

IMPECCABLE: Integrated Modeling PipelinE for COVID Cure by Assessing Better LEads

INTERNATIONAL CONFERENCE ON PARALLEL PROCESSING

ICPP / 2021 / CHICAGO / USA

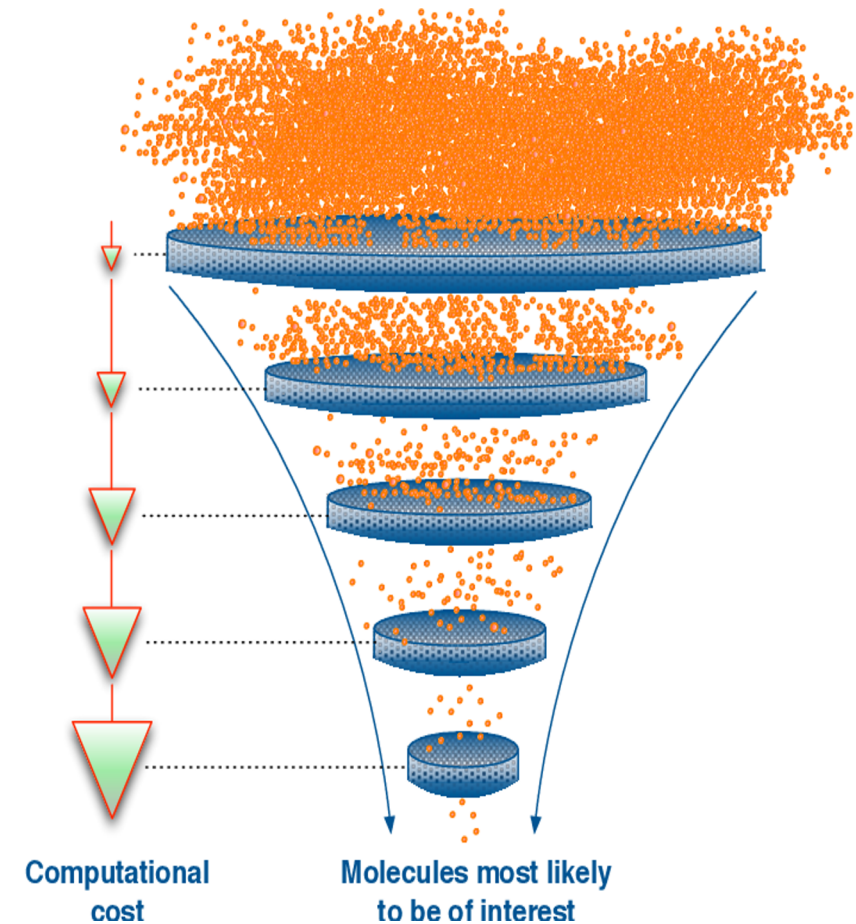


AUGUST 9-12, 2021

Aymen Al Saadi Rutgers University, New Brunswick	Dario Alfe University College London and University of Naples Federico II	Yadu Babuji University of Chicago	Agastya Bhati University College London
Ben Blaiszik University of Chicago and Argonne National Laboratory	Alexander Brace Argonne National Laboratory	Thomas Brettin Argonne National Laboratory	Kyle Chard University of Chicago and Argonne National Laboratory
Ryan Chard University of Chicago	Austin Clyde University of Chicago and Argonne National Laboratory	Peter Coveney* University College London and University of Amsterdam	Ian Foster University of Chicago and Argonne National Laboratory
Tom Gibbs NVIDIA Corporation	Shantenu Jha† Brookhaven National Laboratory, and Rutgers University, New Brunswick	Kristopher Keipert NVIDIA Corporation	Thorsten Kurth NVIDIA Corporation
Dieter Kranzlmüller Leibniz Supercomputing Centre	Hyungro Lee Rutgers University, New Brunswick	Zhuozhao Li University of Chicago	Heng Ma Argonne National Laboratory
Andre Merzky Rutgers University, New Brunswick	Gerald Mathias Leibniz Supercomputing Centre	Alexander Partin Argonne National Laboratory	Junqi Yin Oak Ridge Leadership Computing Facility
Arvind Ramanathan‡ University of Chicago and Argonne National Laboratory	Ashka Shah Argonne National Laboratory	Abraham Stern NVIDIA Corporation	Rick Stevens§ University of Chicago and Argonne National Laboratory
Li Tan Brookhaven National Laboratory	Mikhail Titov Rutgers University, New Brunswick	Anda Trifan Argonne National Laboratory	Aristeidis Tsaris Oak Ridge Leadership Computing Facility
Matteo Turilli Rutgers University, New Brunswick	Huub Van Dam Brookhaven National Laboratory	Shunzhou Wan University College London	David Wifling Leibniz Supercomputing Centre

Overview

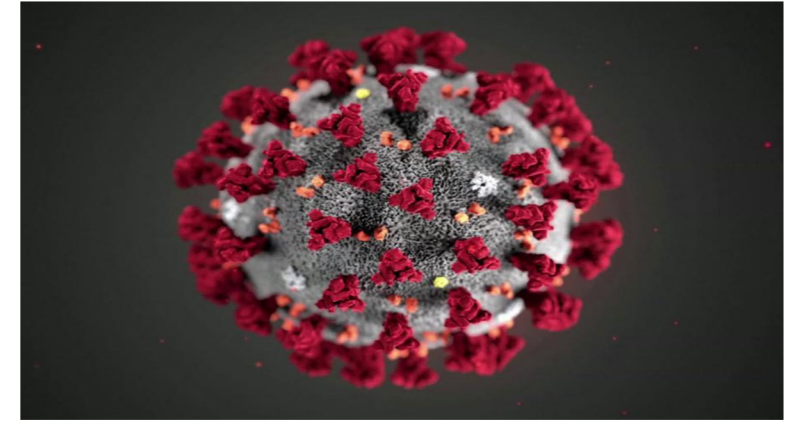
- Drug Discovery & Design is a complex & expensive
 - $O(10)$ years; $O(10^9)$ \$
 - $O(10^{68})$ exhaustive search not an option!
- Integrated performance of multiple stages (methods)
 - Different stages with varying cost vs accuracy
 - Other challenges in constructing: filter ratio
- AI-driven HPC 100-100x *effective performance* of traditional HPC simulations
 - Heterogeneous and adaptive workflows
 - Systems software evolve in response



Ref. Aspuru-Guzik

National Virtual Biotechnology Lab (NVBL)

- National Virtual Biotechnology Lab (NVBL)
 - <https://science.osti.gov/nvbl>
- Aid U.S. policymakers in responding to the COVID-19 pandemic with epidemiological information for decision making
- Accelerate production of critical medical supplies across the nation
- **Supercomputing and artificial intelligence for design of targeted therapeutics**
- Leverage chemical testing & analysis to facilitate new antigen and antibody testing



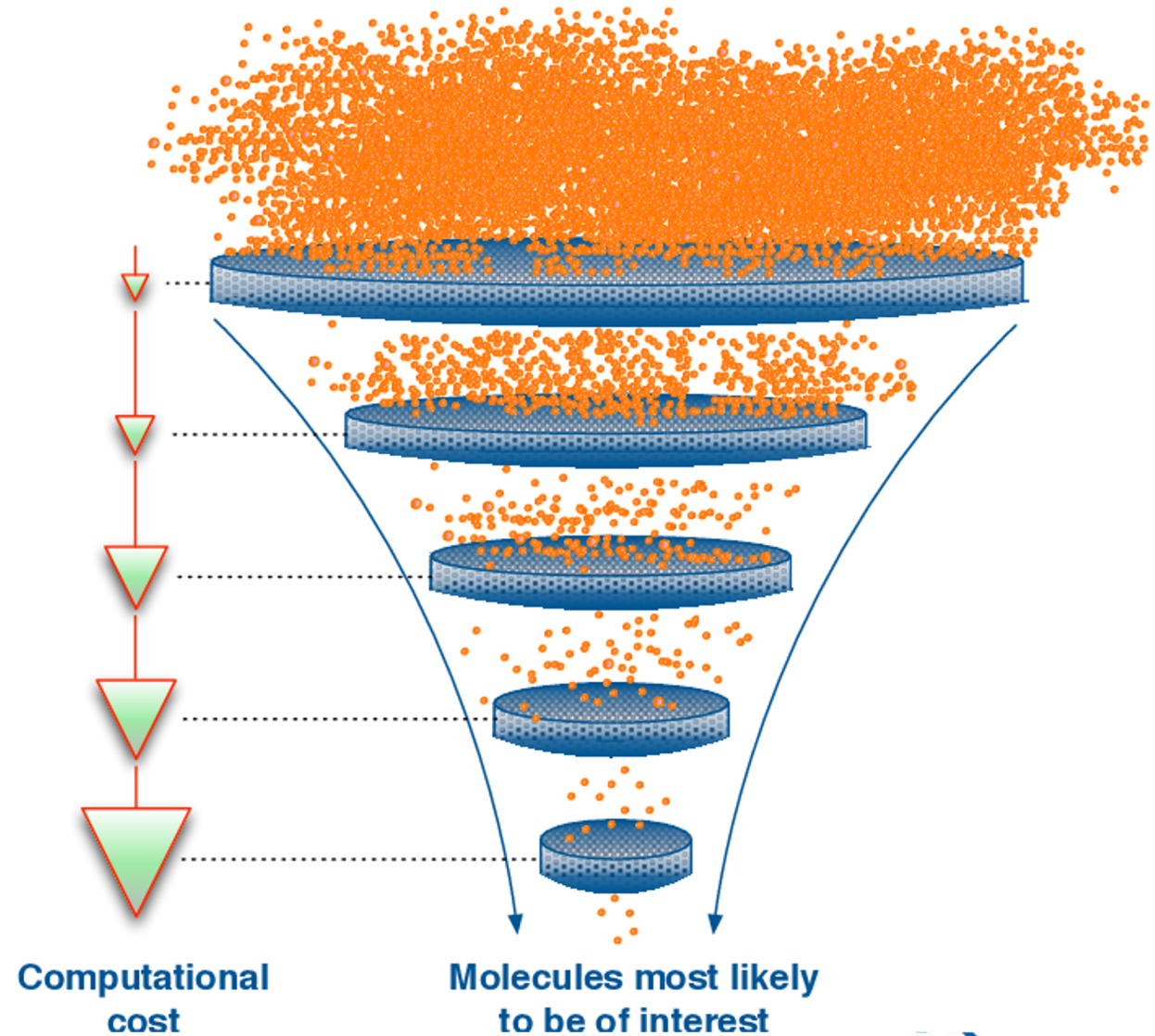
NVBL given US Secretary of Energy Honour Award (2021)

IMPECCABLE: uHigh-Throughput Virtual Screening

Multi-stage campaign employed to select promising drug candidates:

- **WF1:** High-throughput ensemble docking to identify small molecules (“hits”)
- **WF2:** ML-driven Molecular Dynamics for modeling specific binding regions and understanding mechanistic changes involving drugs
- **WF3&4:** Binding Free Energy calculations of promising leads (“Hit-to-Lead” & “Lead Optimization”)

<https://arxiv.org/abs/2010.06574>

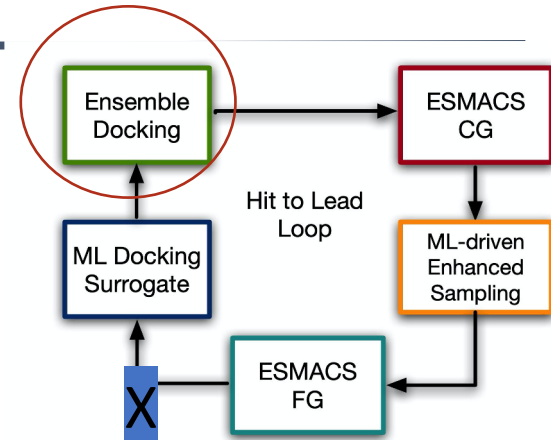
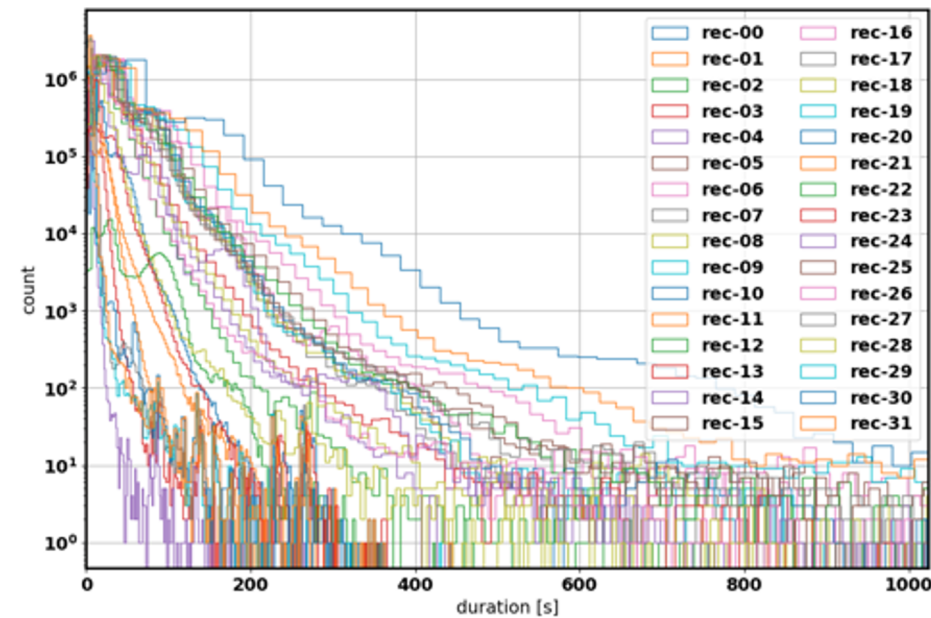
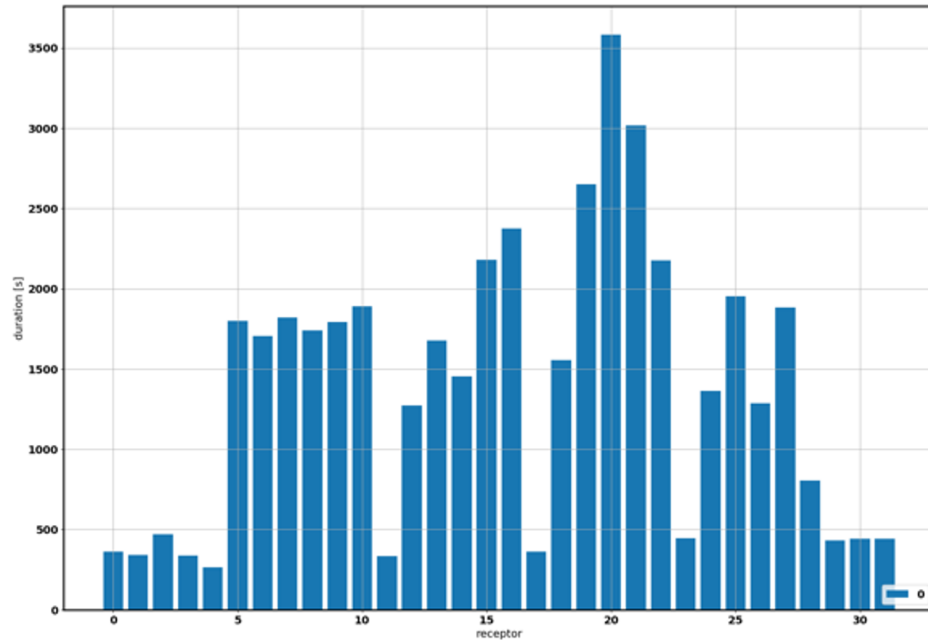


Computational Challenges: Heterogeneity

- **Heterogeneity** of different types and at multiple levels
 - Coupled AI-HPC (WF2)
 - High-throughput function calls (WF1)
 - Ensembles of MPI tasks (WF3/4)
- Spatio-temporal variation within and across WF1

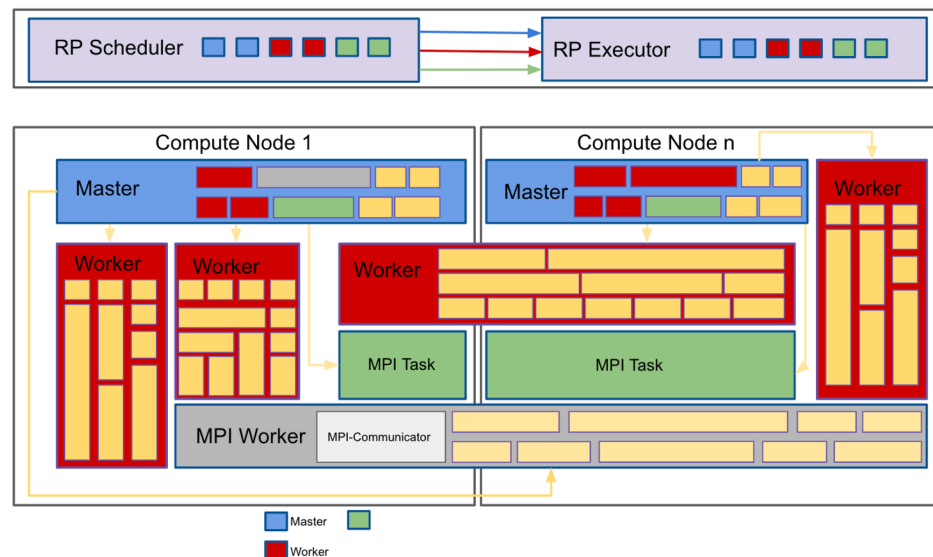
HPC Platform	Facility	Batch System	Node Architecture		GPU	Workflows	Max # nodes utilized
			CPU				
Summit	OLCF	LSF	2 × POWER9 (22 cores)		6 × Tesla V100	WF1-4	2000
Lassen	LLNL	LSF	2 × POWER9 (22 cores)		4 × Tesla V100	WF2,3	128
Frontera	TACC	Slurm	2 × x86_64	(28 cores)	—	WF1	7650
Theta	ALCF	Cobalt	1 × x86_64	(64 cores)	—	WF1	256
SuperMUC-NG	LRZ	Slurm	2 × x86_64	(24 cores)	—	WF3-4	6000 (with failures)

Docking: WF1



- Docking: OpenEye; Library (ORD): 6.25M ligands (drug candidate); 32 targets/receptors
 - Fluctuations in docking execution time library (ORD) for different receptors
 - Long-tailed Tx for different ligands for a given target (receptor)
 - Many work items (function calls) need to be distributed
 - Call duration varies two order of magnitudes (1-100s). Mean duration 8s.

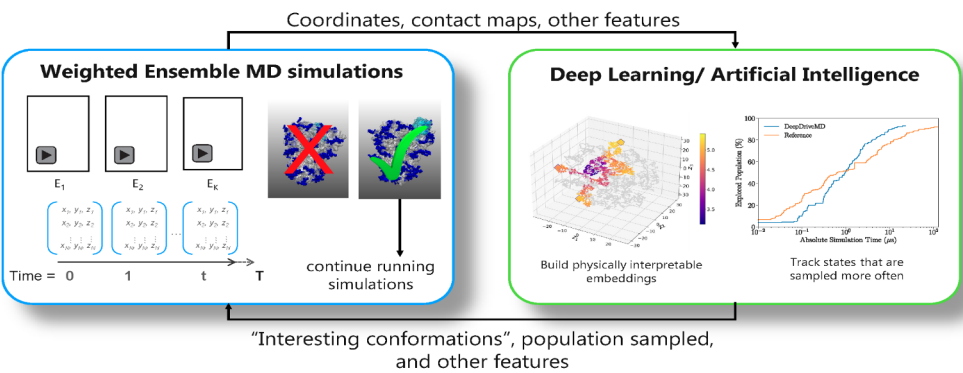
Ensemble Docking (WF1) with RAPTOR



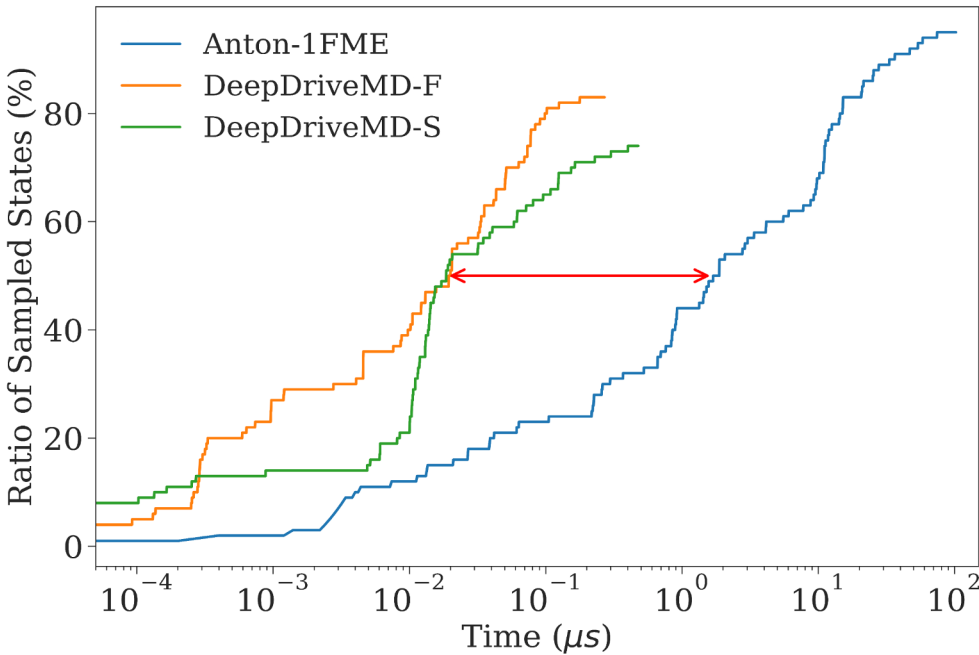
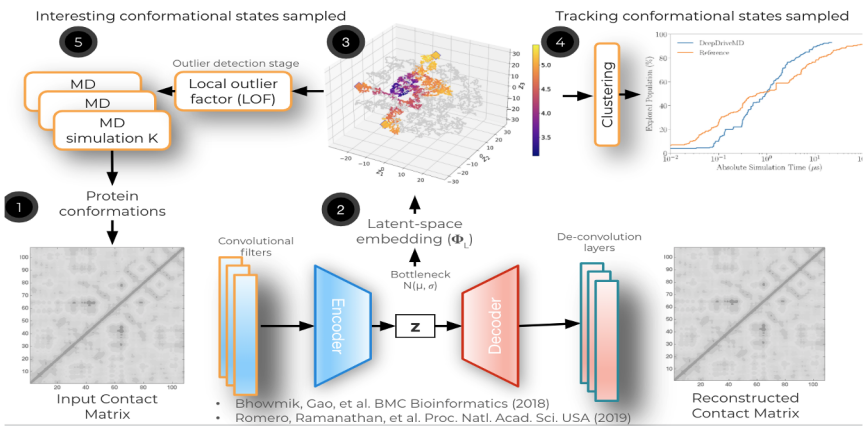
ID	Platform	Application	Nodes	Pilots	Tasks [$\times 10^6$]	Startup [sec]	Utilization avg / steady	Task Time [sec]		Rate [$\times 10^6/h$]	
								max	mean	max	mean
1	Frontera	OpenEye	128	31	205	129	90% / 93%	3582.6	28.8	17.4	5.0
2	Frontera	OpenEye	7600	1	126	81	90% / 98%	14958.8	10.1	144.0	126.0
3	Frontera	OpenEye	8336	1	13	451	63% / 98%	219.0	25.3	91.8	11.0
4	Summit	AutoDock	1000	1	57	107	95% / 95%	263.9	36.2	11.3	11.1

ML-driven Simulations (WF2): 10-100x Sampling

Combining AI with HPC: AI-driven MD simulations -- DeepDriveMD



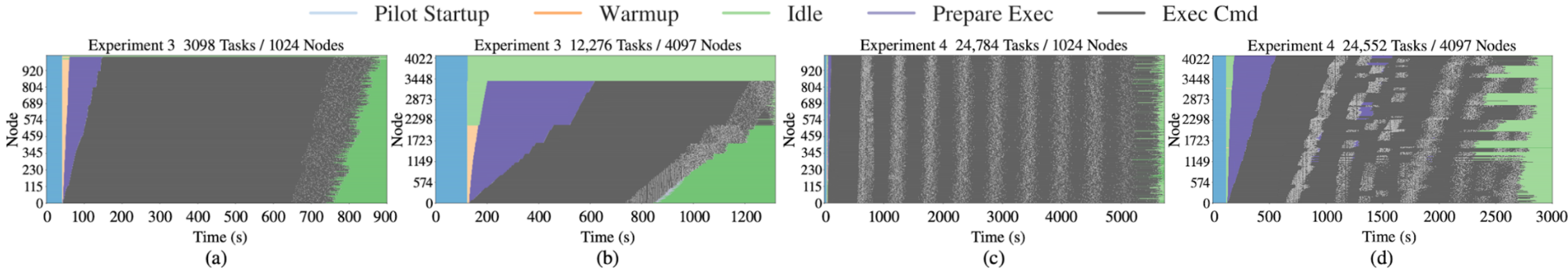
- Learning Everywhere**
- Jha & Fox. In Visionary Track", 15th International Conference eScience (2019), San Diego, California
 - Jha & Fox. 15th International Conference eScience (2019), San Diego, California



Exp.	Ligands	GPUs	Tasks	Platform	Overhead
PLC-1	1	120	250	Summit	334.21 s
PLC-2	1	120	250	Lassen	302.34 s
PLC-3	8	960	2000	Summit	265.05 s
PLC-4	8	960	960	Lassen	314.83 s
PLC-5	51	6120	6120	Summit	254.01 s

Binding Free Energy (WF3 & 4): Heterogeneous Simulations

ID	HPC Platform	#Tasks	#Generations	Task Runtime	#Cores/Task	#GPUs/Task	#Cores/Pilot	#GPUs/Pilot
1	Titan	$2^n; n = [5 - 12]$	1	$828s \pm 14s$	32	-	$2^n; n = [10 - 17]$	-
2	Titan	2^{14}	$2^n; n = [5 - 3]$	-	-	-	$2^n; n = [14 - 16]$	-
3	Summit	3098; 12,276	1	$600s - 900s$	1 - 42	0; 6	43,008; 172,074	6144; 24,582
4	Summit	24,552; 24,784	$\approx 2; 8$	$500s - 600s$	1 - 42	0; 6	-	-
5	Frontera	126×10^6	≈ 300	$1s - 120s$	1	-	392,000	-

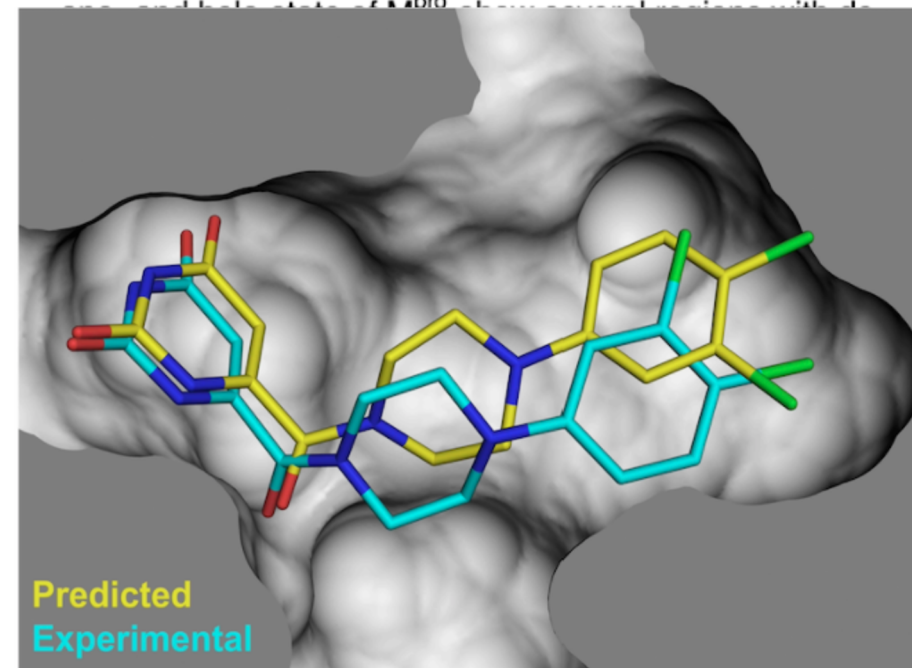


Impacting SARS-CoV-2 Medical Therapeutics

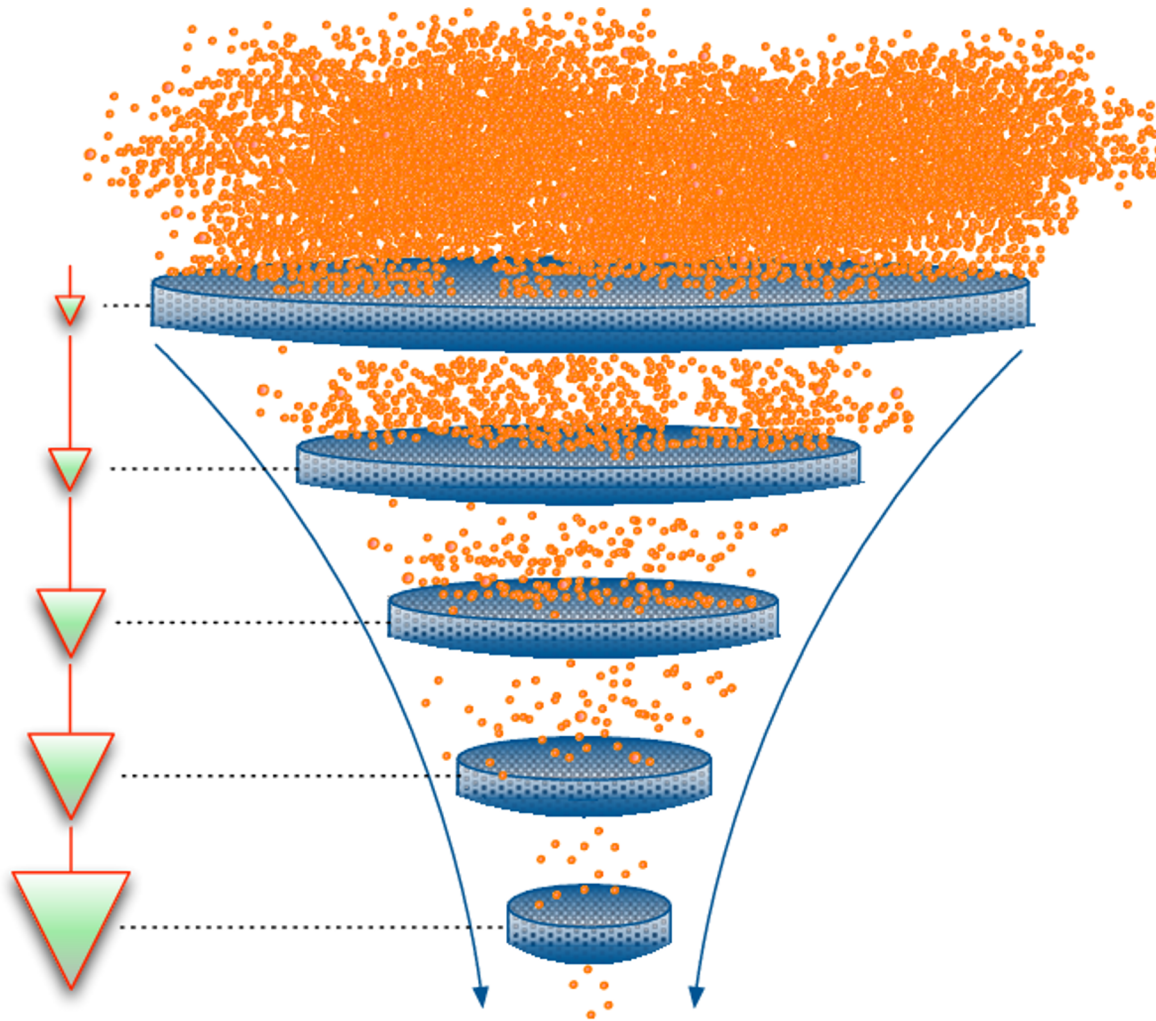


- **Scale of Operation:**
 - $\sim 10^{11}$ Docking calculations
 - $\sim 10^3$ ML-driven MD calculations
 - $\sim 5 \times 10^4$ Binding Free Energy Calculations
 - $\sim 2.5 \times 10^6$ node-hours
- Peak Performance
 - ~ 8000 nodes (Frontera, April. 2021)
 - ~ 4000 nodes on Summit
- Extensible Computational Infrastructure and Capabilities
 - Beyond COVID-19 ?

Fig. 4. Conformational changes upon MCULE-5948770040 binding to M^{pro} indicate changes within distinct regions, both close-to and farther-away from the primary binding site. (a) RMS fluctuations of the



IMPECCABLE: High-Throughput Virtual Screening



Computational
cost

Molecules most likely
to be of interest

ML1

Docking Surrogate

S1

AutoDock-GPU

S3-CG

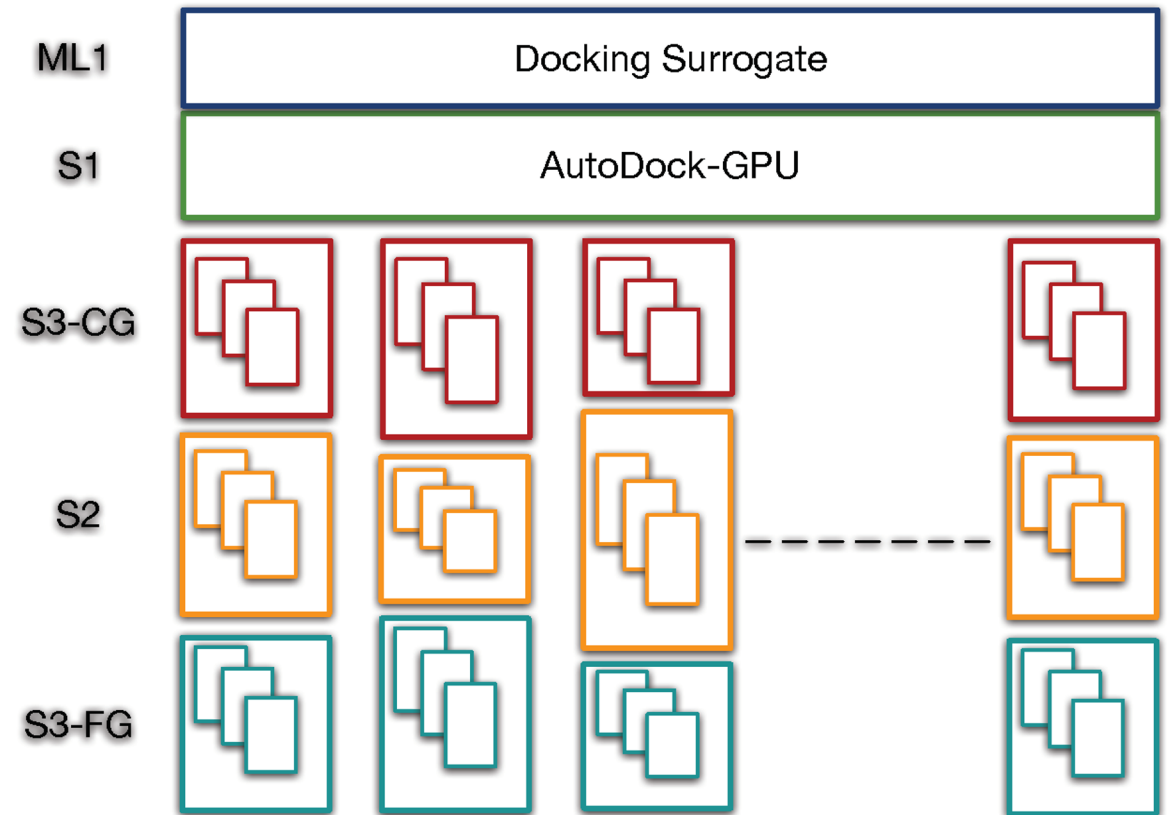
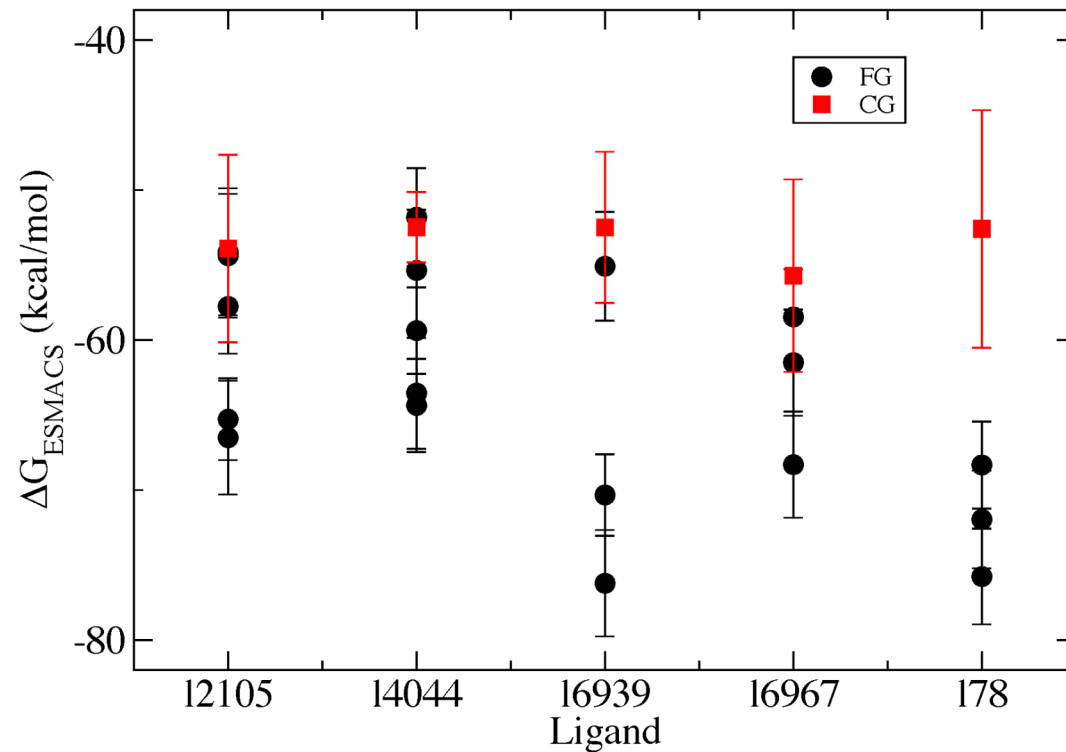
S2

S3-FG

IMPECCABLE: Integrated Modeling Pipeline



Results from S3-CG and S3-FG



IMPECCABLE: Integrated Modeling Pipeline

Why is this challenging?

- **Heterogeneous** at multiple levels
 - Coupled AI-HPC (WF2)
 - High-throughput function calls (WF1)
 - Ensembles of MPI tasks (WF3/4)
 - Spatio-temporal variation within each
- **Collective** versus single-task performance
 - Campaigns are "integrated" workflows: WF1 and WF4 differ by 10^7 x in computational cost
 - Producers of data (WF1) and consumers (ML1)
- **Adaptive Execution** at multiple levels
 - Workload: Task mix varies over campaign
 - Tasks: Run for varying duration

1000x variation in throughput

Table 3: Throughput and performance measured as peak flop per second (mixed precision, measured over short but time interval) per Summit node (6 NVIDIA V100 GPU).

Comp.	#GPUs	Tflop/s	Throughput
ML1	1536	753.9	319674 ligands/s
S1	6000	112.5	14252 ligands/s
S3-CG	6000	277.9	2000 ligand/s
S3-FG	6000	732.4	200 ligand/s

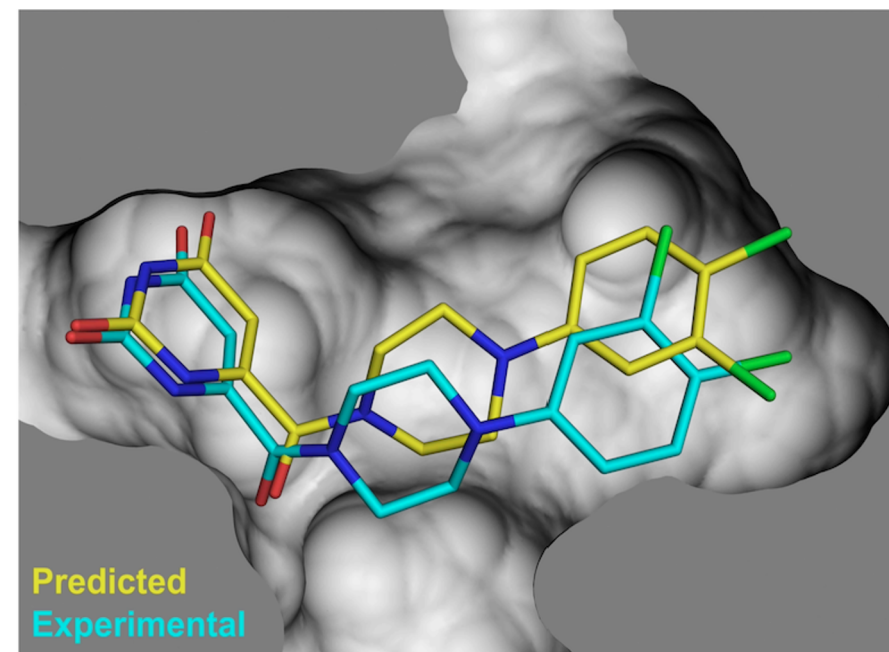
10^7 x variation in cost across workflows

Table 2: Normalized computational costs on Summit.

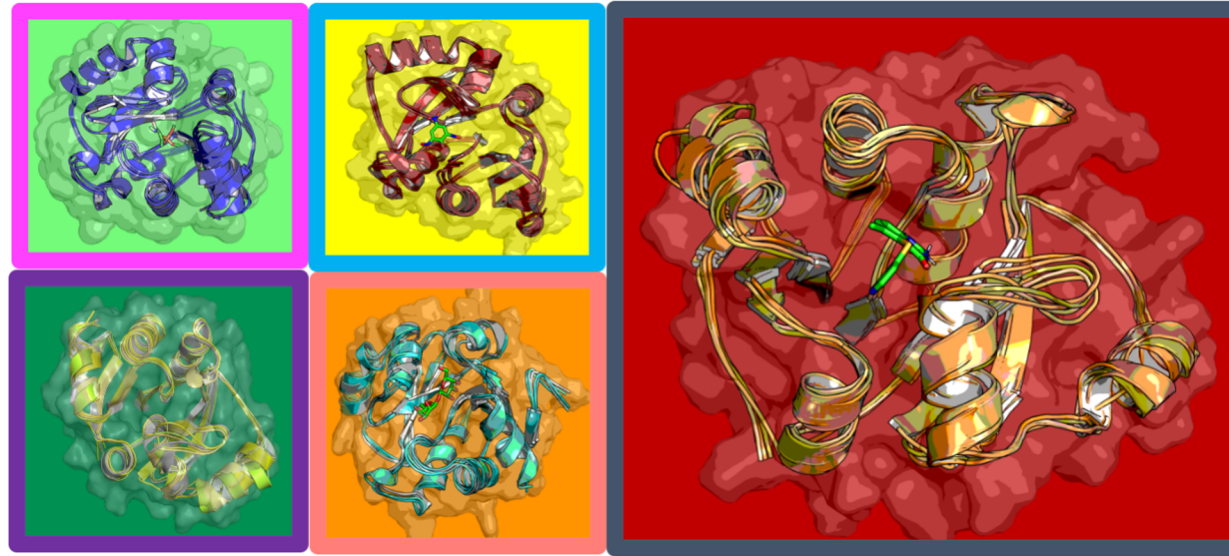
Method	Nodes per ligand	Hours per ligand (approx)	Node-hours per ligand
Docking (S1)	1/6	0.0001	~0.0001
BFE-CG (S3-CG)	1	0.5	0.5
Ad. Sampling (S2)	2	2	4
BFE-FG (S3-FG)	4	1.25	5
BFE-TI (not integrated)	64	10	640

Summary

- Drug Discovery & Design is a complex & expensive
 - Infrastructure, Methodological, Scientific
- Developed 1st gen of AI-HPC infrastructure
 - Sophistication of AI-HPC methods will grow
- Rethink systems software ecosystem
 - Collective perf. of heterogeneous workflows; not just single tasks
 - Advances in adaptive runtime systems for such campaigns



Thank you!



Funding acknowledgement:

- DOE National Virtual Biotechnology Laboratory
- DOE CANDLE ECP
- ECP ExaWorks and ECP ExaLearn
- ASCR Surrogates Benchmarking Initiative
- NSF RADICAL-Cybertools`