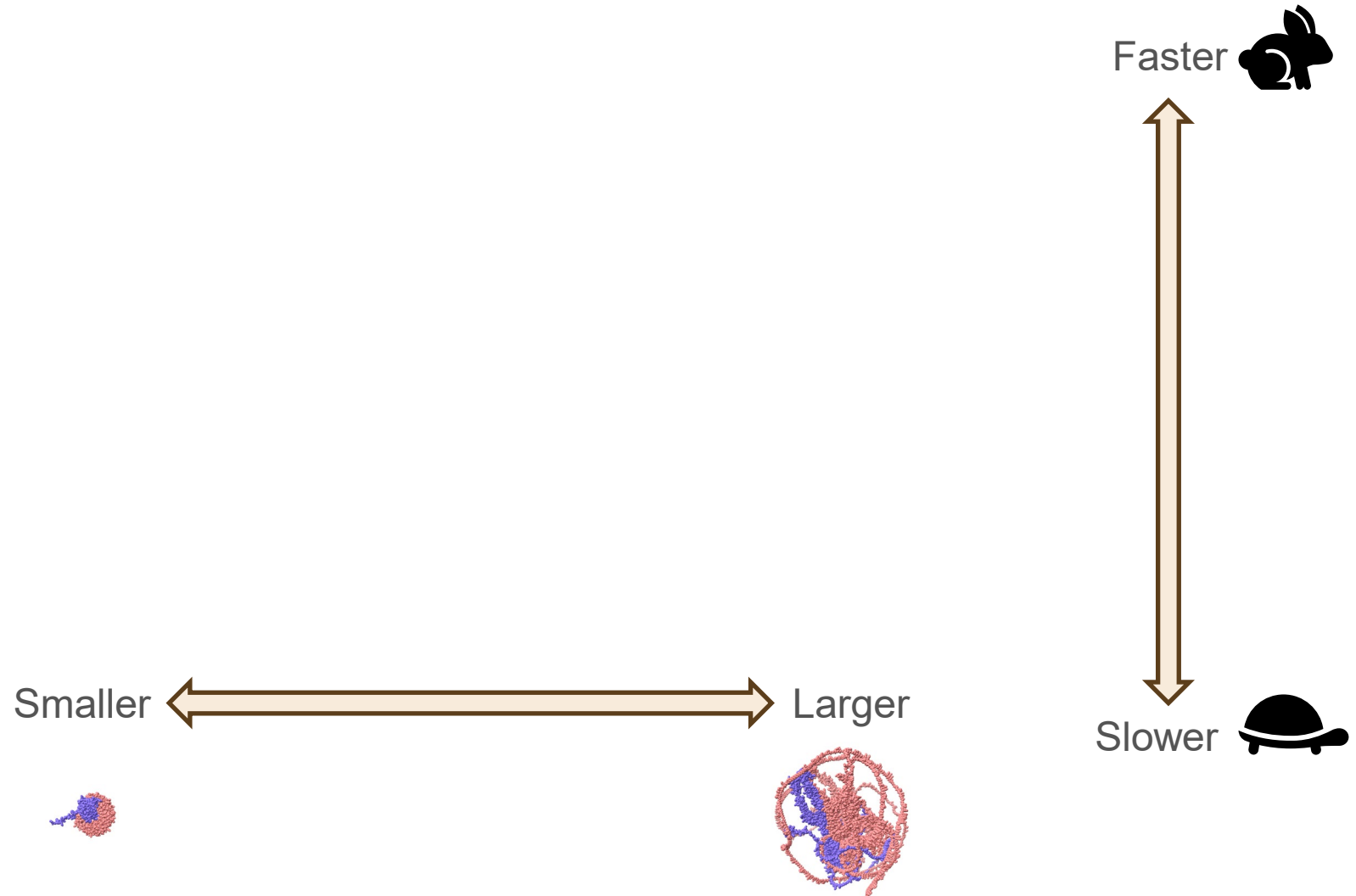
Protein 1	mRnase1	Insulin	Bnl	DI
Protein 2	mRnh1	InsRECD	Btl	N
# Amino acids	606aa	1040aa	1829aa	3537aa
Runtime (3090 GPU)	7 min	24 min	75 min	828 min

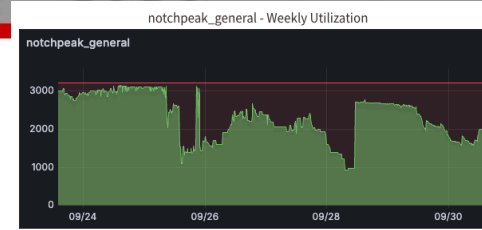




Variables

- Protein size
- **GPU type**
- GPU availability
- Job parallelization

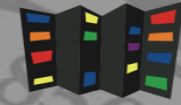




Variables

- Protein size
- GPU type
- **GPU availability**
- Job parallelization

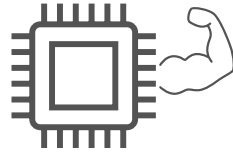




Variables

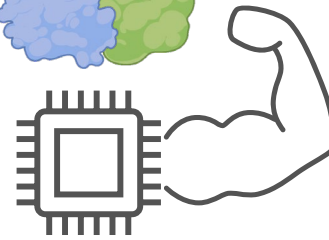
- Protein size
- GPU type
- GPU availability
- **Job parallelization**

Small proteins

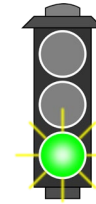


e.g. 2080ti, a30

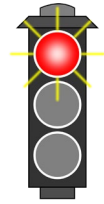
Large proteins



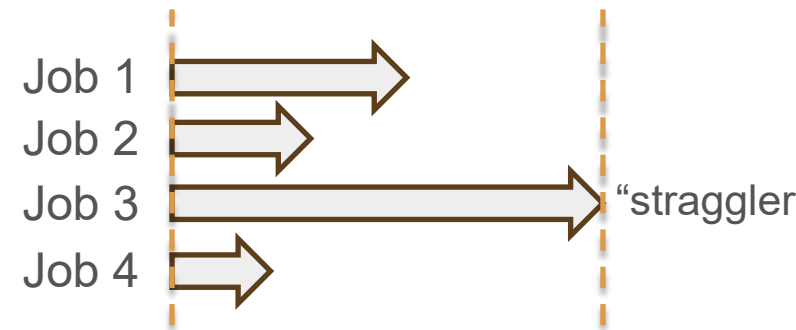
e.g. h100, a100

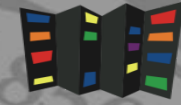


Screen start

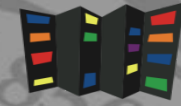


Screen end

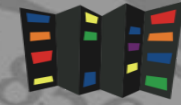




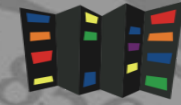
Biomolecular structure program details



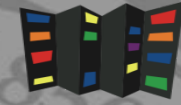
- Installed years ago when had JBOD NFS file servers
 - very slow MSA search
 - only marginally better with spinner local disk
 - copy database indices into RAM disk (~20-30GB), ~8x speedup over NFS file server, ~2x over VAST
 - custom separate MSA and inference into 2 jobs, MSA on CPUs, inference on GPUs
 - alias `run_alphafold.sh` command with paths to the databases via Lmod module
 - <https://www.chpc.utah.edu/documentation/software/alphafold.php#alphafold>



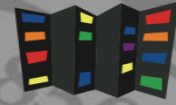
- Installed in a container
- 2 jobs, 2nd job automatically submitted by the first job, 2nd job waits till 1st job finishes
- 1st job uses CPUs and more RAM (DBs on RAM disk)
- 2nd job uses GPU, no need for DBs
- Detailed instructions
 - <https://www.chpc.utah.edu/documentation/software/alphafold.php#alphafold>
- Uses fasta file for input



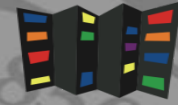
- Databases are not indexed - all on VAST
 - a bit slower than indices in RAM, but freer RAM
- AF3 has options to split the MSA and inference - use that
- Parameters license restriction
 - require users to request access to the parameters, when granted we add them to "alphafold3" group which has access to the centrally installed parameters
- <https://www.chpc.utah.edu/documentation/software/alphafold.php#alphafold3>



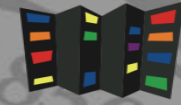
- Installed in a container
- 2 jobs, like in AlphaFold 2
- Detailed instructions
 - <https://www.chpc.utah.edu/documentation/software/alphafold.php#alphafold3>
- Not just proteins, uses json file for input
- Need to agree to [Model Parameters](#) and [Outputs](#) terms of use by filling out [weights access form](#) and then e-mailing helpdesk@chpc.utah.edu with the approval



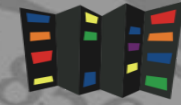
- Discovered when trying to find faster MSA alternative
 - different, faster MSA program
 - MSA search done at remote ColabFold server by default
 - we run our own ColabFold (mmseqs2) servers
 - use LocalColabFold for running jobs
 - <https://www.chpc.utah.edu/documentation/software/alphafold.php#colabfold>



- One researcher got banned from ColabFold server and asked us to set one up locally
- For good performance need some databases in RAM (~0.7 TB) and the rest on SSD drive (~1 TB)
- We found 10yo donated 32 core server with 1 TB RAM and 1 TB SSD
 - copy the non-RAM databases to local SSD, RAM databases on network file server symlinked to local SSD
 - "jobs" directory that caches past jobs on network file server
- New server with 64 cores, 2 TB RAM



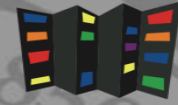
- Installed in Miniforge environment
- One job on GPUs
- Detailed instructions
 - <https://www.chpc.utah.edu/documentation/software/alphafold.php#colabfold>
- Uses fasta file for input
- Use CHPC's MSA server running [mmseqs2](#)
 - `colabfold_batch --host-url=http://colabfold02.int.chpc.utah.edu:8088`
- Generally faster than Alphafold (because of the faster MSA)
- Uses Alphafold 2 for inference



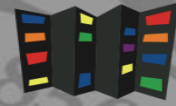
- Open source alternative to Alphafold 3
- Easy installation into venv with pip
- Uses mmseqs2 for MSA search, like ColabFold
- Can use the ColabFold server



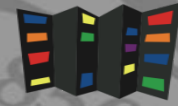
- Fully open source
- Installed as Miniforge3 environment
- Single GPU job
- Can work with multiple biomolecule types
- Can use mmseqs2 for MSA
- Detailed instructions
 - <https://www.chpc.utah.edu/documentation/software/alphafold.php#boltz>
- Uses fasta or yaml file for input
- Use lmi4boltz for very large sequences



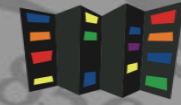
- Alphafold, RFdiffusion, ProteinMPNN in Colab notebook
- User requested to run Colab on their hardware
- Set up Colab container for local runtime
- Run SLURM job with Jupyter using the Colab container
- SSH tunnel from personal desktop to compute node
- Connect to Jupyter from personal web browser
- Modify the Colab notebook for shared user env.
- <https://www.chpc.utah.edu/documentation/software/google-colab.php>



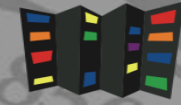
- Some tools, esp. from the Baker lab are available as Colab notebooks
 - <https://colab.research.google.com/github/sokrypton/ColabFold/blob/main/AlphaFold2.ipynb>
 - <https://colab.research.google.com/github/sokrypton/ColabDesign/blob/main/rf/examples/diffusion.ipynb>
- Free Colab has 15 min execution limit
 - can use "local runtime" to run on CHPC
 - start Jupyter with Colab environment as a SLURM job
 - create SSH tunnel to this job from local laptop/desktop
 - after opening Colab notebook choose the "local runtime"
- Detailed instructions
 - <https://www.chpc.utah.edu/documentation/software/google-colab.php>



- Most notebooks require local installation of the dependencies
 - install inside of the Colab notebook or manually in cluster terminal
 - may need to change hard coded paths in the notebook (e.g. /usr/local) to paths in user's CHPC home directory
 - not for the faint hearted but for the most part doable
- Caveats
 - the Colab notebook resource monitor shows the whole node, not its part allocated to the job
 - GPU and disk resource monitor are not correct
 - the Colab notebook's file manager does not work - use one that comes with Jupyter



- More like a workflow
- User installable, and requires newer OS
 - typical example of packages coming from RosettaCommons/Baker Lab
 - workaround to set it all up in a container
- In the next slide we show installation instructions modified for running in an Apptainer container
- Better use new Apptainer instructions (we may install it as a module)
 - <https://github.com/RosettaCommons/RFantibody?tab=readme-ov-file#apptainer-installation-hpc>



- Shell into a container that has set RF Antibody base dependencies

```
cd <directory where you want to install the program>
```

```
ml apptainer
```

```
apptainer shell --nv -B ./:/home
```

```
/uufs/chpc.utah.edu/common/home/u0101881/containers/singularity/containers/rfantibody/rfa.sif
```

- In the container, get a RF Antibody fork with fixed examples and download the weights

```
git clone https://github.com/katyachemistry/RFantibody.git
```

```
cd RFantibody
```

```
git checkout fixed
```

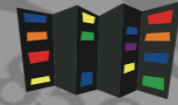
```
bash ./include/download_weights.sh
```

- For the dependencies setup, need to run "poetry install" outside of the bash script, otherwise it creates new venv.

```
nano include/setup.sh
```

comment out

```
#poetry install &&
```



- then run the dependency installation:

```
bash include/setup.sh  
poetry install
```

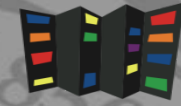
- The program is installed, now you can run the example

```
bash /home/scripts/examples/rfdiffusion/antibody_pdbdesign.sh
```

This needs to be run on a GPU which a decent size memory, it crashed on my desktop which only has 2 GB GPU.

- If you need to run other simulation, you don't need to do steps 1-3, just shell into the container in the same directory you have the program installed

```
apptainer shell --nv -B ./:/home  
/uufs/chpc.utah.edu/common/home/u0101881/containers/singularity/containers/rfantibod  
y/rfa.sif
```



- Recent install / in the works
- RFdiffusion is designed for user space - have to modify for shared environment
- RosettaFold All Atom - similar functionality to Alphafold 3 or Boltz
- Overall issue with tools from RosettaCommons is that they are designed to be installed by user, which requires modifications if it's installed by us for everyone