

Creating a platform for rapid computational antibody design via machine learning, HPC, and laboratory experimentation

Thomas Desautels
LLNL

LLNL ML4I Workshop
August 12, 2021

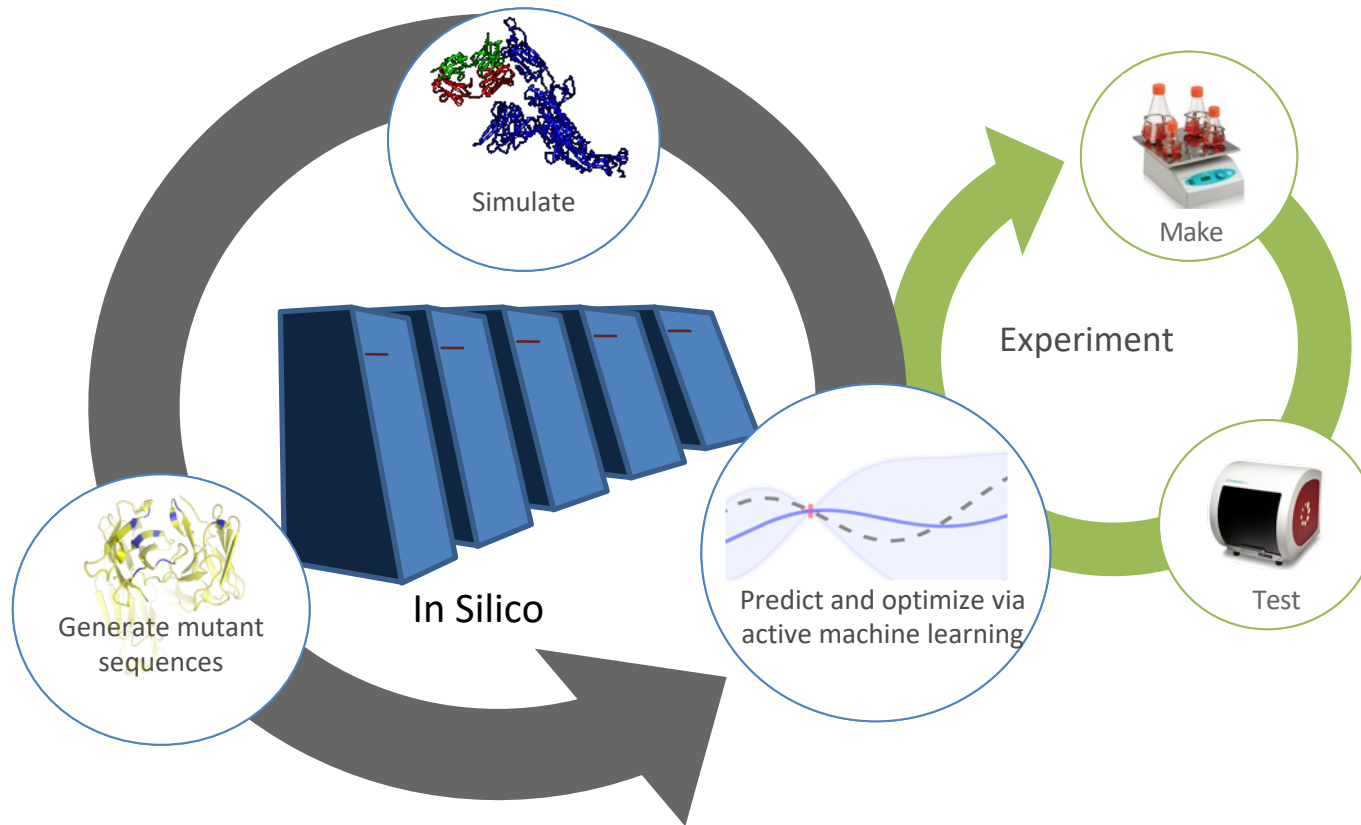
LLNL: Daniel Faissol, Adam Zemla, Ed Lau, Fangqiang Zhu, John Goforth, Denis Vashchenko, Mary Silva, Rebecca Haluska, Drew Bennett, Emilia Grzesiak, Alexander Ladd, Brent Segelke, Feliza Bourguet, Victoria Lao, Monica Borucki, Dina Weilhammer, Jacky Lo, Nicole Collette, Magdalena Franco, Kathryn Arrildt
Sandia NL: Brooke Harmon, Oscar Negrete, Max Stefan



The problem: how can you *rapidly* respond to a new pathogen?

- Premises:
 - New pathogens can emerge with little warning
 - The immune system may need assistance to effectively counter a new pathogen
 - Vaccine antigens and therapeutic antibodies are the most important protein design targets
 - Basically the only things that have worked at all for COVID
- Ordinarily, vaccines and therapeutic antibodies take years or decades to reach market
- In the long-term, we want a system for scalable, high-confidence, *in silico* design that could accelerate delivery of a countermeasure that is (1) effective, (2) manufacturable, and (3) safe.
- In a familiar LLNL plan, do design & certification as much as possible in the computer
 - Critically, this can enable *preemptive* design against emerging virus variants or novel members of families of pathogens

Our approach to countermeasure design combines simulation and ML-driven decision-making with laboratory experimentation



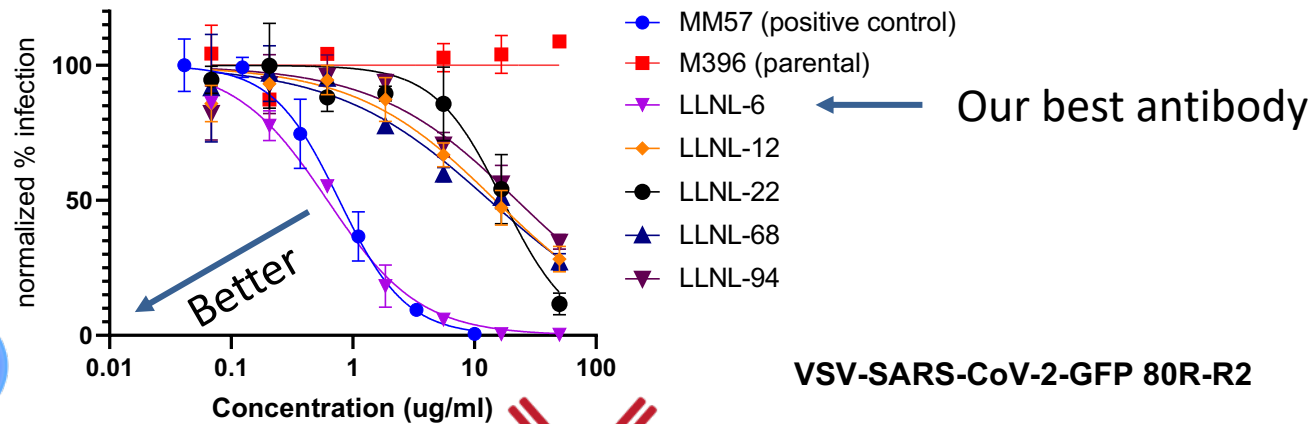
- Select & simulate computationally until promising candidates are found
- Send best candidates for laboratory testing
- If necessary, re-design from most promising candidates identified in the laboratory

We've executed and validated rapid antibody design against SARS-CoV-2: novel to our knowledge

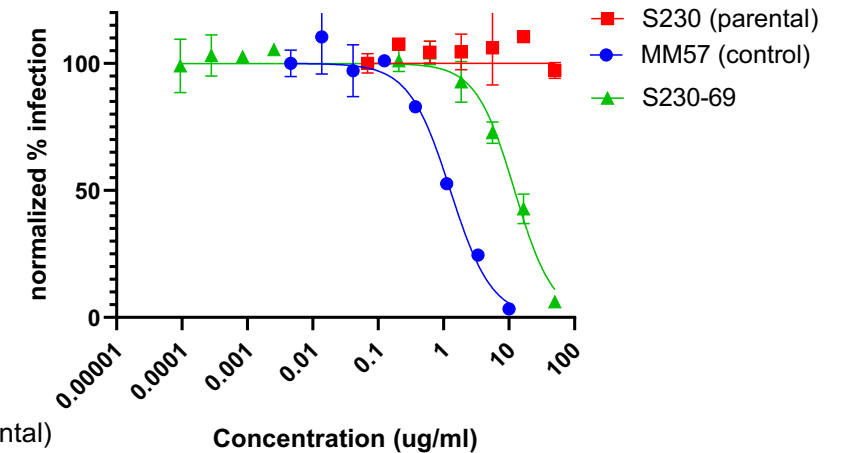
- From Jan 2020 to present, designed several neutralizing antibodies for SARS-CoV-2



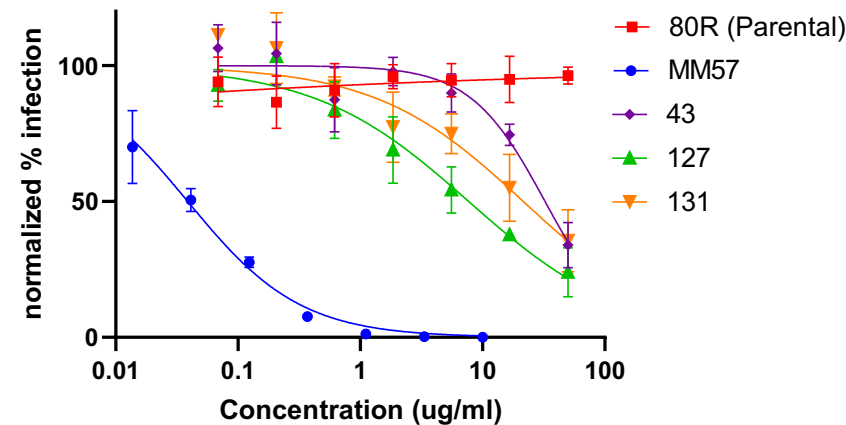
VSV-SARS-CoV-2-GFP M396-R2 hits



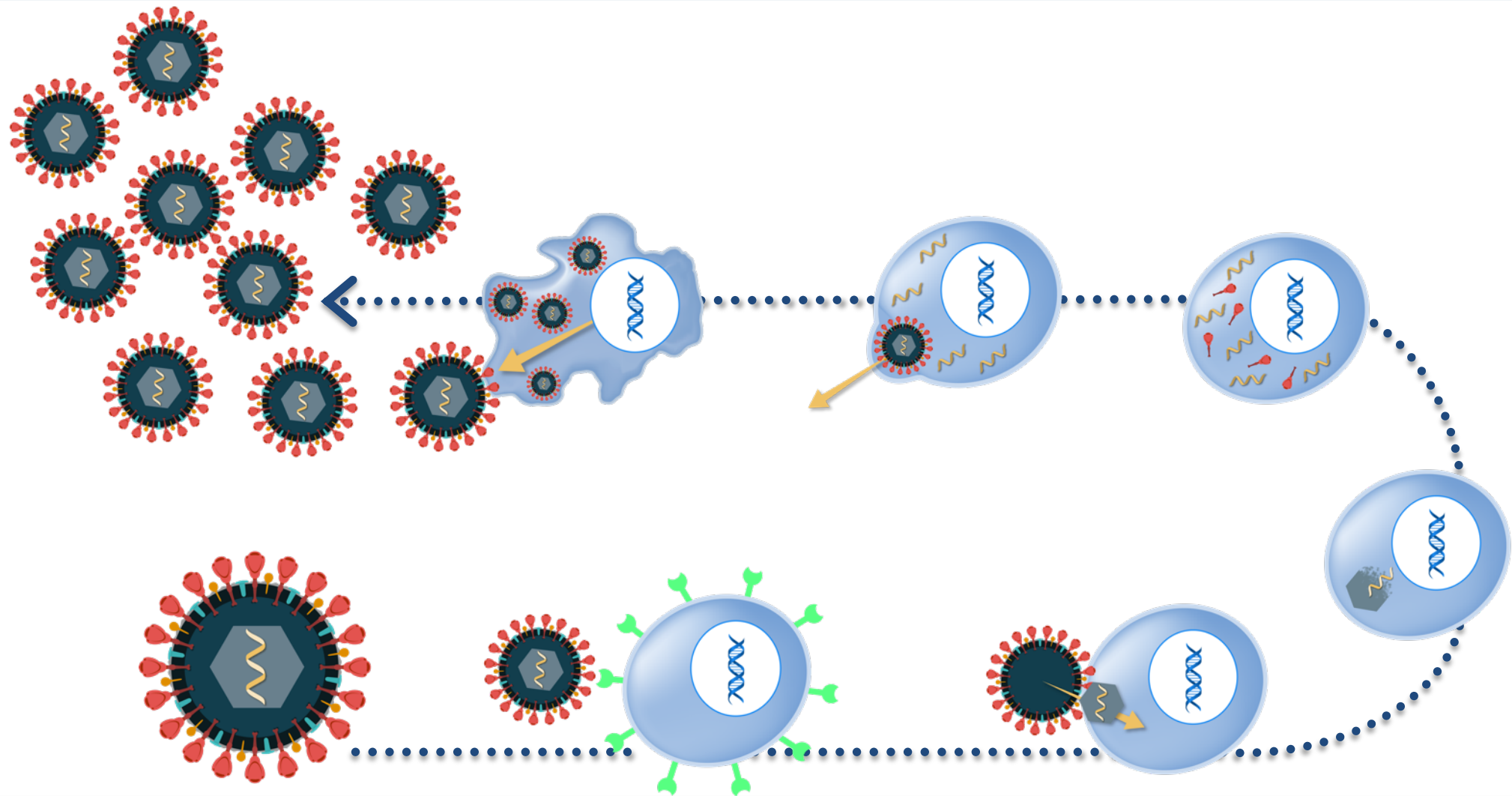
VSV-SARS-CoV-2-GFP S230-69



VSV-SARS-CoV-2-GFP 80R-R2



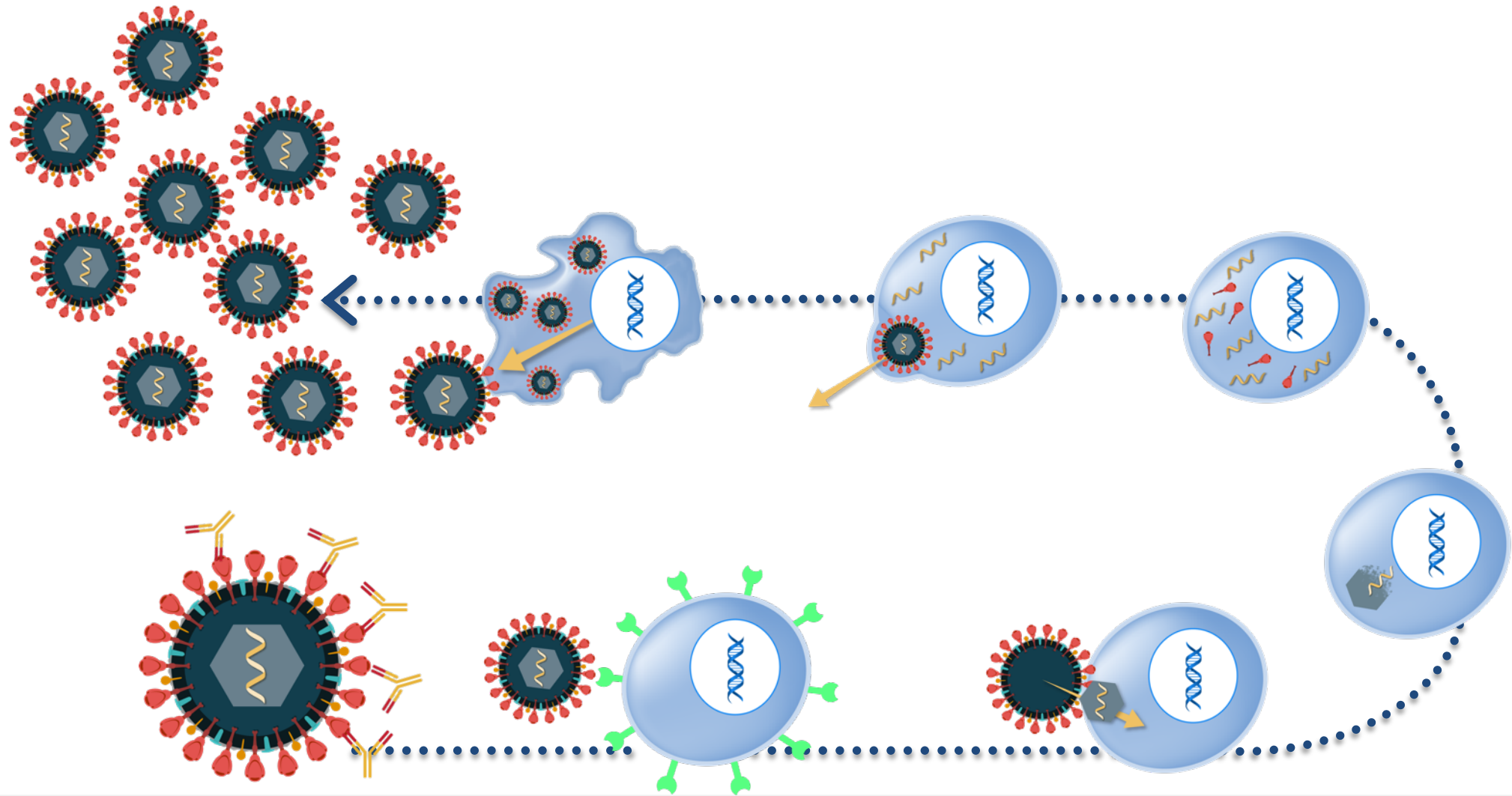
Viruses reproduce by entering and hijacking host cells



Lawrence Livermore National Laboratory
LLNL-PRES-825249



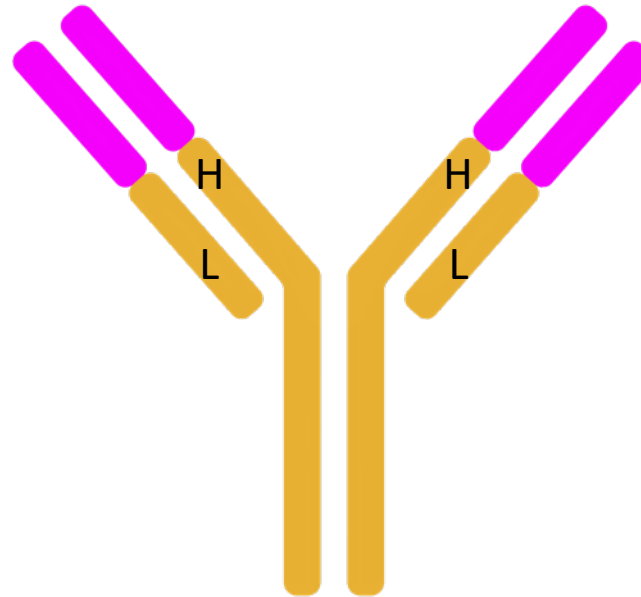
Neutralizing antibodies can stop viral entry



Proteins are described by their amino acid sequence: Antibody design becomes finding a suitable sequence

> m396 heavy chain

QVQLQQSGAEVKKPGSSVKV SCKASG **GTFS**
SYTISWVRQAPGQGLE **WMGGITPIL**GIANY
AQKFQGRVTITTTDESTSTAYMELSSLRSEDTA
VYYCARD**TMGGM**MDVWGQGTTVTVSSAS
TKGPSVFPLAPSSKSTSGGTSALGCLVKDYFP
EPVTVSWNSGALTSGVHTFPAVLQSSGLYSLS
SVVTVPSSSLGTQTYICNVNHNKPSNTKVDKK
VEPKSCDKTSPLFVHHHHHHG DYKD
DDDKG



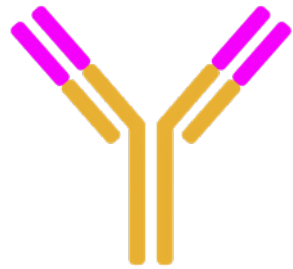
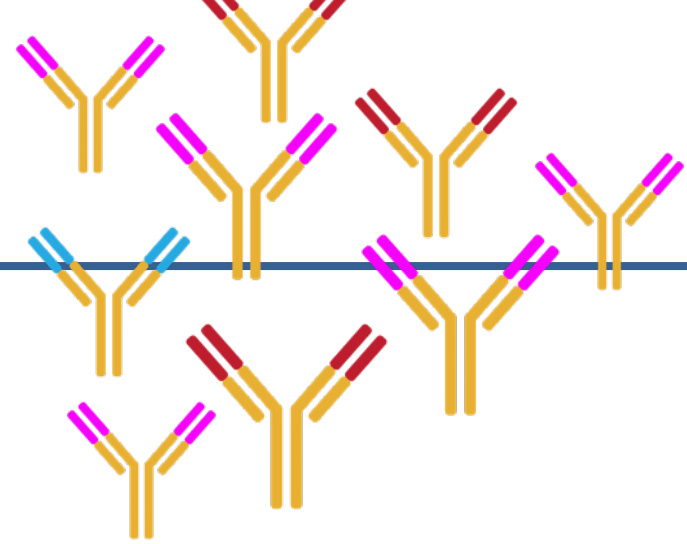
> m396 light chain

SYELTQPPSVSVAPGKTARITCGGN **NIGSKSV**
HWYQQKPGQAPV**LVVYDDSDRPS**GIPERFS
GSNSGNTATLTISRVEAGDEADYYC **QVWDSS**
SDYVFGTGTKVTVLGQPKANPTVTLFPPSSE
EFQANKATLVCLISDFYPGAVTVAWKADGSP
VKAGVETTKPSKQSNNKYAASSYLSLTPEQW
KSHRSYSCQVTHEGSTVEKTVAPTECS

m396 neutralizes **SARS-CoV-1**, but not **SARS-CoV-2**; can its sequence be modified to bind a target antigen and neutralize a new virus?



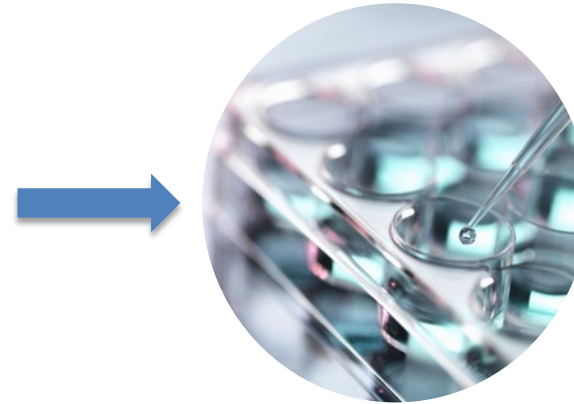
The design space is vastly larger than what we can simulate or test



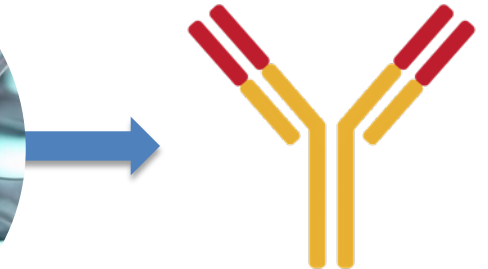
CoV-1 +
changes
 $\sim 10^{30}$



Computer
Simulations
1,000,000

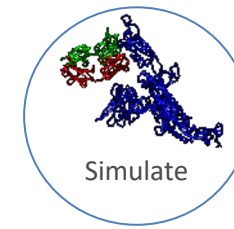


Laboratory
Experiments
100-1,000

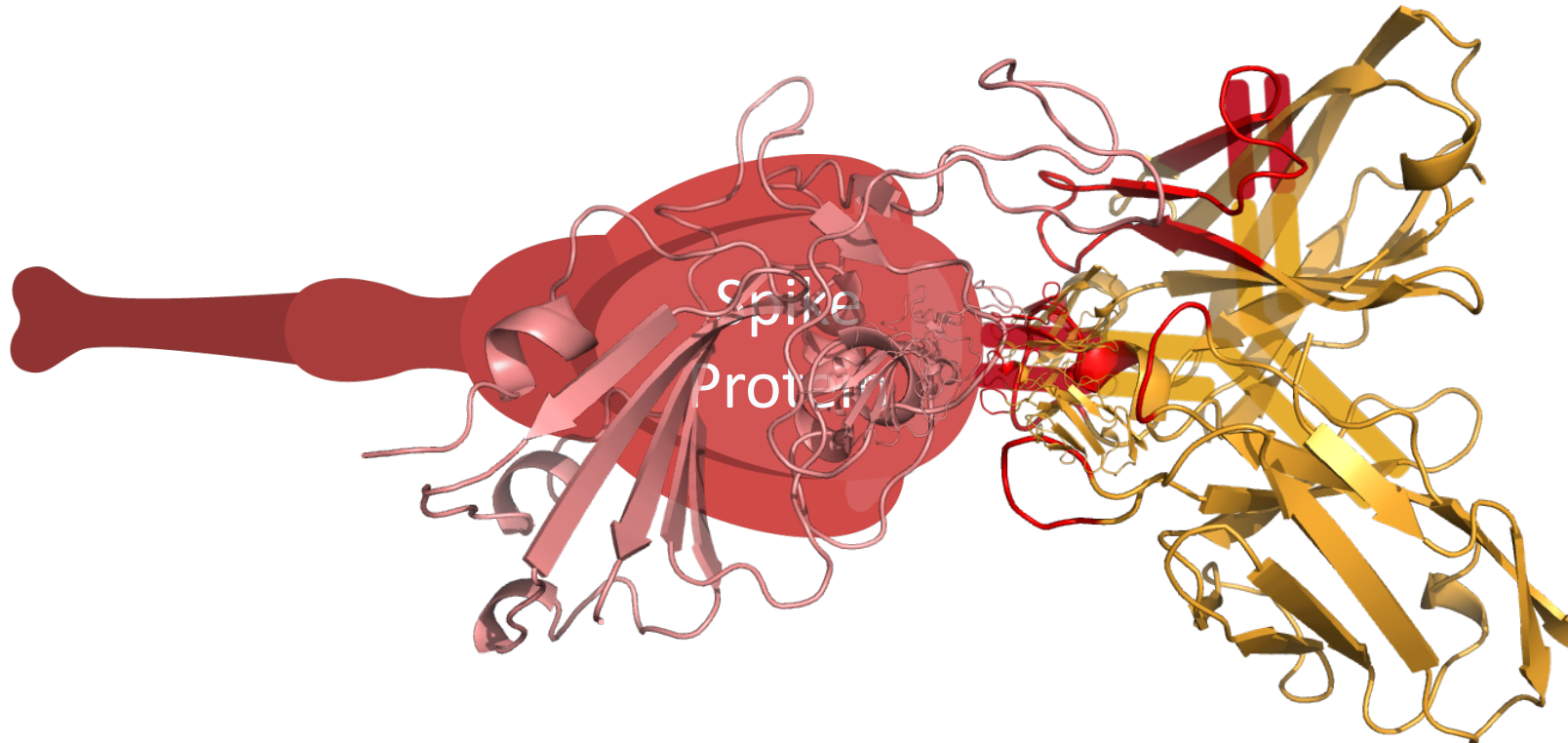


CoV-2
Need just one!

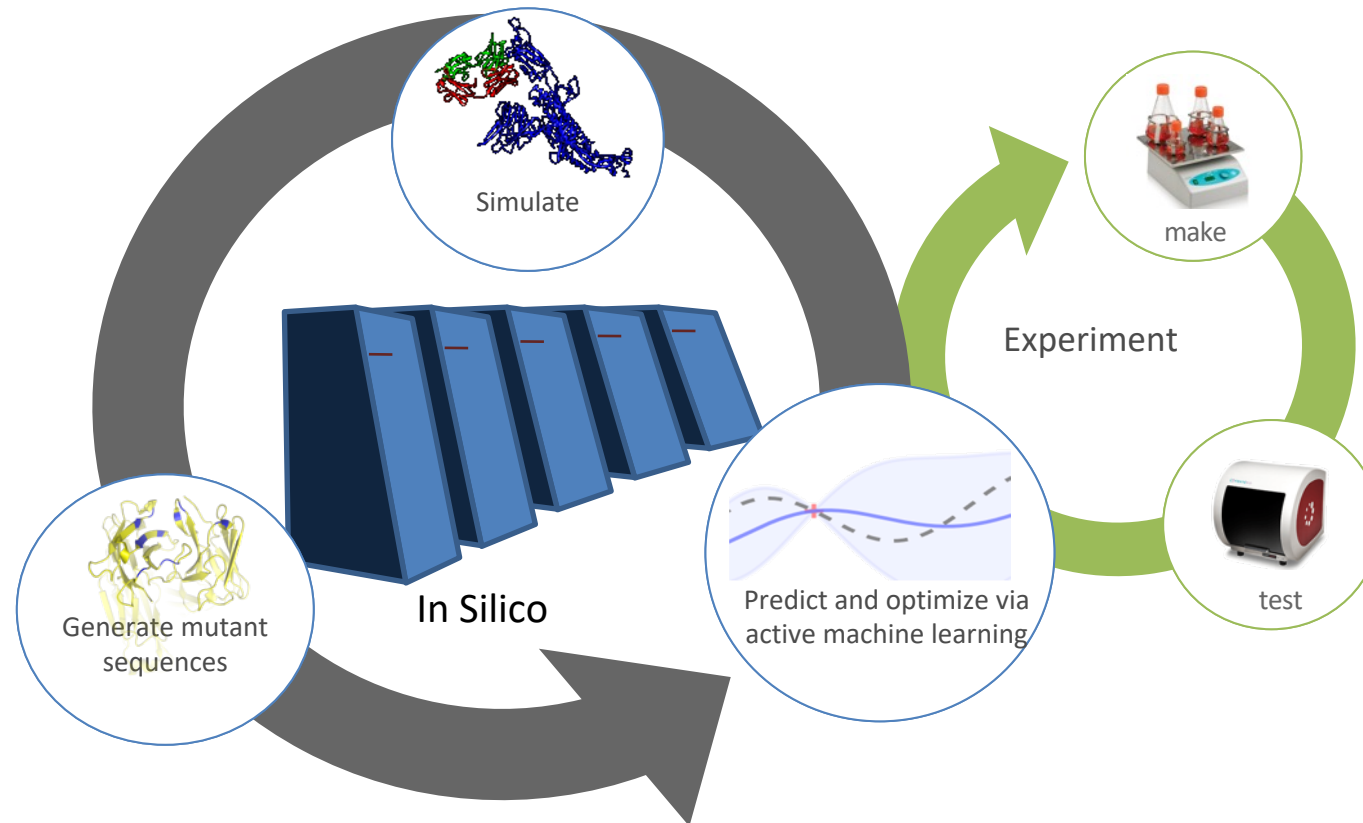
Strong binding is our main target; neutralization objective *may* follow



- In **simulation** and in the **laboratory**, we can ask questions like:
 - How strongly does the antibody bind its target? dG (binding free energy) or K_D (rate const.)
 - How does this change as we mutate the antibody? ddG (mutational change in dG)



Platform software and active machine learning support these simulation and experimental tools

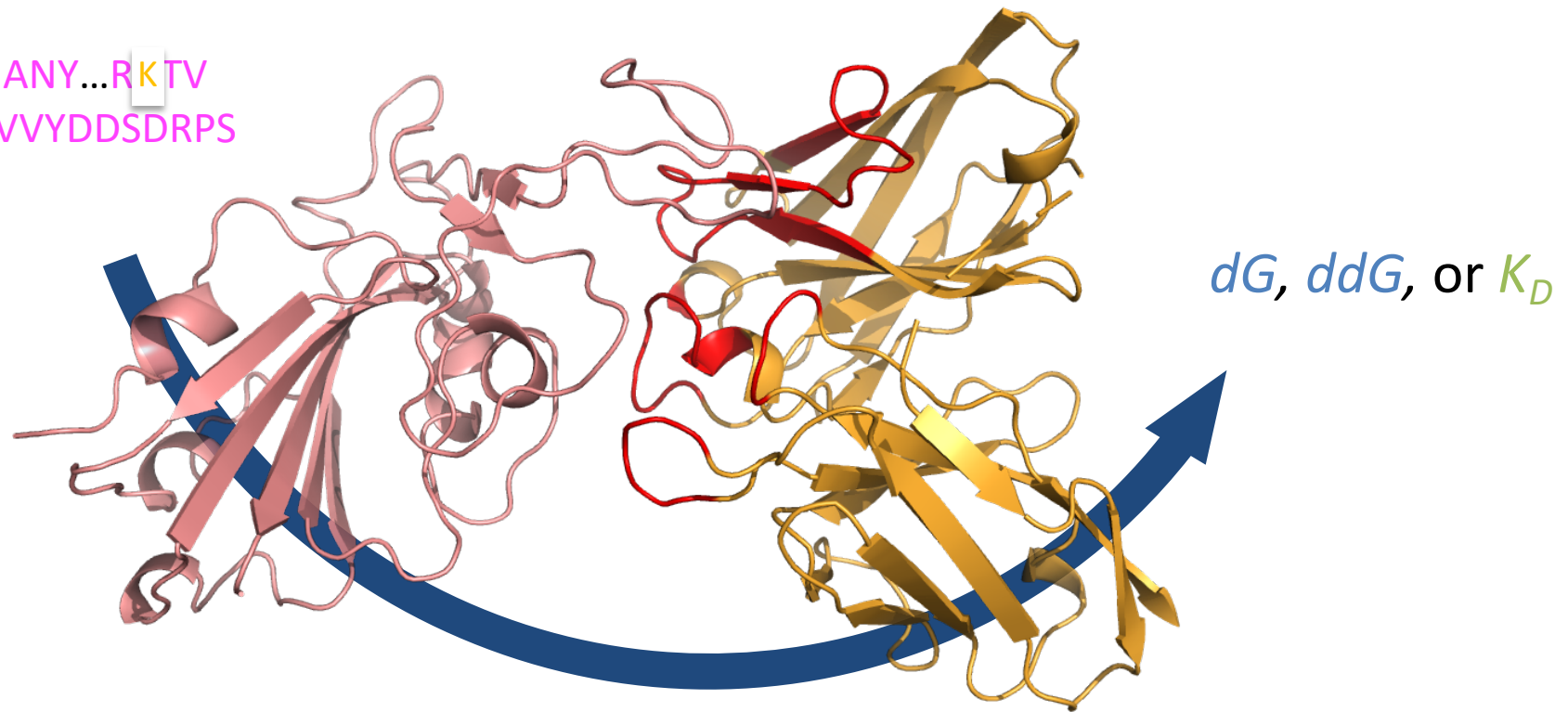


Pose the design problem as active learning

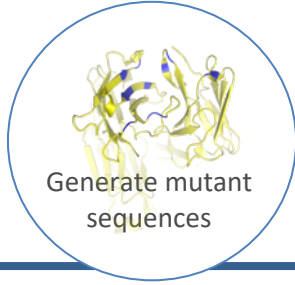
- Improve the antibody sequence by iteratively selecting antibodies from a discrete set and evaluating them

> m396 mutable residues

...GTFSSYTIS...WMGGSPILGIAN...RKT
MGGMDV.../...NIGSKSVH...LVVYDDSDRPS
...QVWDSSSDY



Enumerate many antibody designs

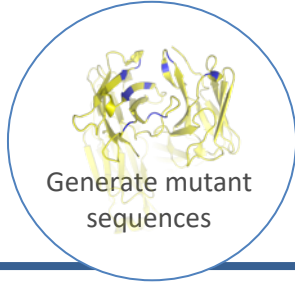


> m396 mutable residues

...GTFSSYTIS...WMGGSPILGIANY...RKTV
MGGMDV.../...NIGSKSVH...LVVYDDSDRPS
...QVWDSSSDY

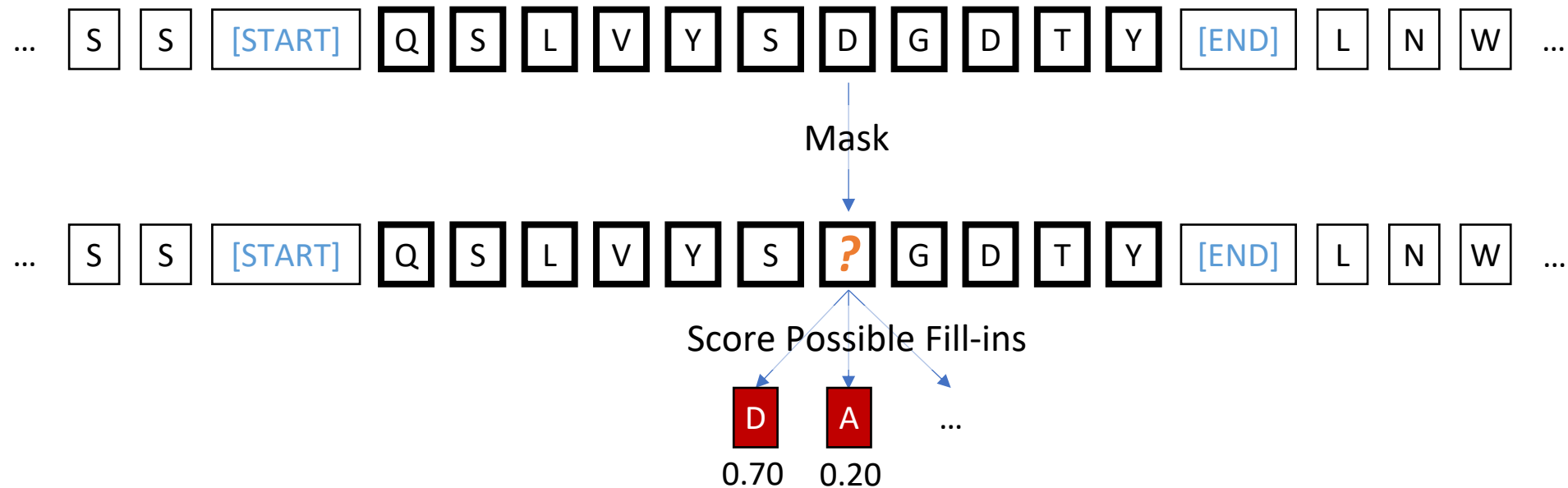
- Generators for novel sequences have so far been mostly tabular
 - Based on frequency of “typical” mutational “swaps”
 - OR based on expensive, high-fidelity calculations of single changes to **template antibody** in hypothesized complex with **SARS-CoV-2 spike**.
- This works all right, but can lead you to unrealistic sequence designs
 - Downstream problems in manufacturability, etc. are major concerns

More realistic antibody sequences via language modeling

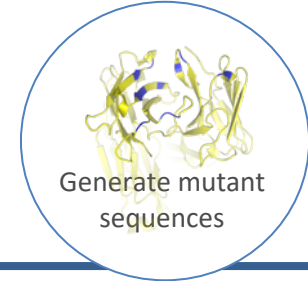


- Use a transformer model to learn to fill “masked” amino acids in the antibody sequence

Annotated L1 from s230



Our models learn to produce reasonable antibodies

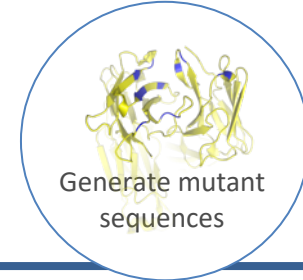


mask and predict 3 central amino acids of s230's L1 "loop"

Mask

<input type="checkbox"/> S	S	—	Go
<input type="checkbox"/> L	L	—	Go
<input checked="" type="checkbox"/> V	[MASK]	V, L, A, I, G (0.866, 0.070, 0.023, 0.023, 0.007)	Go
<input checked="" type="checkbox"/> Y	[MASK]	Y, H, F, S, N (0.530, 0.342, 0.044, 0.027, 0.018)	Go
<input checked="" type="checkbox"/> S	[MASK]	S, T, R, G, N (0.898, 0.034, 0.031, 0.016, 0.007)	Go
<input type="checkbox"/> D	D	—	Go
<input type="checkbox"/> G	G	—	Go
<input type="checkbox"/> D	D	—	Go

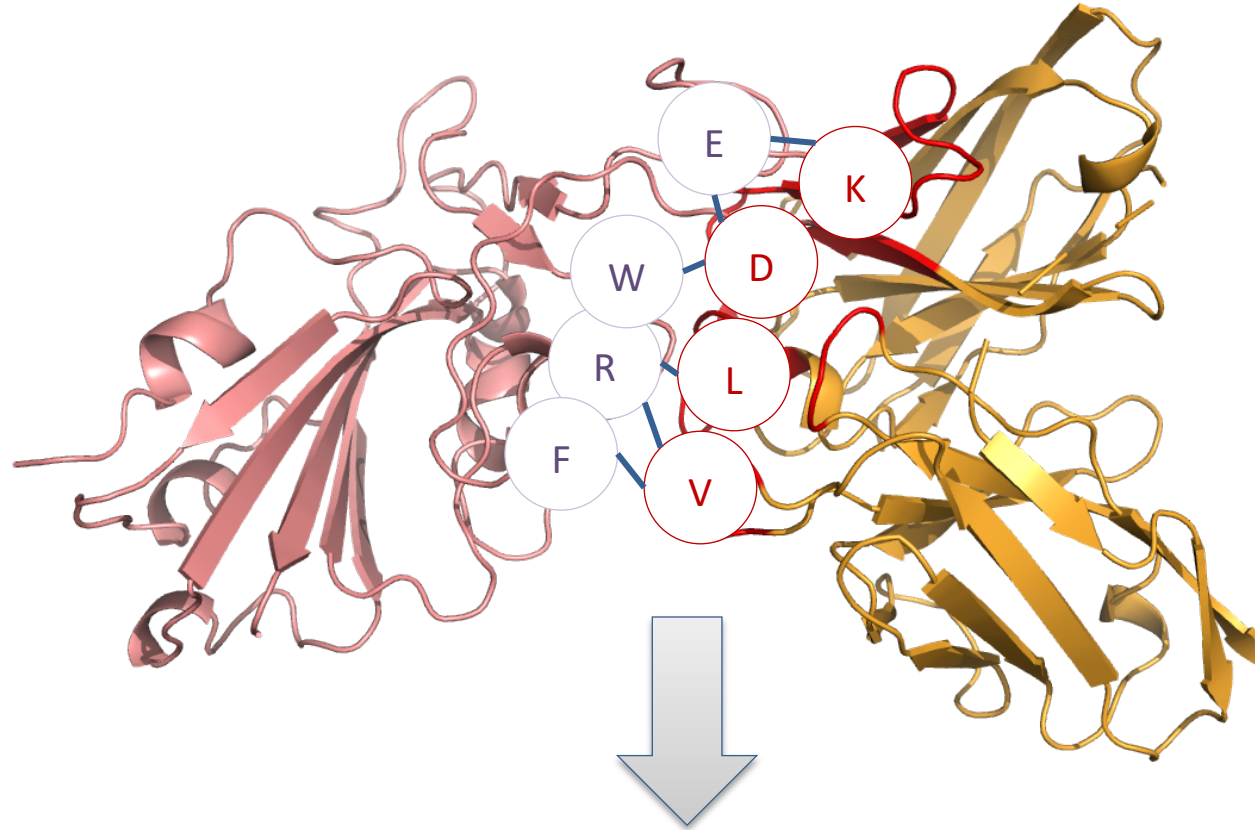
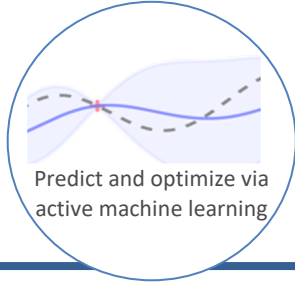
Our models learn to produce reasonable antibodies



mask and predict all 16 amino acids of s230's L1 "loop"

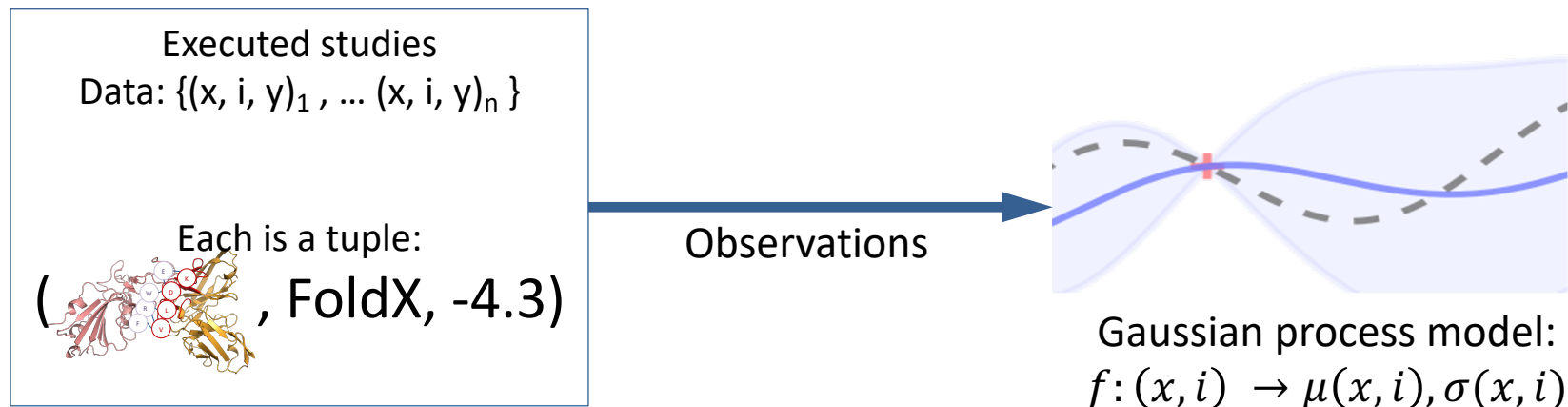
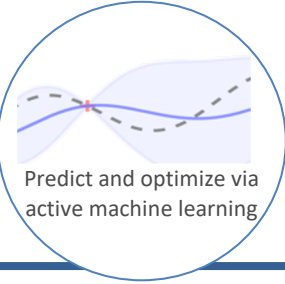
<input type="checkbox"/> C		C		Predict
<input type="checkbox"/> [START_L1]	Mask	[START_L1]		
<input checked="" type="checkbox"/> R	Structurally-identified contacts	[MASK]	R, F, W, T, S (0.815, 0.074, 0.035, 0.032, 0.019)	
<input checked="" type="checkbox"/> S		[MASK]	S, F, A, T, Y (0.958, 0.019, 0.011, 0.005, 0.003)	
<input checked="" type="checkbox"/> S		[MASK]	S, T, R, N, G (0.919, 0.022, 0.021, 0.016, 0.012)	
<input checked="" type="checkbox"/> Q		[MASK]	Q, L, H, R, L (0.891, 0.034, 0.030, 0.017, 0.011)	
<input checked="" type="checkbox"/> S		[MASK]	S, C, R, T, N (0.870, 0.073, 0.031, 0.009, 0.005)	
<input checked="" type="checkbox"/> L		[MASK]	L, F, R, I, F (0.976, 0.006, 0.006, 0.004, 0.002)	
<input checked="" type="checkbox"/> V		[MASK]	V, L, A, I, E (0.814, 0.115, 0.018, 0.017, 0.015)	
<input checked="" type="checkbox"/> Y		[MASK]	H, Y, F, S, N (0.483, 0.391, 0.026, 0.023, 0.019)	
<input checked="" type="checkbox"/> S		[MASK]	S, F, T, G, N (0.817, 0.057, 0.047, 0.023, 0.021)	
<input checked="" type="checkbox"/> D		[MASK]	D, L, G, A, E (0.906, 0.044, 0.016, 0.011, 0.007)	
<input checked="" type="checkbox"/> G		[MASK]	G, L, V, A, E (0.970, 0.009, 0.004, 0.004, 0.004)	
<input checked="" type="checkbox"/> D		[MASK]	N, S, D, K, T (0.884, 0.035, 0.019, 0.017, 0.011)	
<input checked="" type="checkbox"/> T		[MASK]	T, I, S, P, N (0.947, 0.020, 0.011, 0.008, 0.007)	
<input checked="" type="checkbox"/> Y		[MASK]	Y, F, H, S, C (0.827, 0.067, 0.050, 0.016, 0.009)	
<input checked="" type="checkbox"/> L		[MASK]	L, F, V, S, I (0.965, 0.022, 0.007, 0.003, 0.001)	
<input checked="" type="checkbox"/> N		[MASK]	N, S, H, T, Y (0.820, 0.044, 0.032, 0.032, 0.018)	
<input type="checkbox"/> [END_L1]		[END_L1]		
<input type="checkbox"/> W		W		

To predict how an antibody sequence will bind, we use a structure-based representation of the interactions

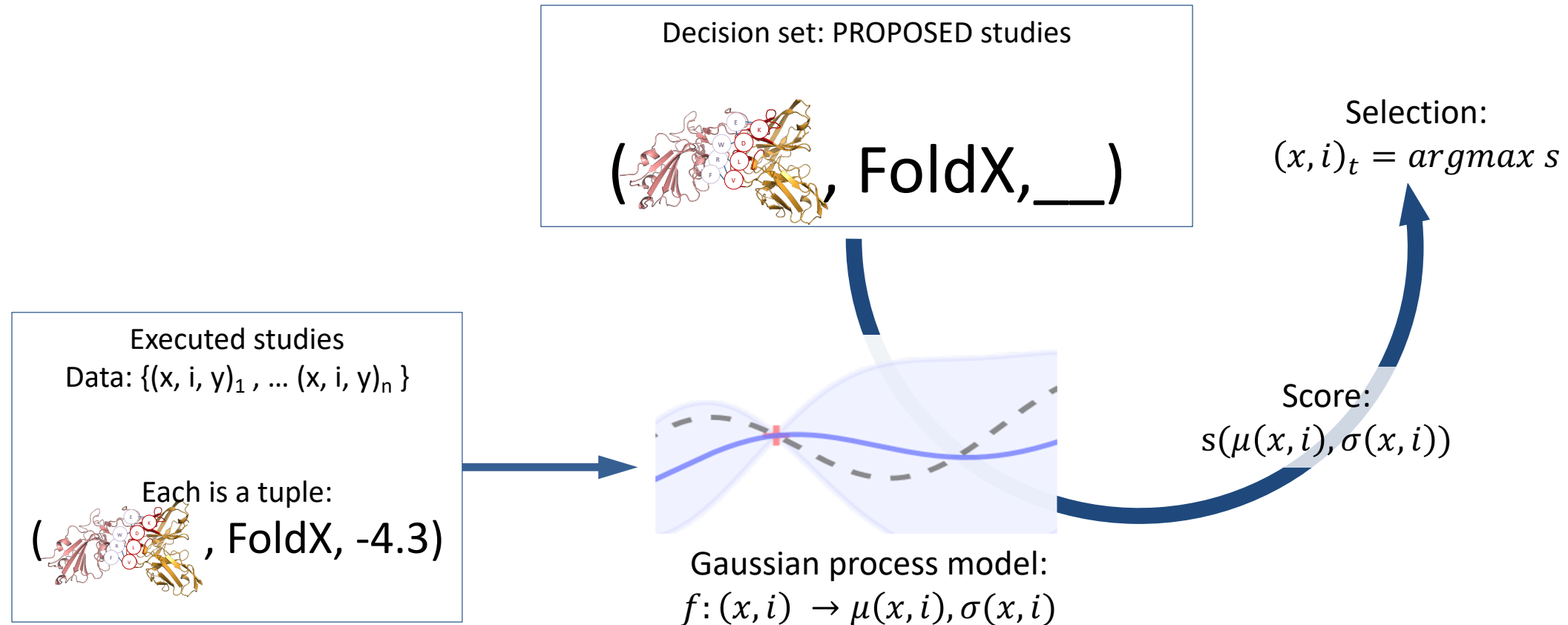
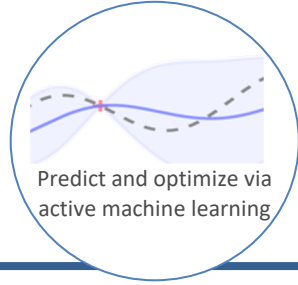


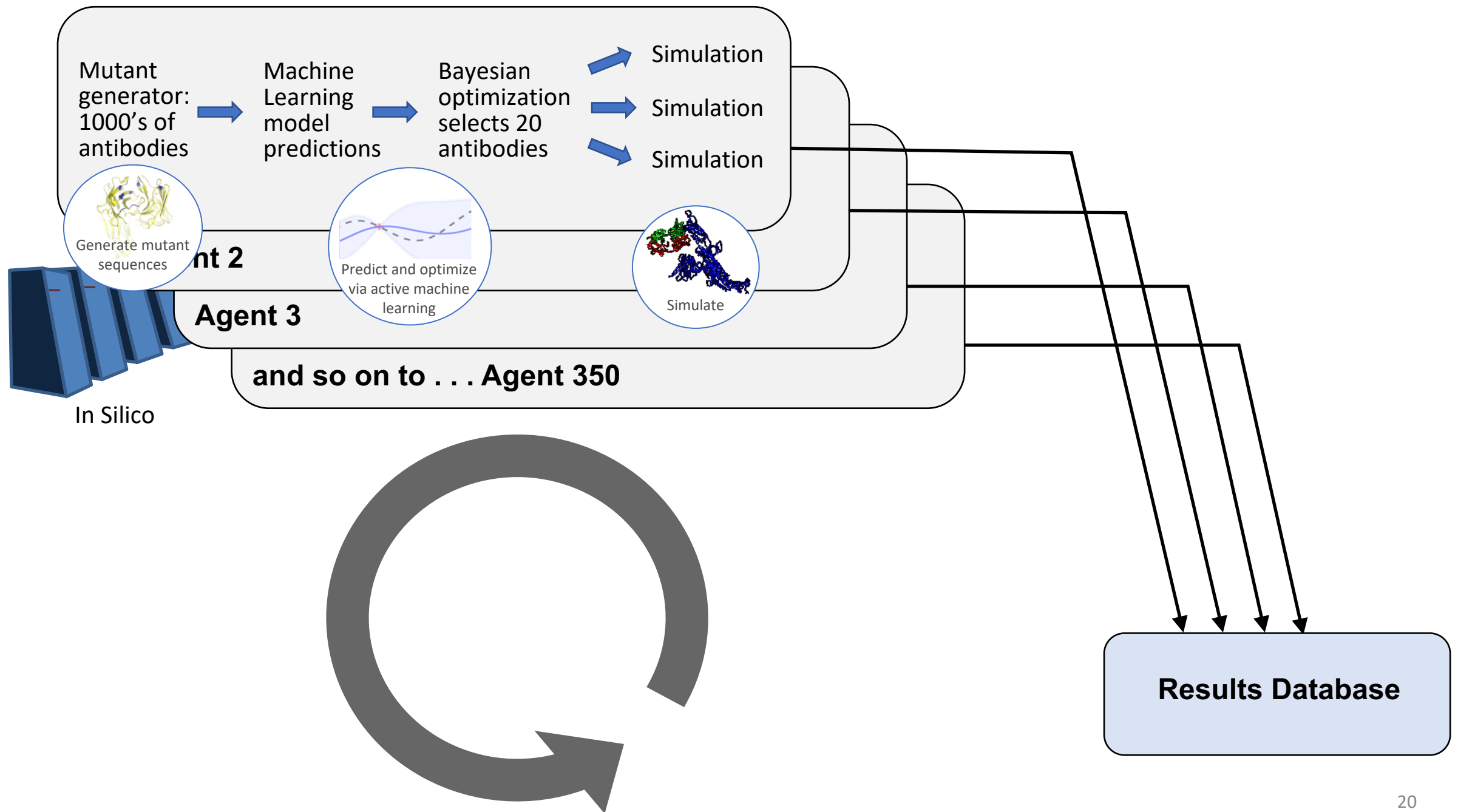
$\mathbf{x} = [0, 1, 0, 2, 0, 0, 1, \dots]$
Vector of interaction type counts

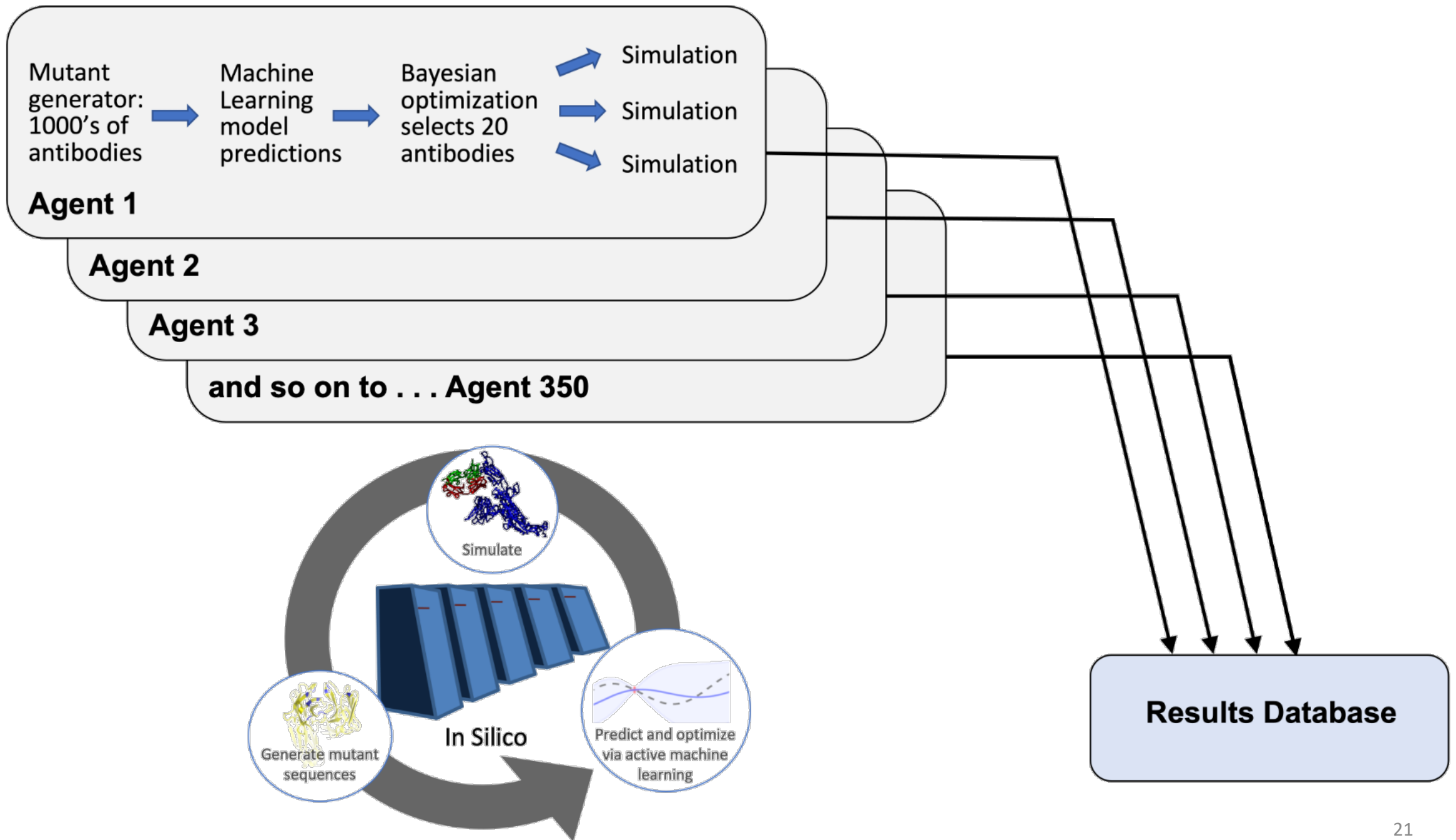
Represented in feature space, binding free energy estimates feed into a multi-fidelity Gaussian process model

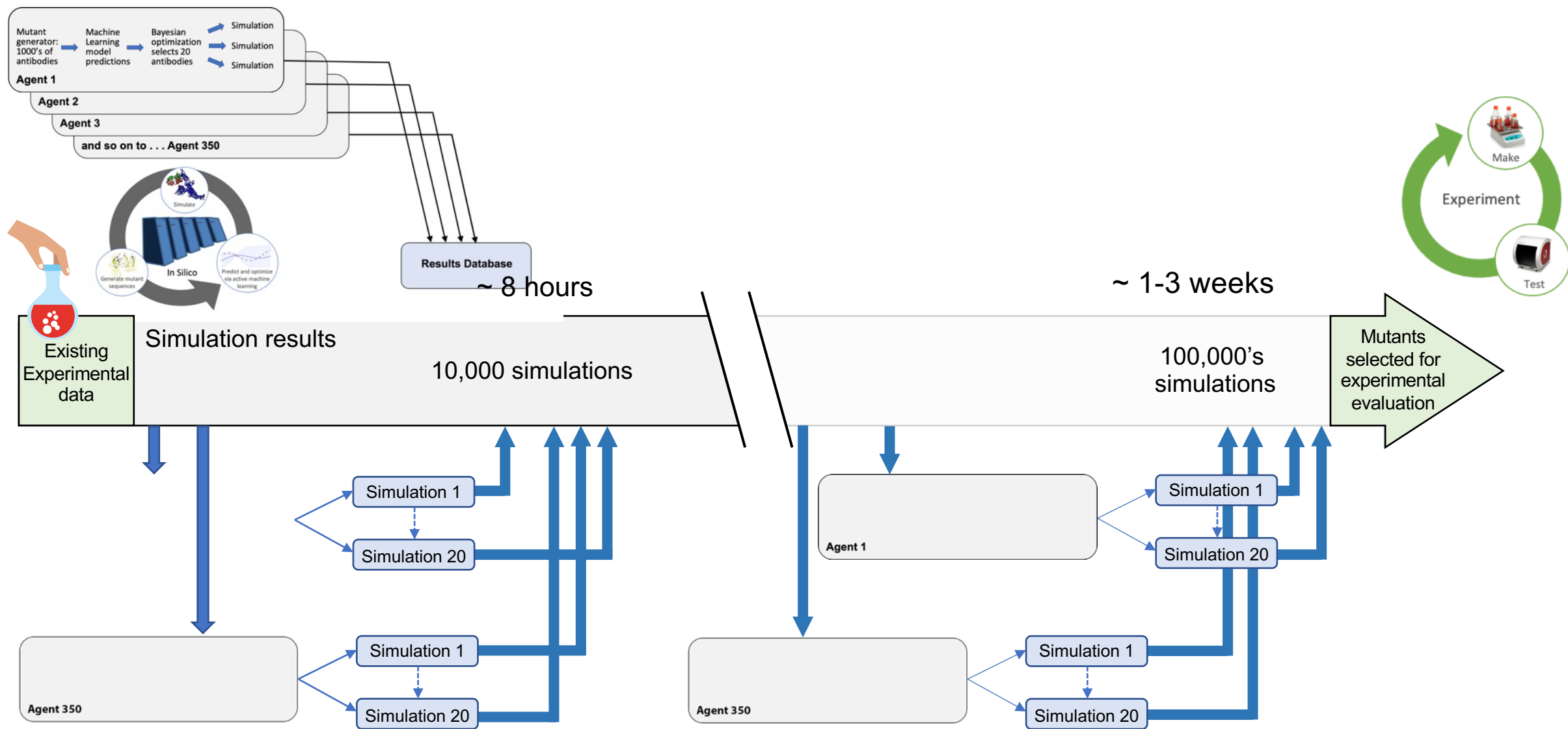


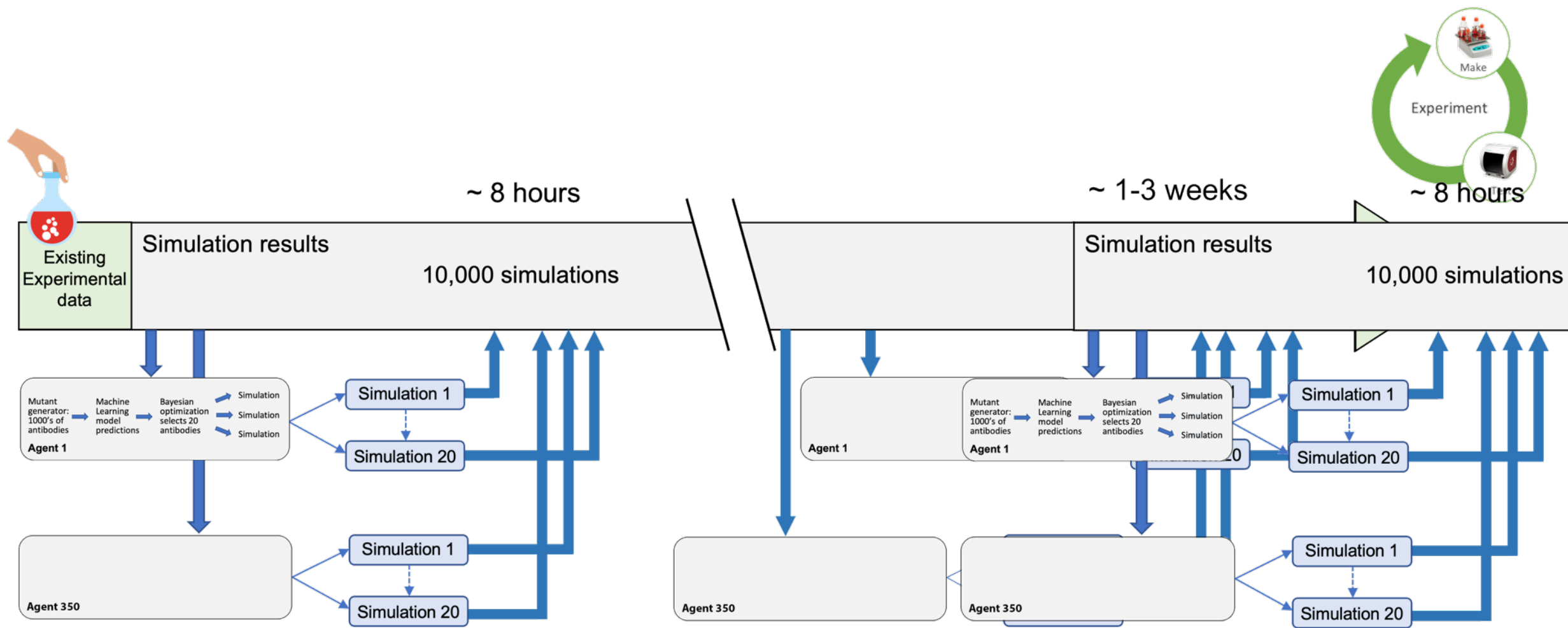
The next set of simulations is selected via Bayesian optimization using the Gaussian process model







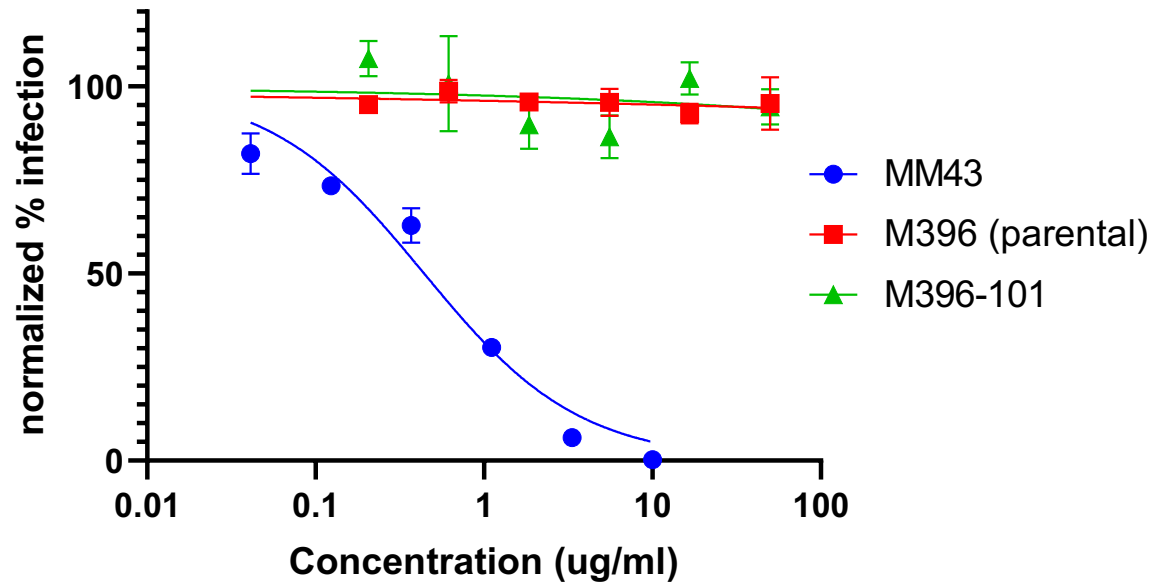




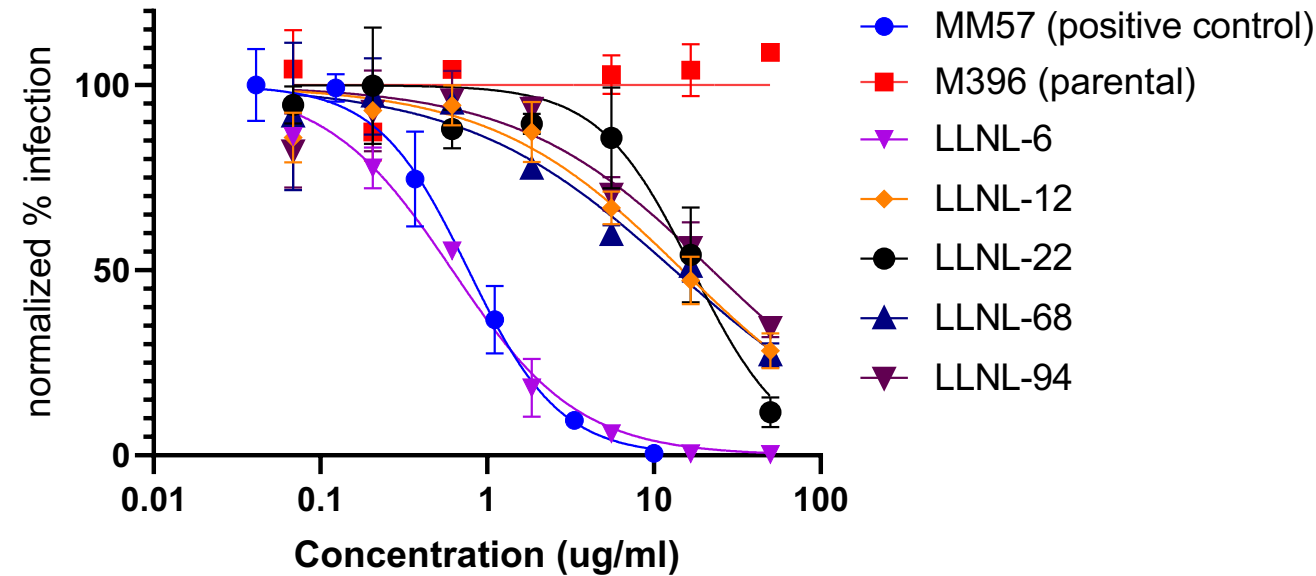
Several of our m396-derived antibodies inhibit VSV-SARS-CoV-2 virus



mAb101-Iteration 1
rVSV-SARS-CoV-2 S-GFP



VSV-SARS-CoV-2-GFP M396-R2 hits

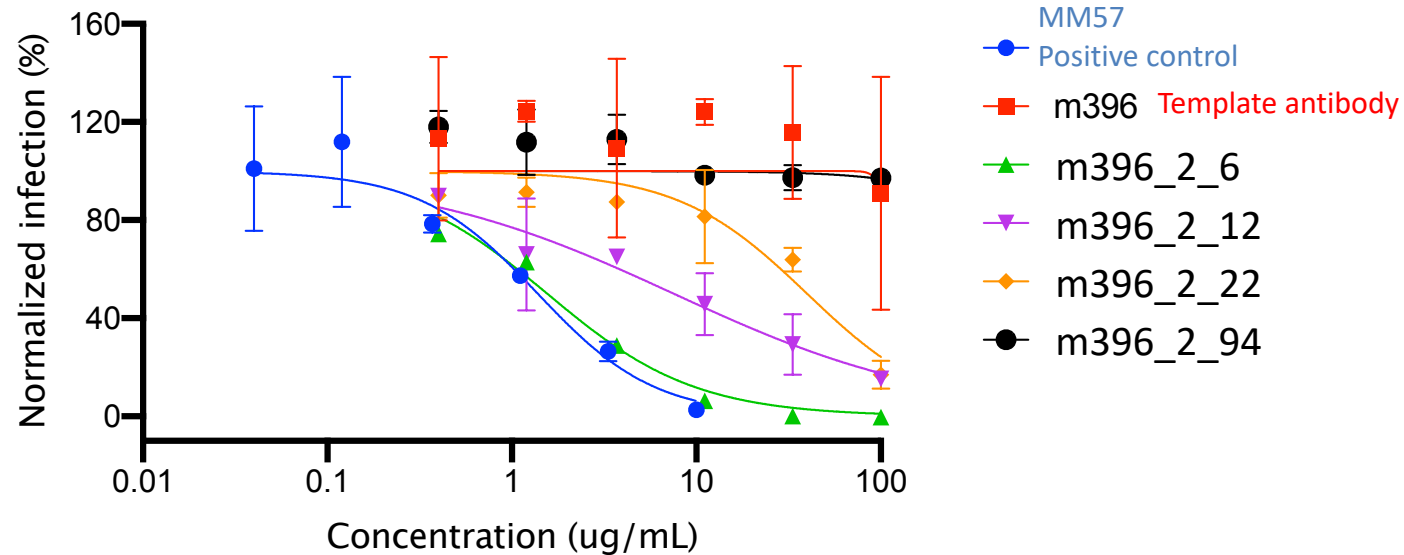


Recall:
down and left is more virus
neutralized needing less
antibody = better

We confirmed these m396-derived antibodies neutralize authentic SARS-CoV-2 virus in our BSL-3 facility



SARS-CoV-2-GFP Neutralization



This work is the product of a growing multidisciplinary team

- LLNL:
Daniel Faissol, Adam Zemla, Ed Lau, Fangqiang Zhu, John Goforth, Denis Vashchenko, Mary Silva, Rebecca Haluska, Claudio Santiago, Sam Nguyen, Drew Bennett, Emilia Grzesiak, Brent Segelke, Feliza Bourguet, Victoria Lao, Monica Borucki, Dina Weilhammer, Jacky Lo, Nicole Collette, Kathryn Arrildt, and Magdalena Franco (now ThermoFisher)
- Sandia NL:
Brooke Harmon, Oscar Negrete, Max Stefan
- Generous computer time and support from LC!
 - Catalyst, early access to Mammoth
 - Workflow enablement (database) and Sina (database interface) groups are critical to our ongoing success

PyTorch, GPyTorch, BioPython

Maestro, Sina, Improv

FoldX, RosettaFlex





This document was prepared as an account of work sponsored by an agency of the United States government. Neither the United States government nor Lawrence Livermore National Security, LLC, nor any of their employees makes any warranty, expressed or implied, or assumes any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States government or Lawrence Livermore National Security, LLC. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States government or Lawrence Livermore National Security, LLC, and shall not be used for advertising or product endorsement purposes.