

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

Thomas Desautels
Staff Scientist, LLNL

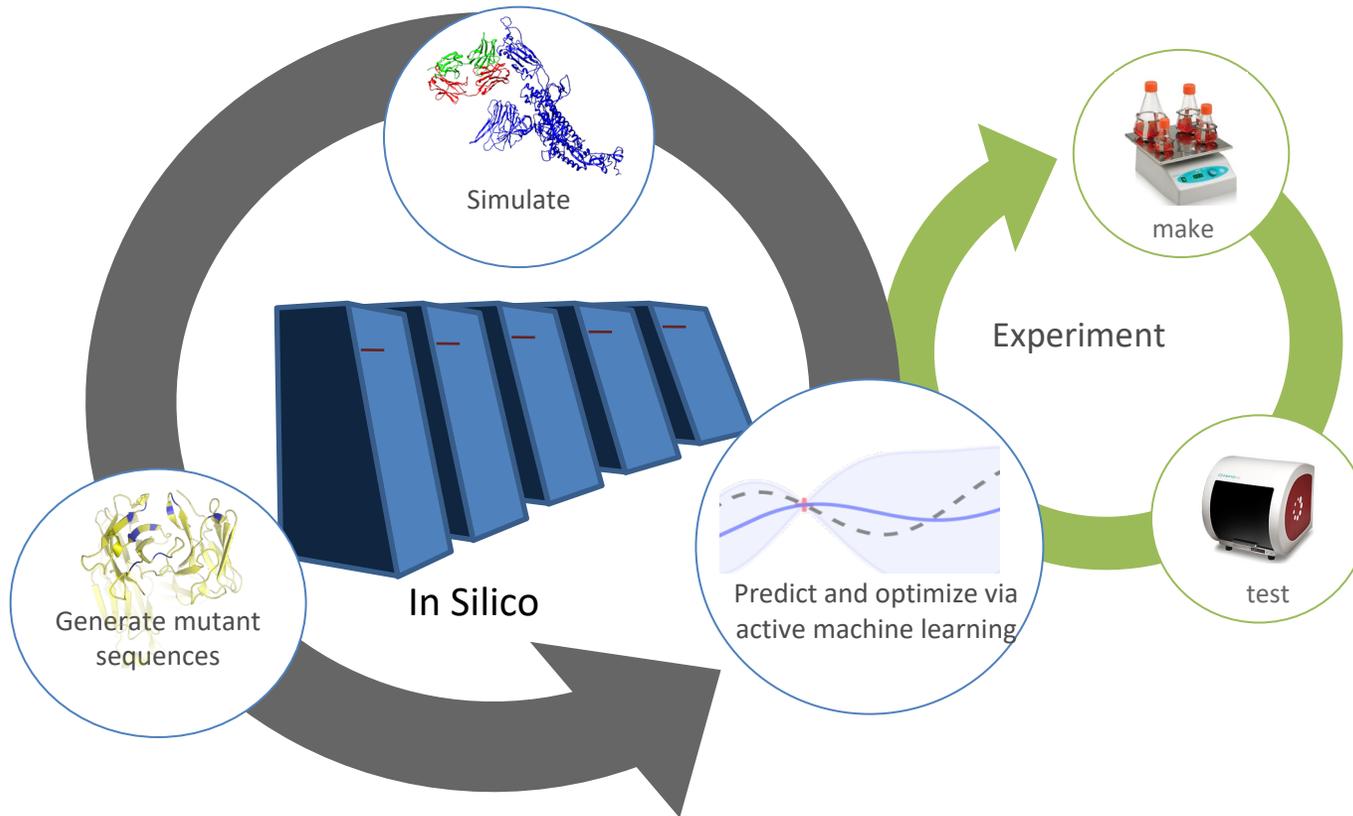
LLNL Data Science Institute Workshop:
AI in Healthcare

March 25, 2021

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

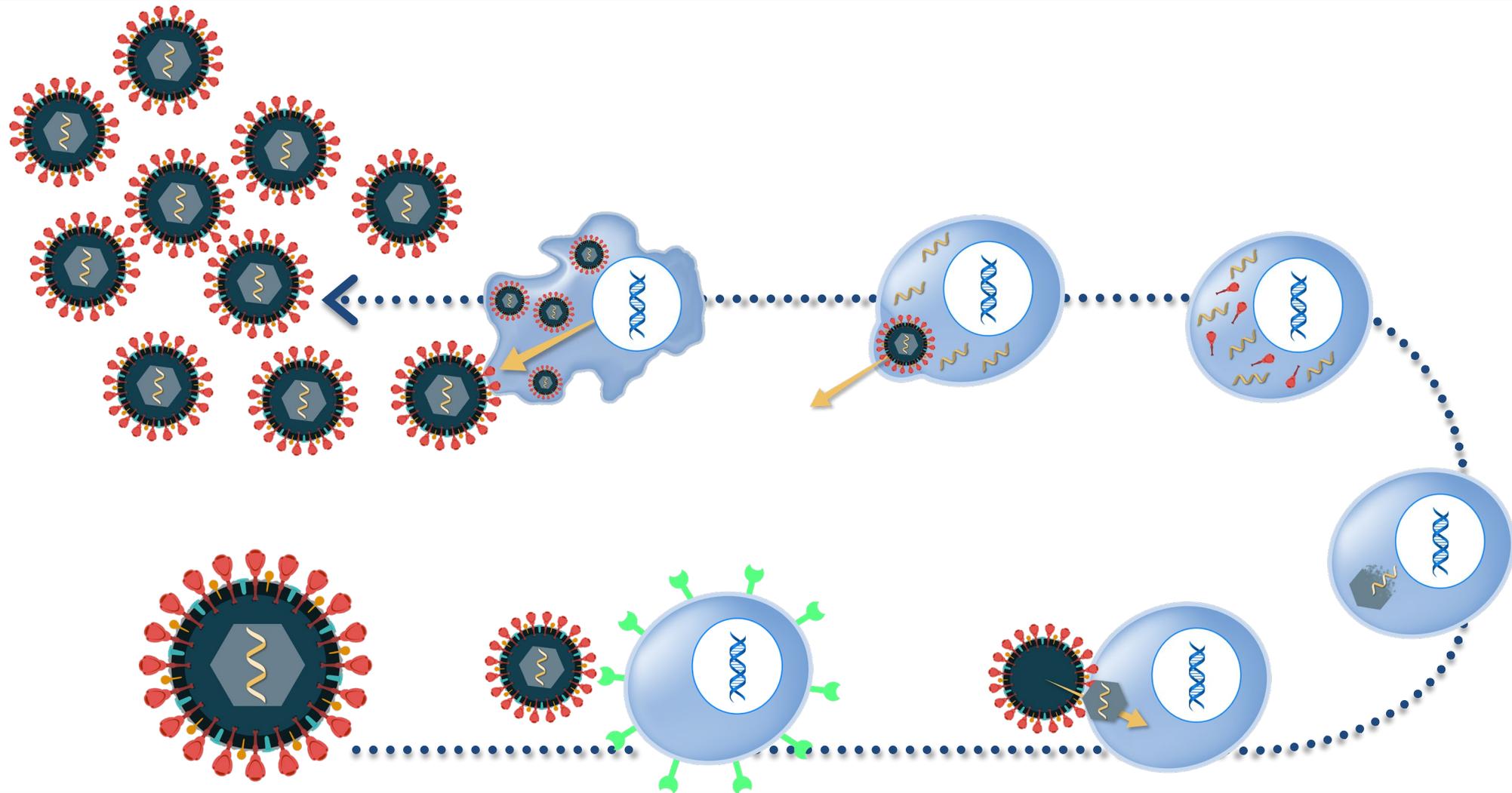


LLNL is developing a machine-learning-driven platform for rapid, rational, design of therapeutic antibodies and vaccine antigens

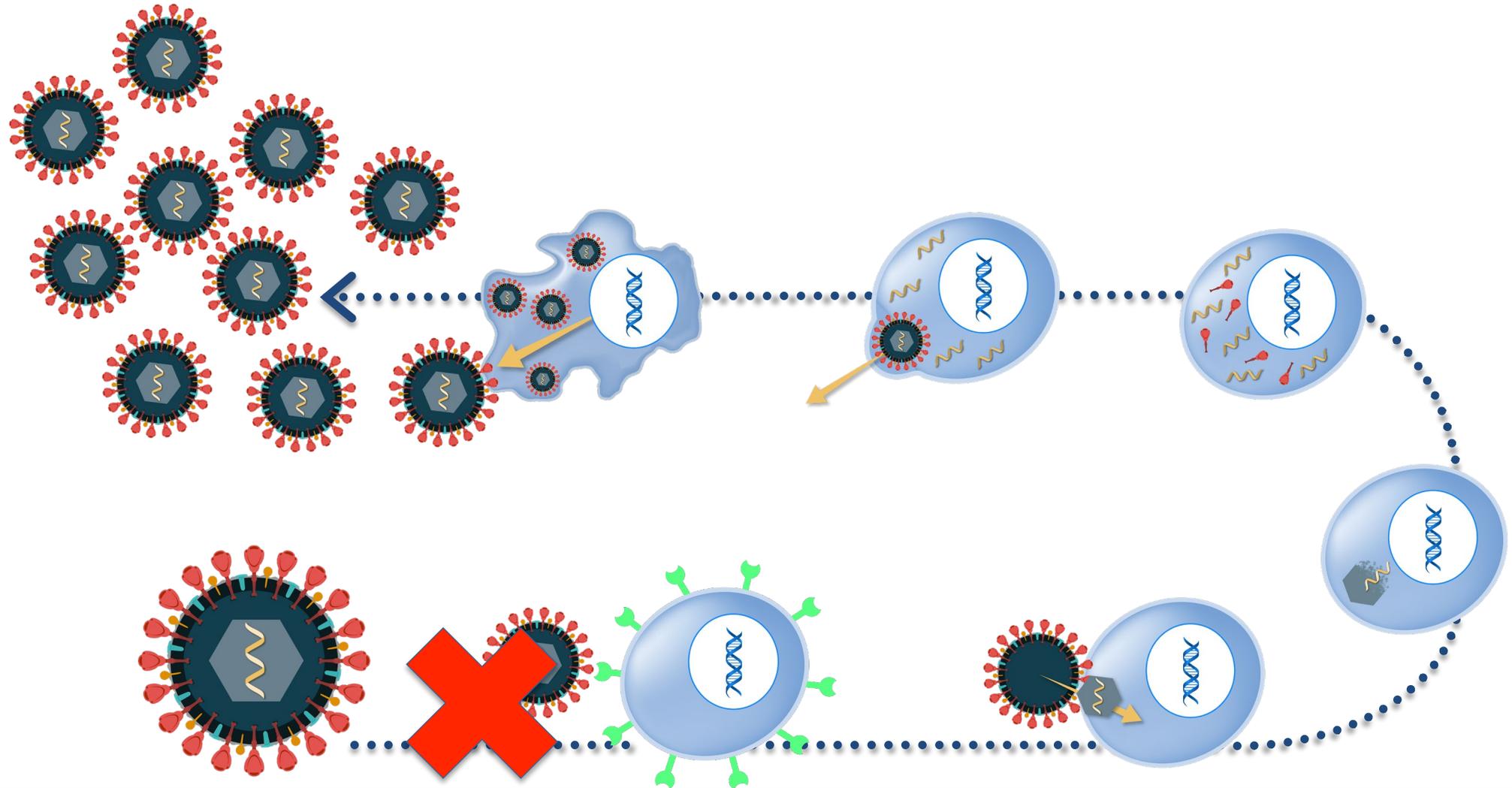


- Our platform is a novel autonomous, ML-driven system that performs in-silico biomolecular design at scale using HPC
- Active learning and Bayesian optimization approach harnesses predictions from molecular simulations and bioinformatic predictors together with high-value, directed experiments
 - > 1 million simulations performed w/ 3 million core hours on HPC to date
- Approach does not require starting from survivor serum sample isolation or library screening
- Demonstrated feasibility in pilot with a major pharma company by re-designing antigens for multiple antibody targets.
- Efforts currently focused on designing monoclonal antibodies against SARS-CoV-2

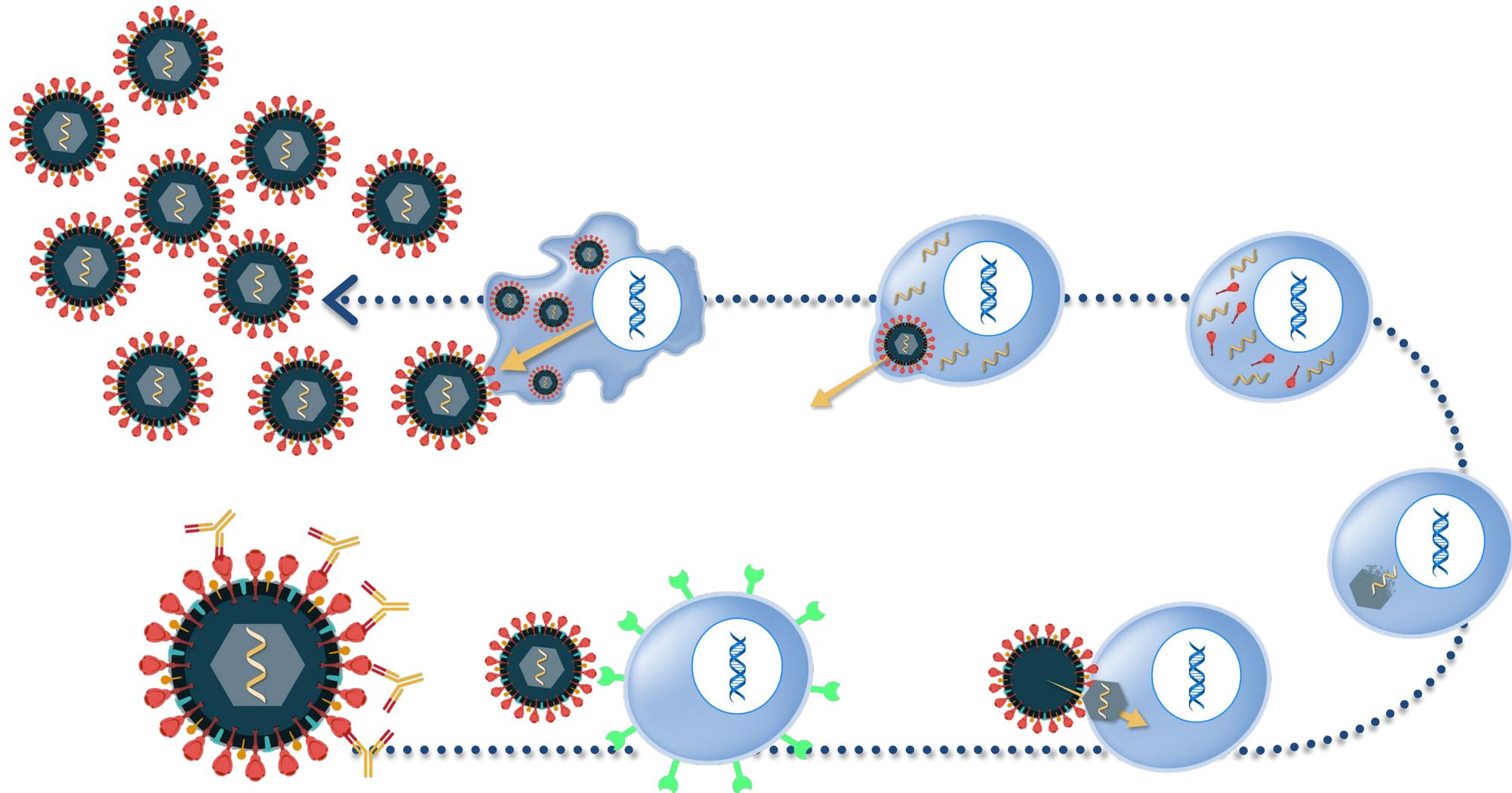
Viruses reproduce by entering and hijacking host cells



If we could stop viral entry, we could stop the viral cycle



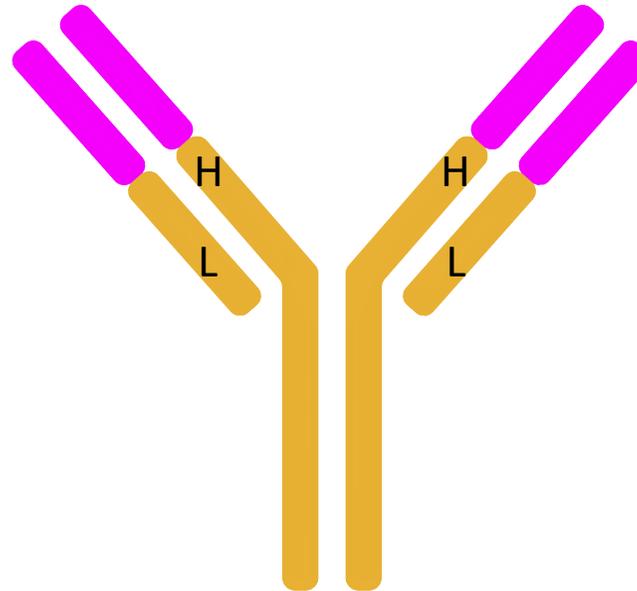
Neutralizing antibodies can stop viral entry



Antibodies are human proteins that specifically and sensitively recognize pathogen targets

> m396 heavy chain

QVQLQQSGAEVKKPGSSVKV SCKASGGTFS
SYTISWVRQAPGQGLEWMGGITPILGIANY
AQKFQGRVTITTTDESTSTAYMELSSLRSEDTA
VYYCARDTVMGGMDVWGQGTTVTVSSAS
TKGPSVFPLAPSSKSTSGGTSALGCLVKDYFP
EPVTVSWNSGALTSGVHTFPAVLQSSGLYSLS
SVVTVPSSSLGTQTYICNVNHKPSNTKVDKK
VEPKSCDKTSPLFVHHHHHHG DYKD
DDDKG



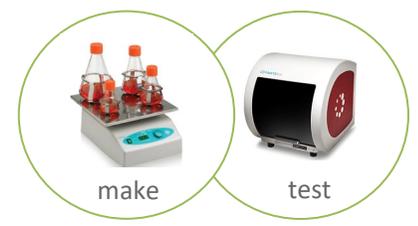
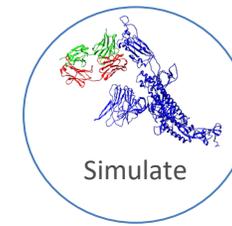
> m396 light chain

SYELTQPPSVSVAPGKTARITCGGNIGSKSV
HWYQQKPGQAPVLVVYDDSDRPSGIPERFS
GSNSGNTATLTISRVEAGDEADYYCQVWDSS
SDYVFGTGTKVTVLGQPKANPTVTLFPPSSE
EFQANKATLVCLISDFYPGAVTVAWKADGSP
VKAGVETTKPSKQSNNKYAASSYLSLTPEQW
KSHRSYSCQVTHEGSTVEKTVAPTECS

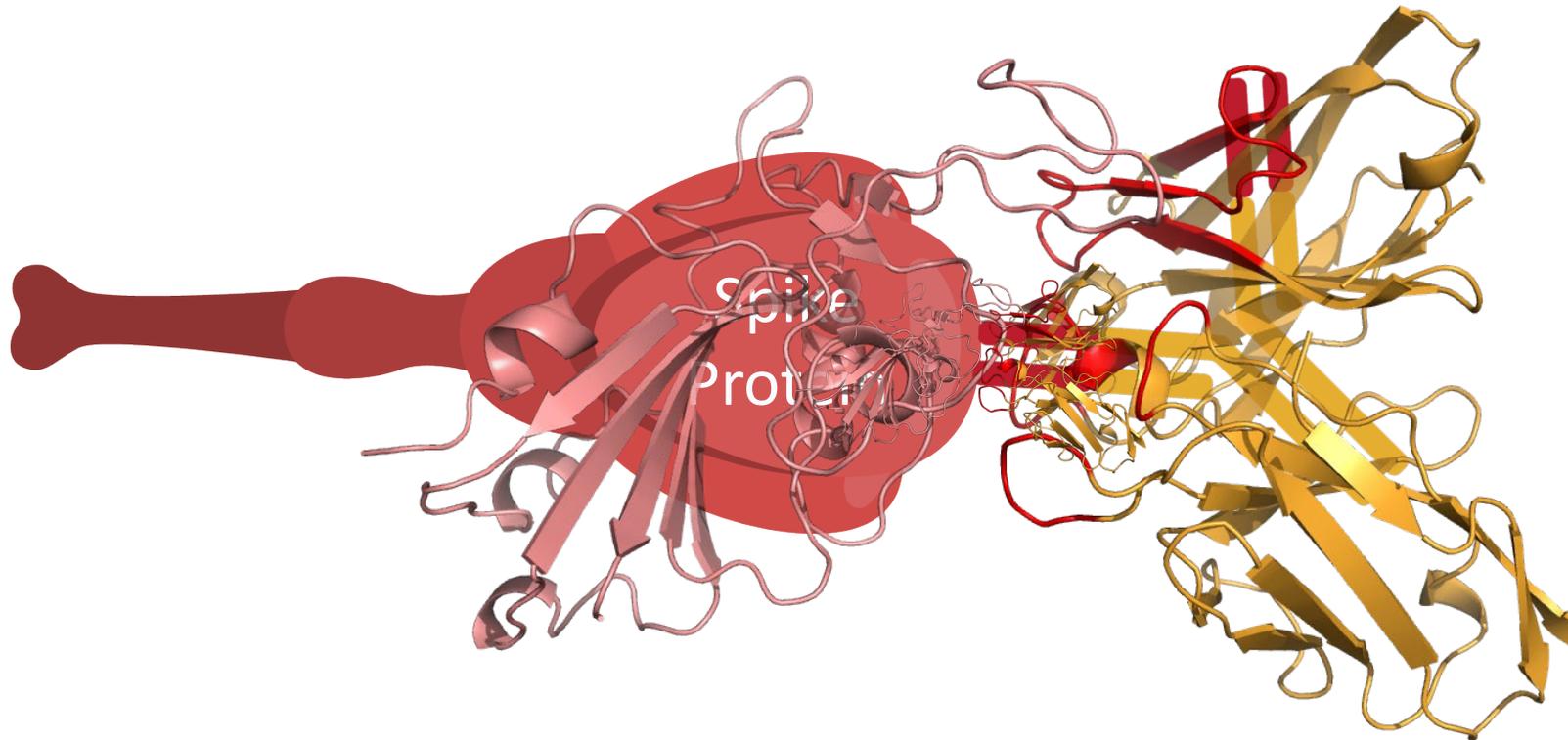
Sequence encodes 3D structure; structure mediates function



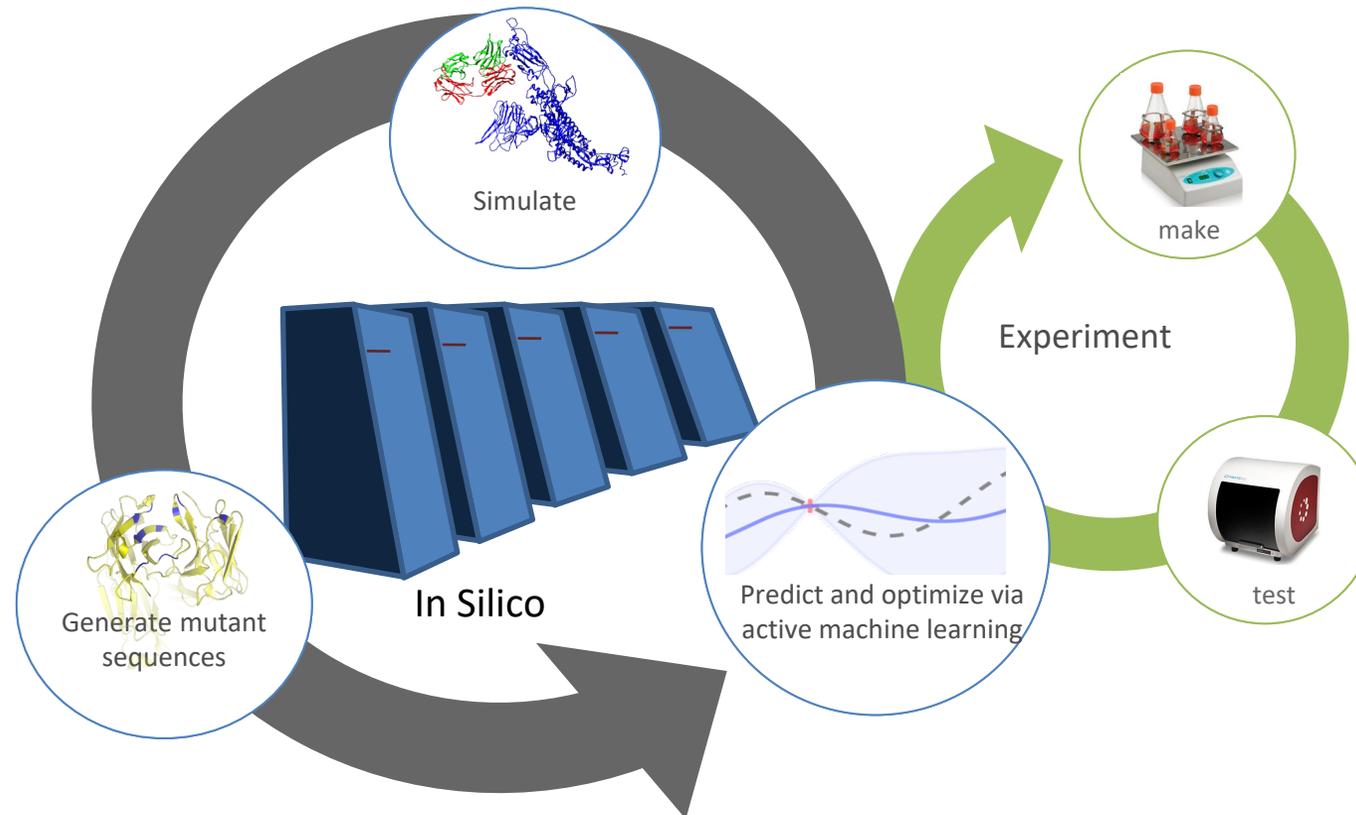
Strong binding is our main target



- 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

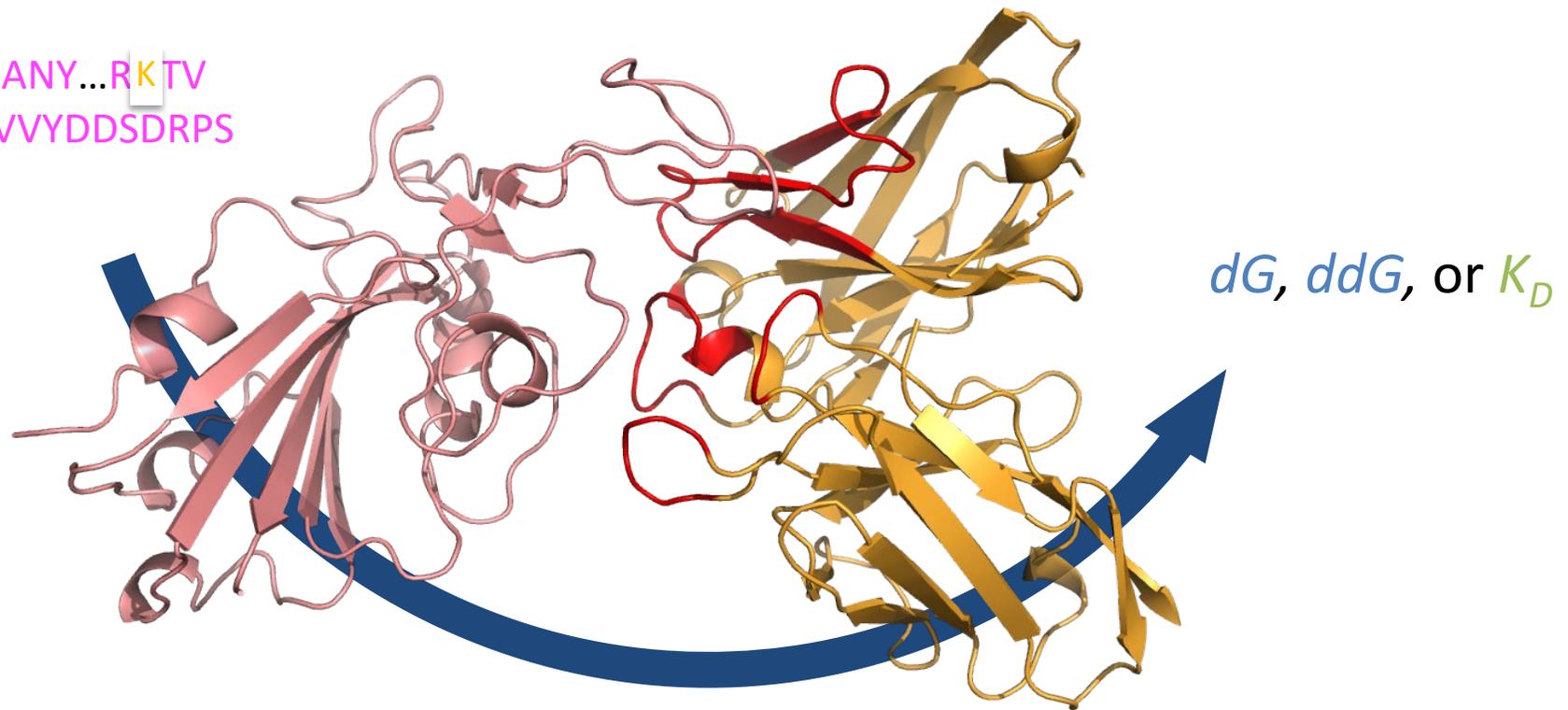


Pose this problem as active learning

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

> m396 mutable residues

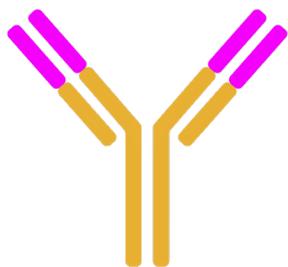
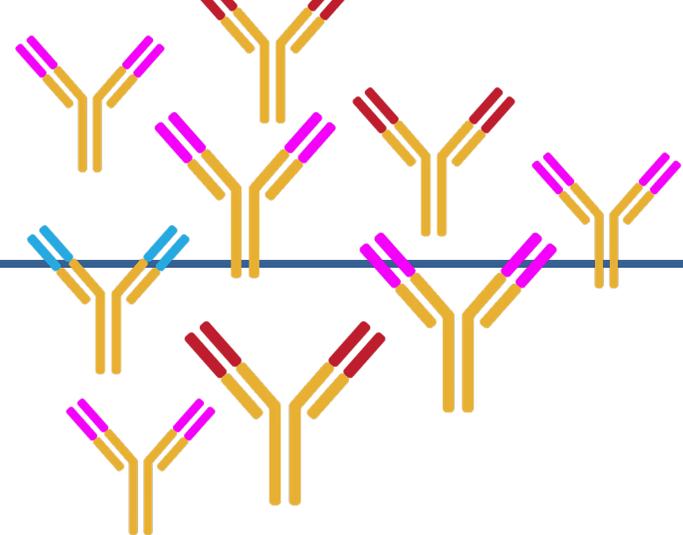
...GTFSSYTIS...WMGGSPILGIANY...RKTV
MGGMDV.../...NIGSKSVH...LVVYDDSDRPS
...QVWDSSSDY



How do you get started?

- “*De novo*” antibody design is a major challenge for the field.
- In the last year, we’ve started from template antibodies that neutralize the closely related SARS-CoV-1 virus (early 2000’s), but don’t neutralize SARS-CoV-2.
- Likely several mutational steps away from any related, SARS-CoV-2 neutralizers, impractical to search this in the lab
- **HOWEVER**, we have good reason to believe that the interaction we’re trying to “restore” should be neutralizing for SARS-CoV-1 too.

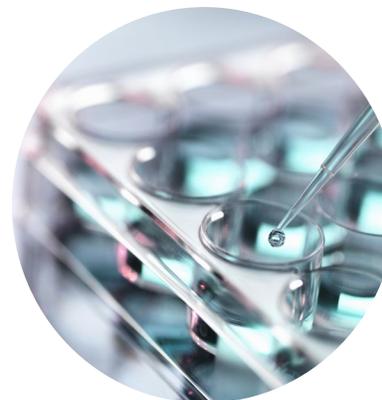
The design space is vastly larger than what we can simulate or test



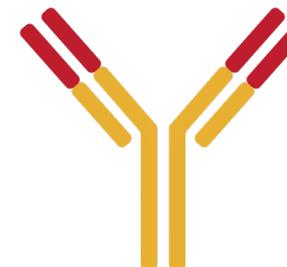
CoV-1 +
changes
 $\sim 10^{30}$



Computer
Simulations
1,000,000

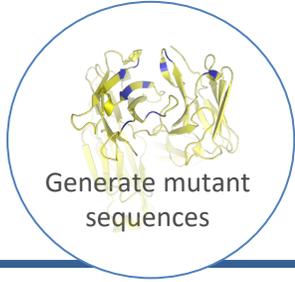


Laboratory
Experiments
100-1,000



CoV-2
Need just one!

Enumerate many antibody designs

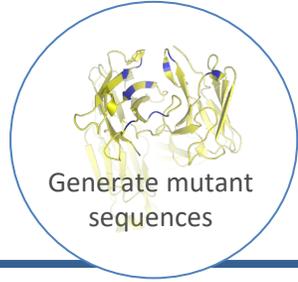


> m396 mutable residues

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

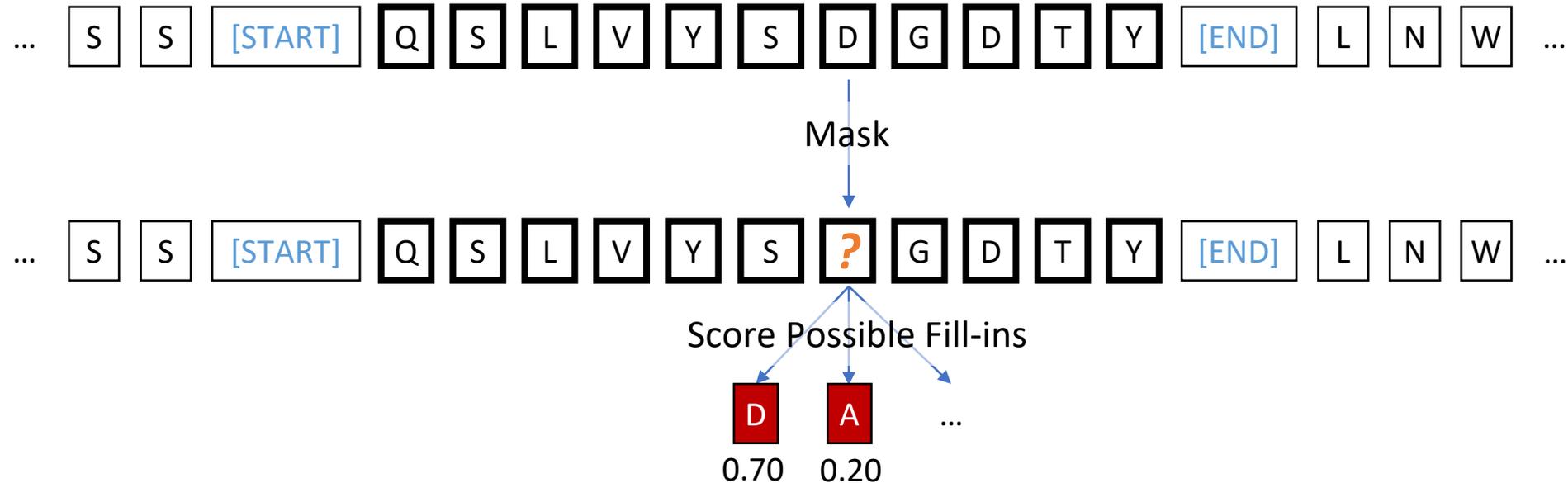
- Generators for novel sequences have so far been mostly tabular
 - Based on frequency of “typical” mutational “swaps”
 - OR based on expensive, high-fidelity calculations of single changes to **template antibody** in hypothesized complex with **SARS-CoV-2 spike**.
- This works all right, but can lead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

More realistic antibody sequences via language modeling

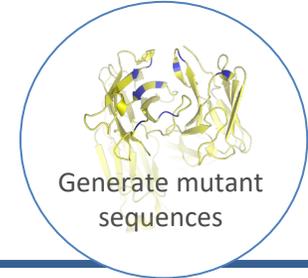


- Use a transformer model to learn to fill “masked” amino acids in the antibody sequence

Annotated L1 from s230



Our models learn to produce reasonable antibodies



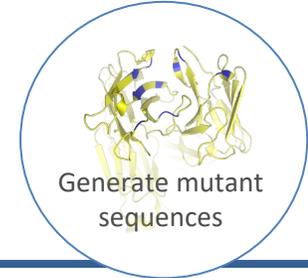
mask and predict 3 central amino acids of s230's L1 "loop"

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

The interface shows a 'Mask' operation on a protein sequence. On the left, a list of amino acids has checkboxes: S, L, V, Y, S, D, G, D. The V, Y, and S options are checked and highlighted with a blue box. An arrow points from this box to the 'Mask' column in the table. The table shows the masked sequence with predicted amino acids and their probabilities. The predicted amino acids are V, L, A, I, G for the first masked position, Y, H, F, S, N for the second, and S, T, R, G, N for the third. The probability for Y is 0.530, which is underlined. A blue arrow points from the V, L, A, I, G prediction back to the V, Y, S checkboxes.

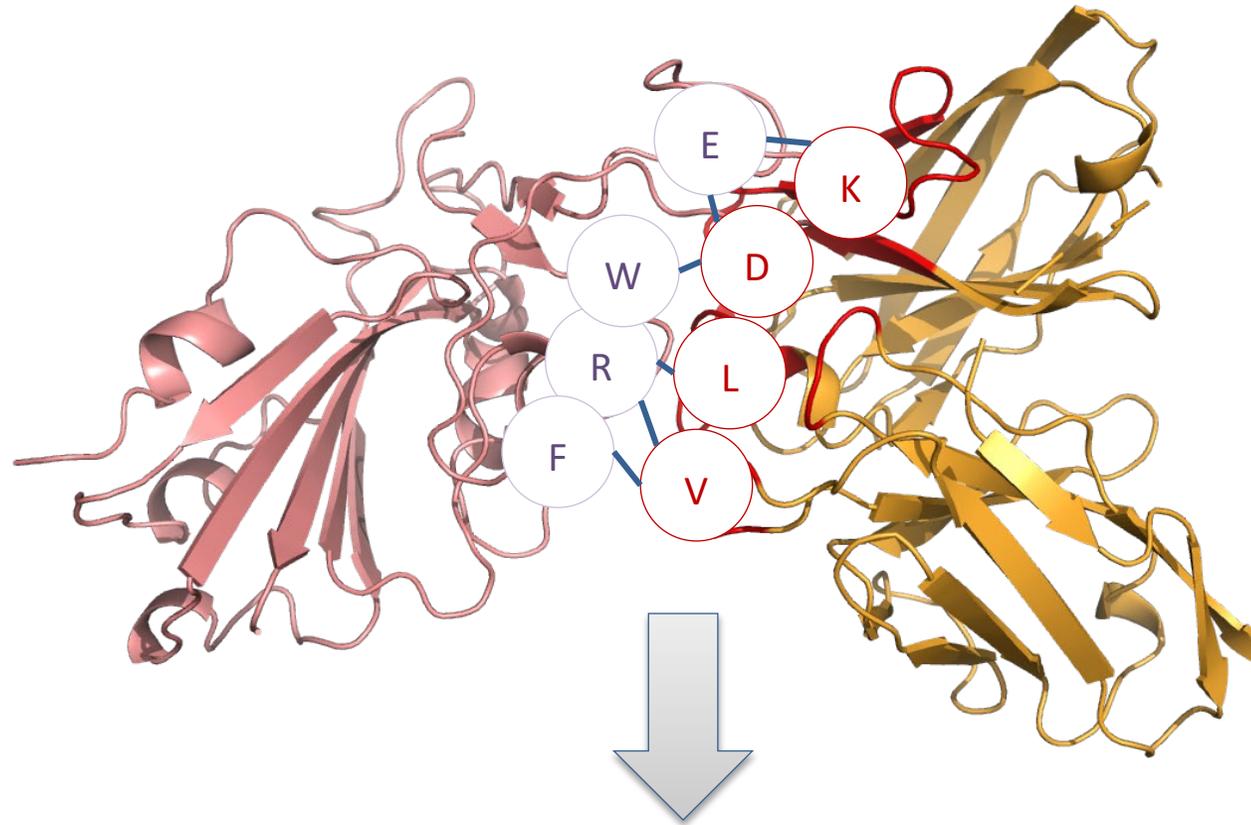
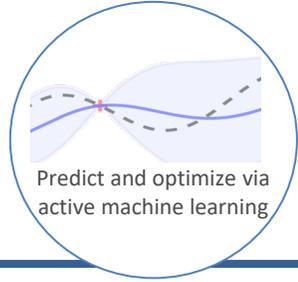
Our models learn to produce reasonable antibodies



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

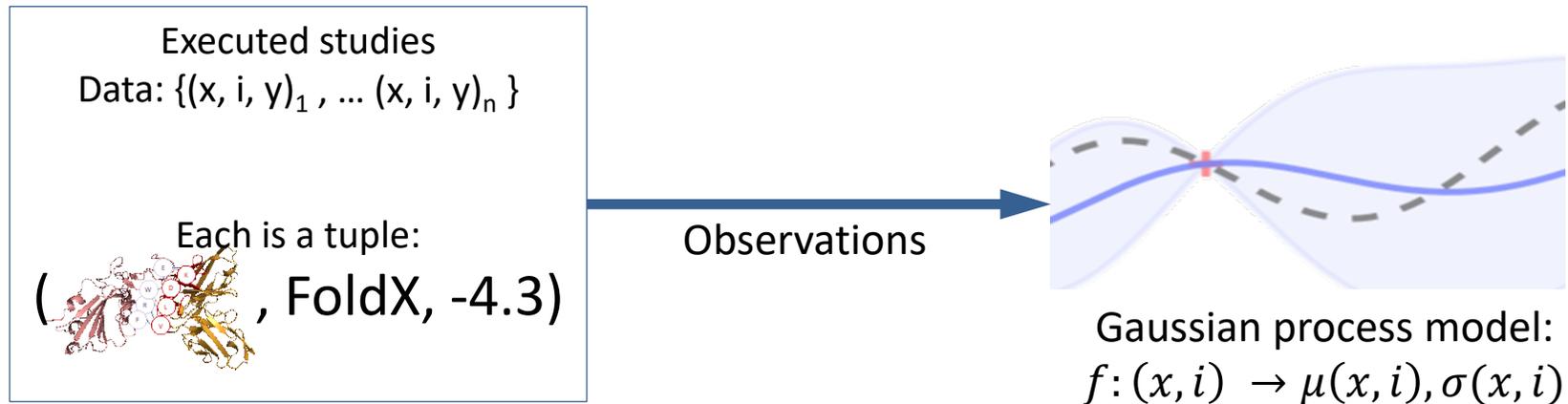
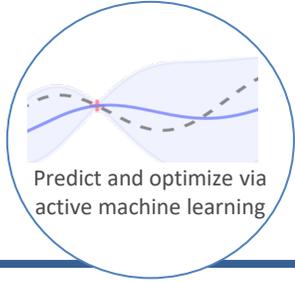
<input type="checkbox"/> C	Mask	C	Predict
<input type="checkbox"/> [START_L1]		[START_L1]	
<input checked="" type="checkbox"/> R	[MASK]	[MASK]	R, F, W, T, S (0.815, 0.074, 0.035, 0.032, 0.019)
<input checked="" type="checkbox"/> S	[MASK]	[MASK]	S, F, A, T, Y (0.958, 0.019, 0.011, 0.005, 0.003)
<input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

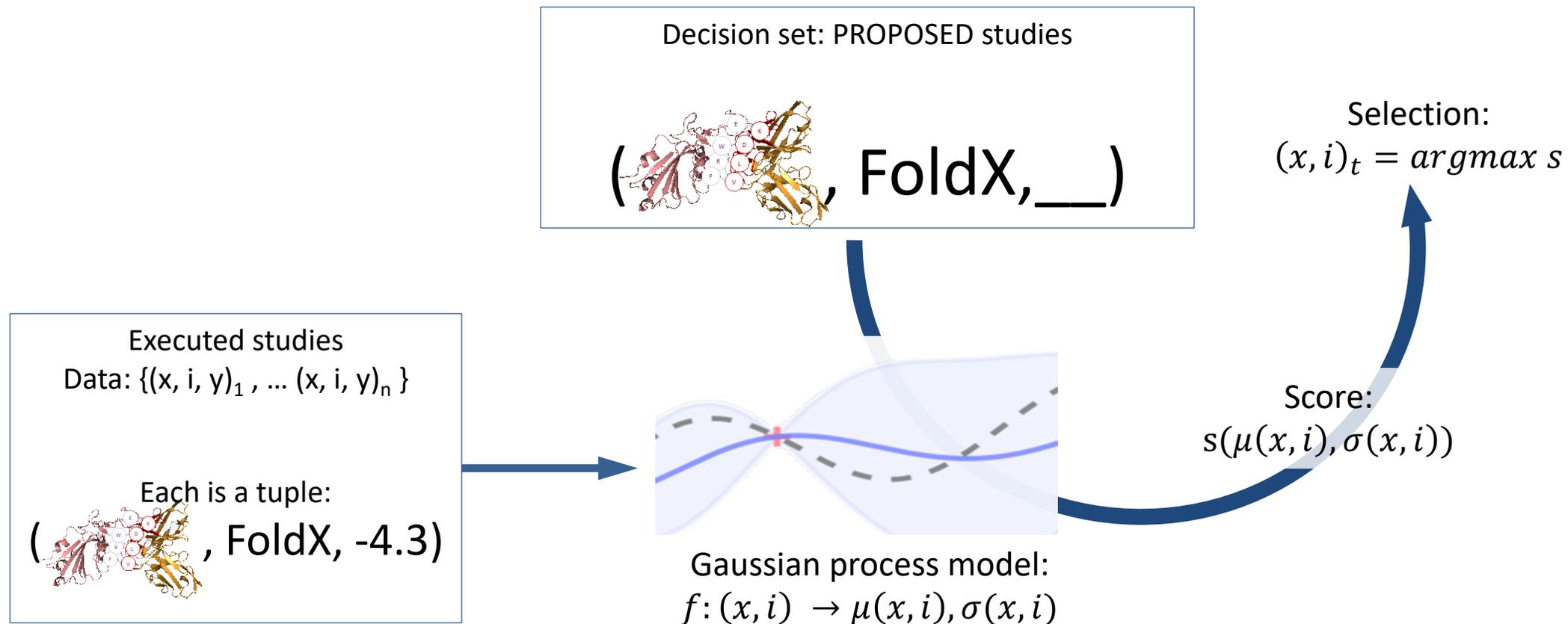
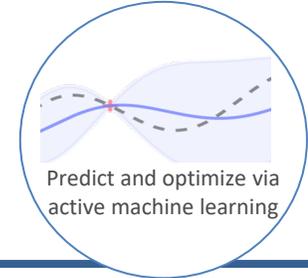


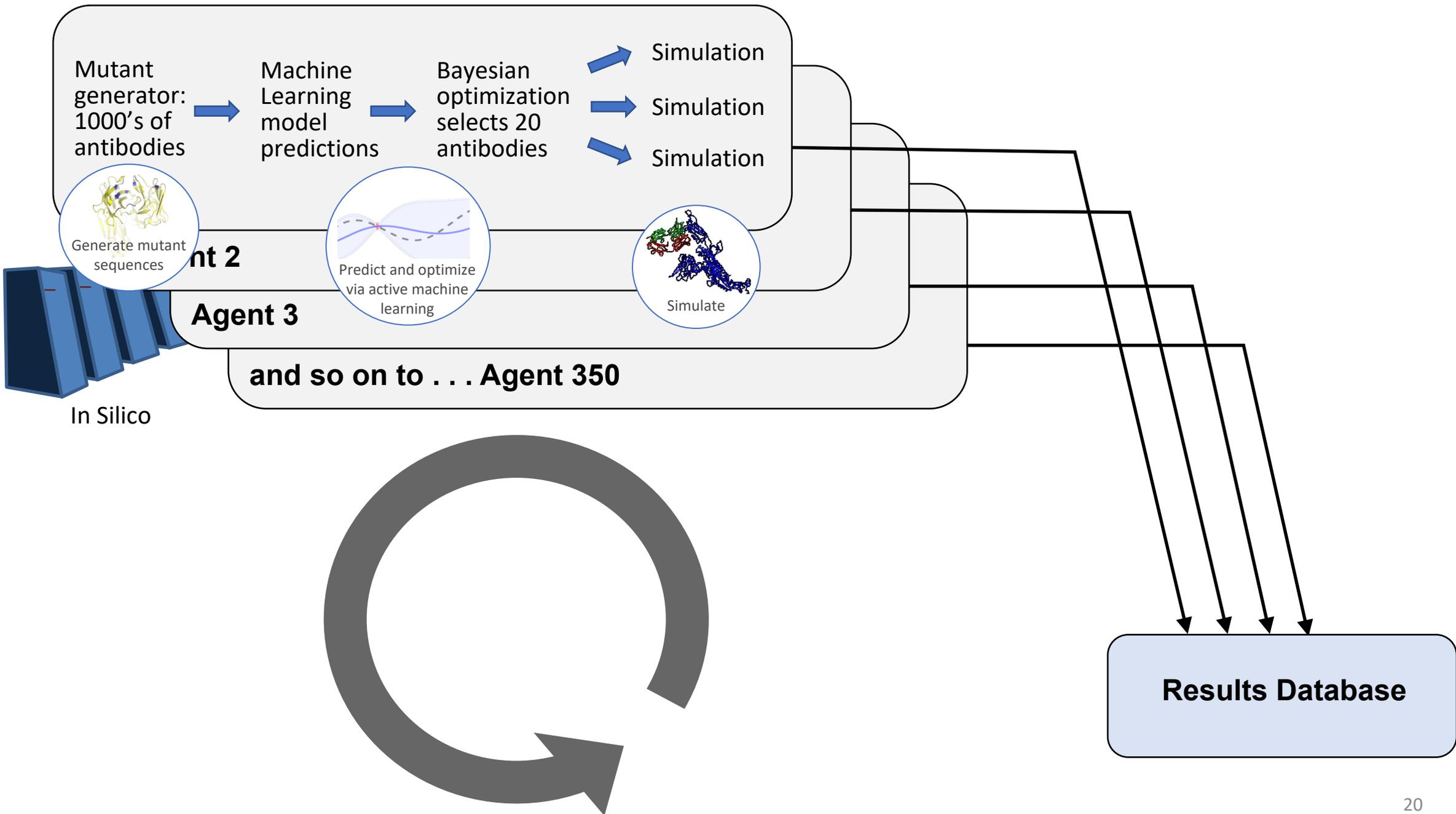
$\mathbf{x} = [0, 1, 0, 2, 0, 0, 1, \dots]$
Vector of interaction type counts

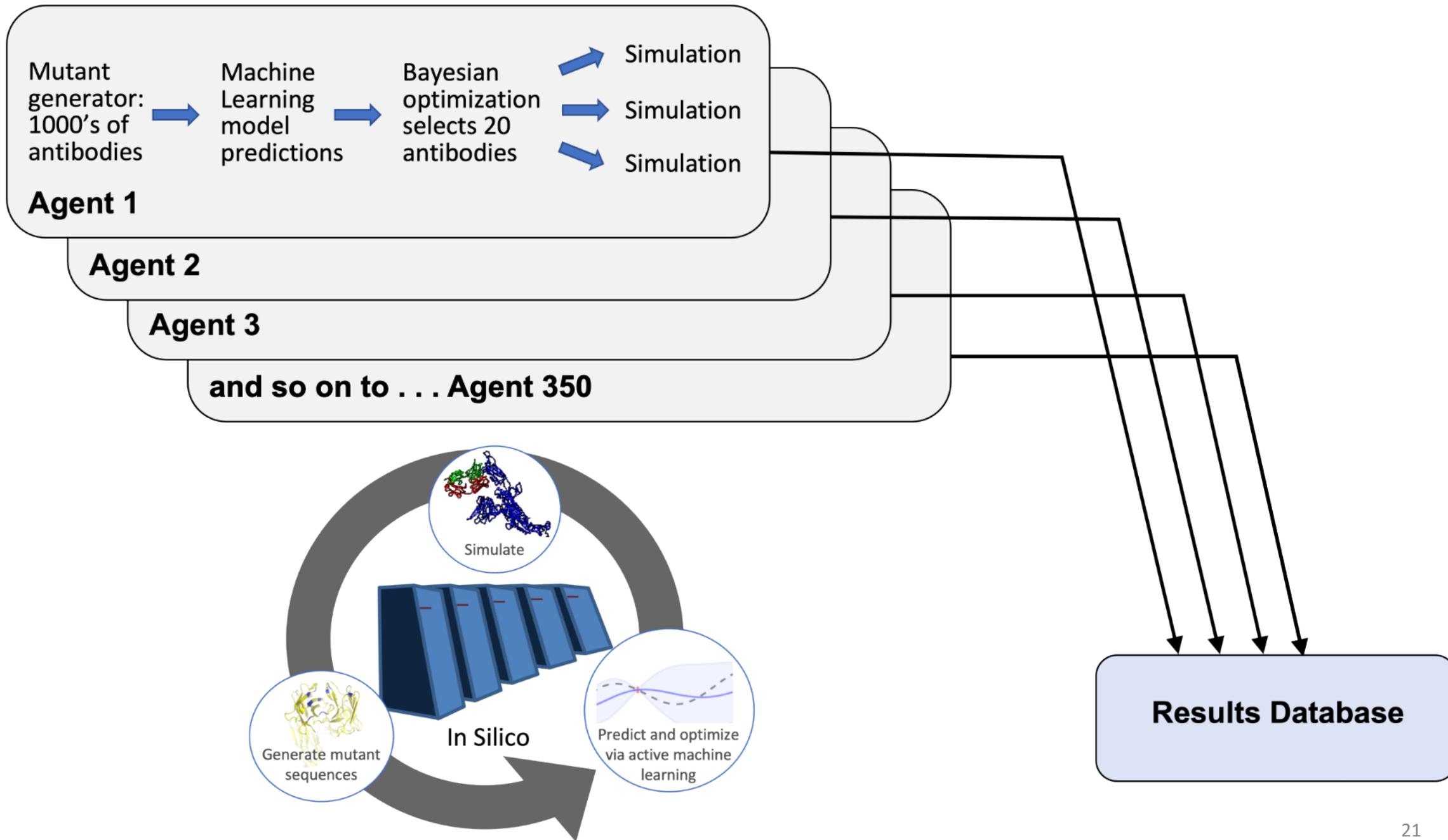
Represented in feature space, binding free energy estimates feed into a multi-fidelity Gaussian process model

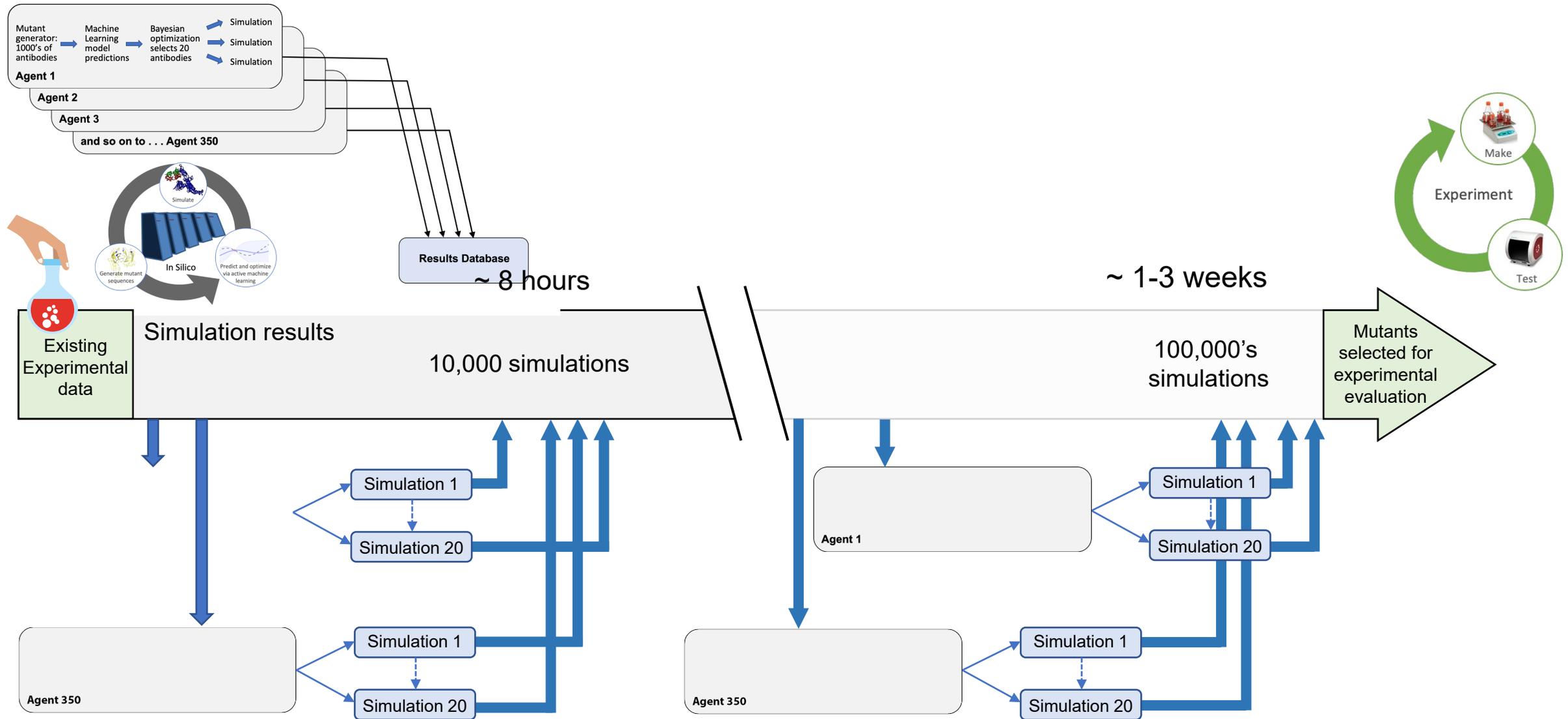


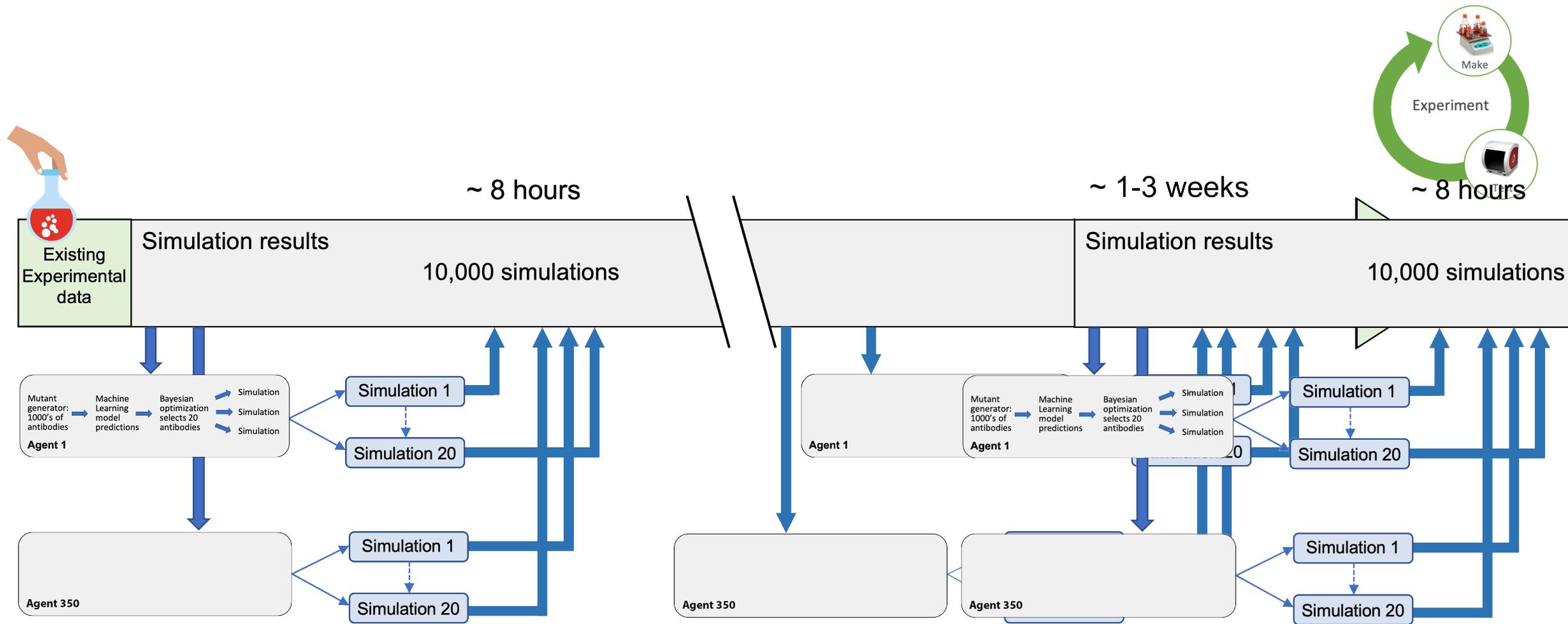
The next set of simulations is selected via Bayesian optimization using the Gaussian process model



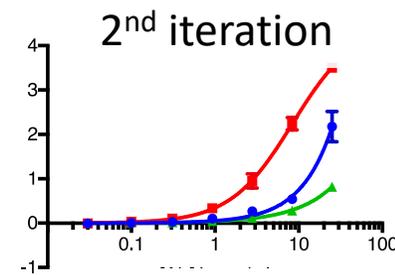
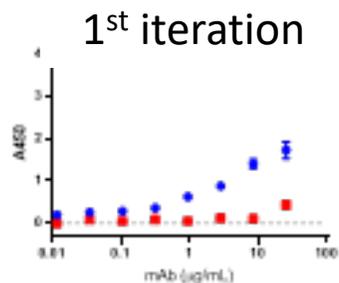
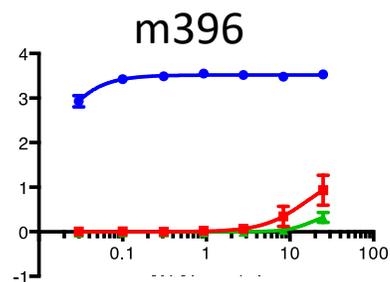




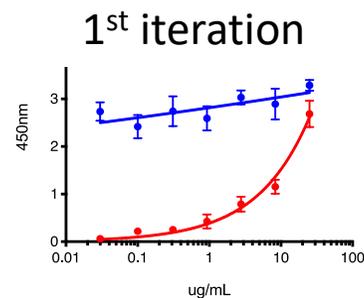
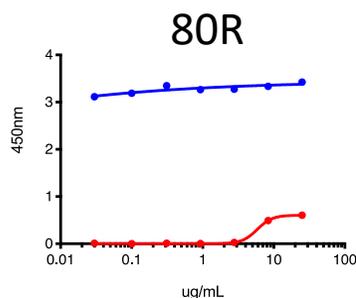




For two systems, we designed binders without having received *any* experimental feedback

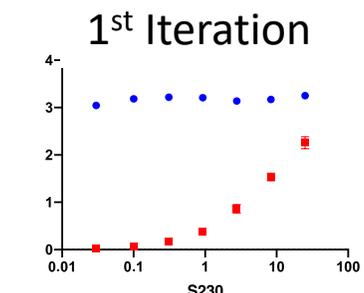
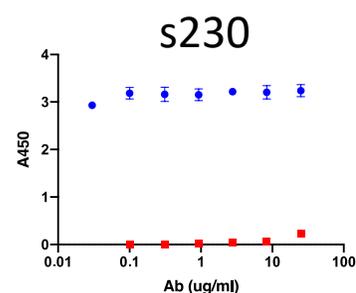


SARS-1
SARS-2
BSA (neg. control)



ELISA:

Binding → more absorbance
Binding signal at lower concentration is better



Each iteration is ~100 designs; select antibodies shown

This work is the product of a multidisciplinary team

- LLNL:

Daniel Faissol, Adam Zemla, Ed Lau, Fangqiang Zhu, John Goforth, Denis Vashchenko, Mary Silva, Rebecca Haluska, Claudio Santiago, Sam Nguyen, Brent Segelke, Feliza Bourguet, Victoria Lao, Monica Borucki, Dina Weilhammer, Jacky Lo, Nicole Collette, and Magdalena Franco (now ThermoFisher)

- Sandia NL:

Brooke Harmon, Oscar Negrete, Max Stefan

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.