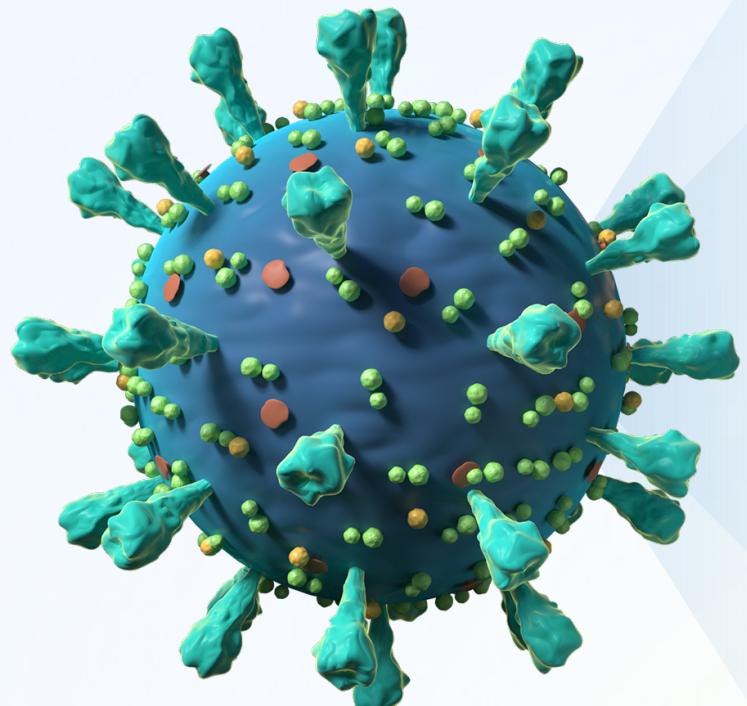




# Running Ahead of Evolution

AI based simulation for predicting future high-risk  
SARS-CoV-2 variants

Peng Cheng Laboratory AI Team



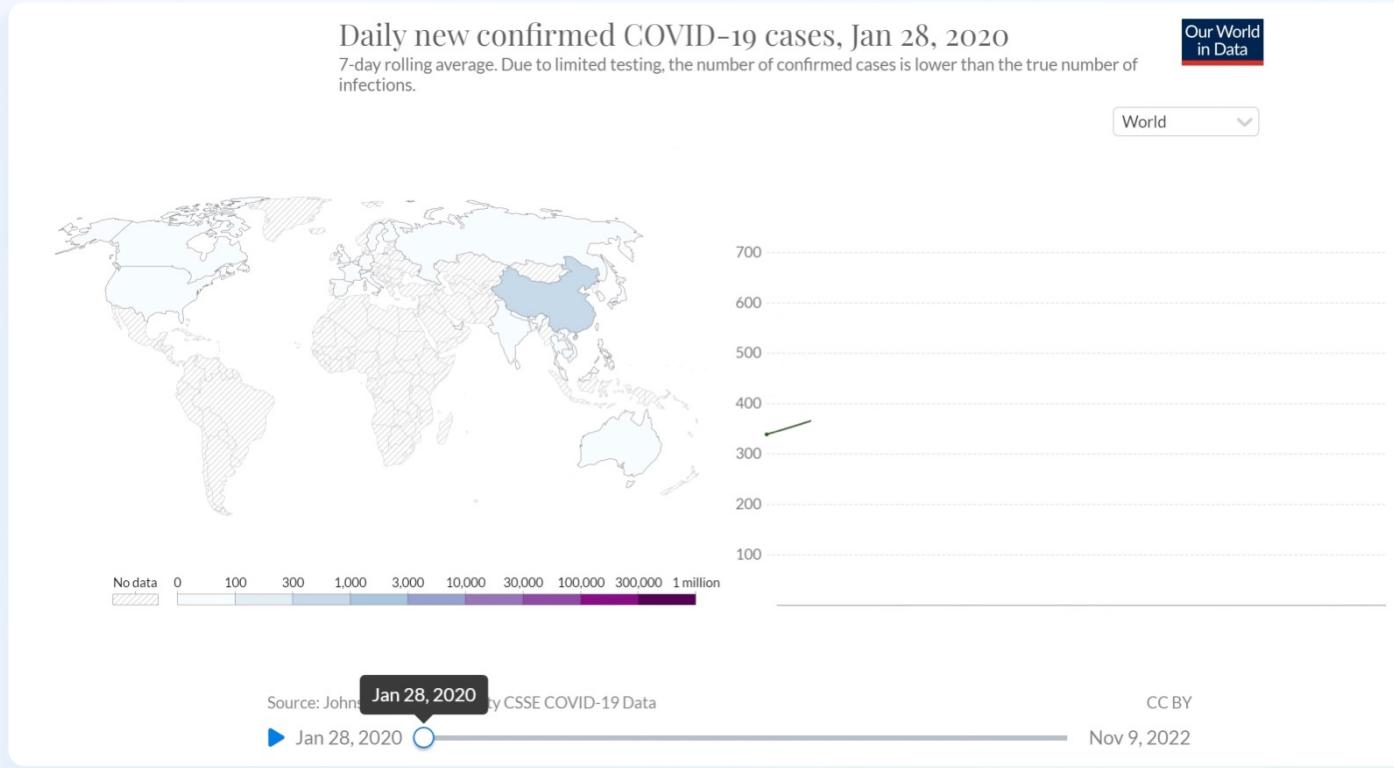
## Background & Highlight

Technical challenges  
& Innovations

Science results  
& Key performance

Summary  
& Prospect

# COVID-19: A world-wide pandemic



Never-ending story

## Timeline of outbreak



# Running ahead of evolution



# SARS-CoV-2 is the most sequenced virus, EVER.



**hCoV-19 data sharing via GISAI**

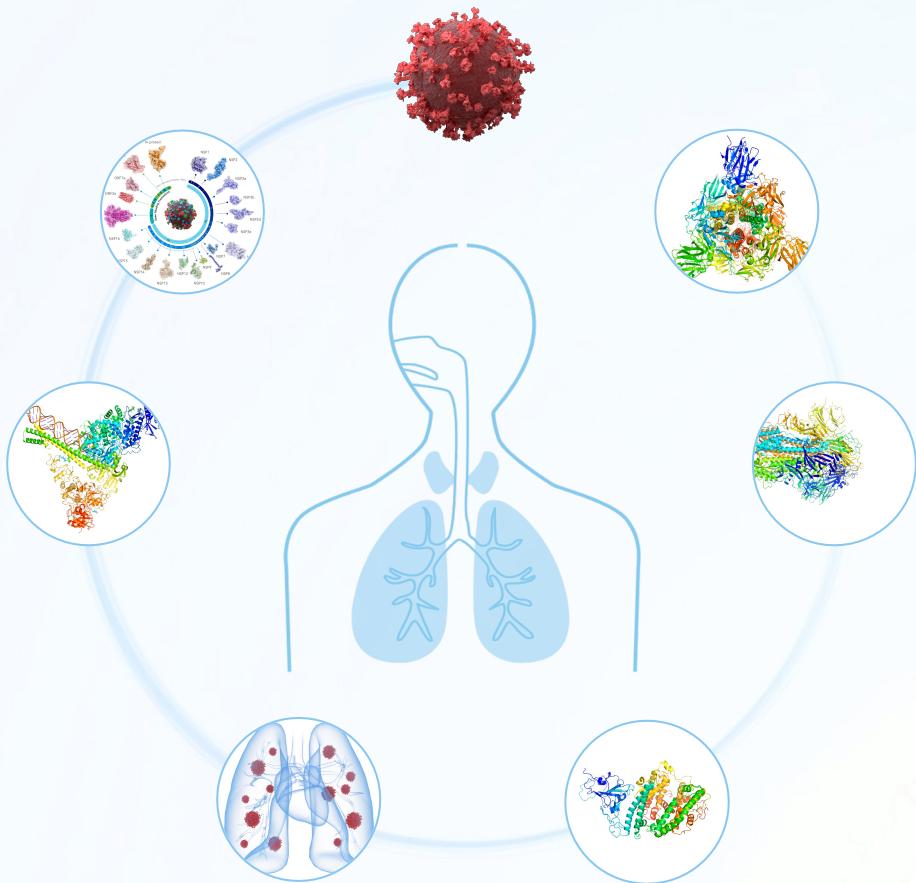
## **13,890,985**

---

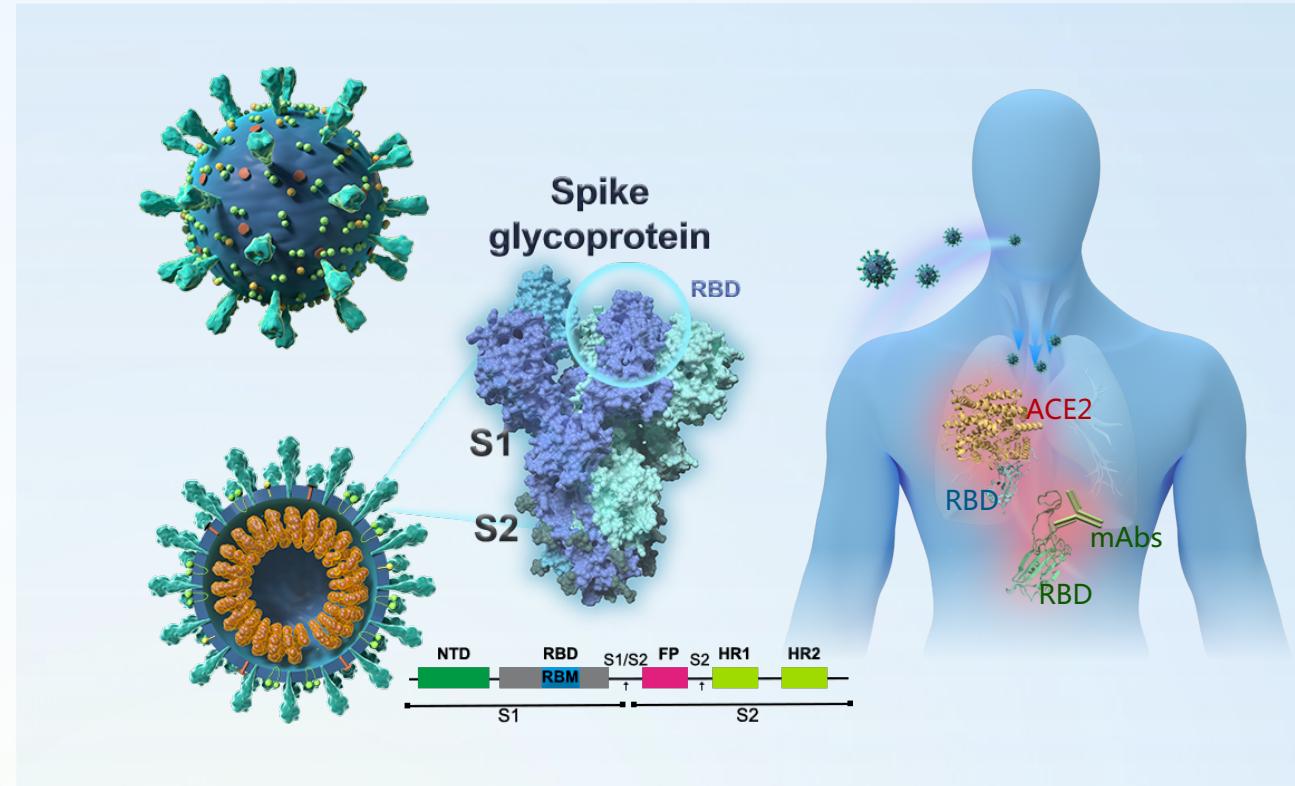
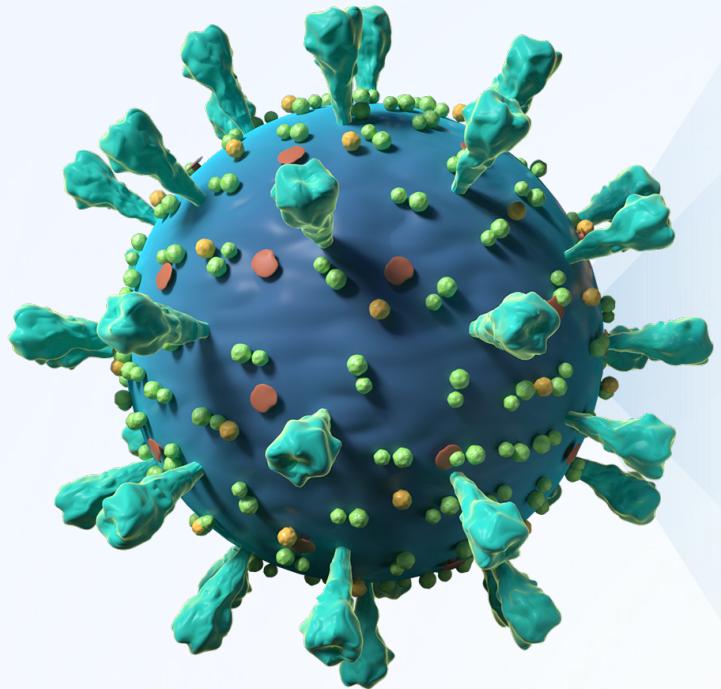
Genome sequence submissions

Can the viral mutational landscape and trends be inferred from these  
13.8 million viral genome sequences?

# The infection process of SARS-CoV-2



# The most mutated region of SARS-CoV-2



The **RBD region** of the Spike protein is **an area of concern** because it has a **high mutation rate**, which can significantly **affect binding to hACE2 and antibodies**.

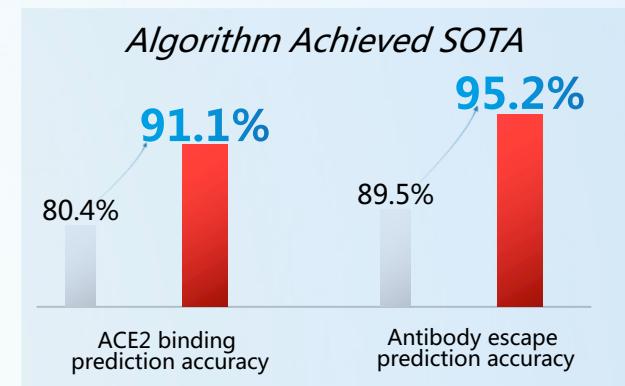
# Highlight: SARS-CoV-2 RBD mutation simulation and variants prediction

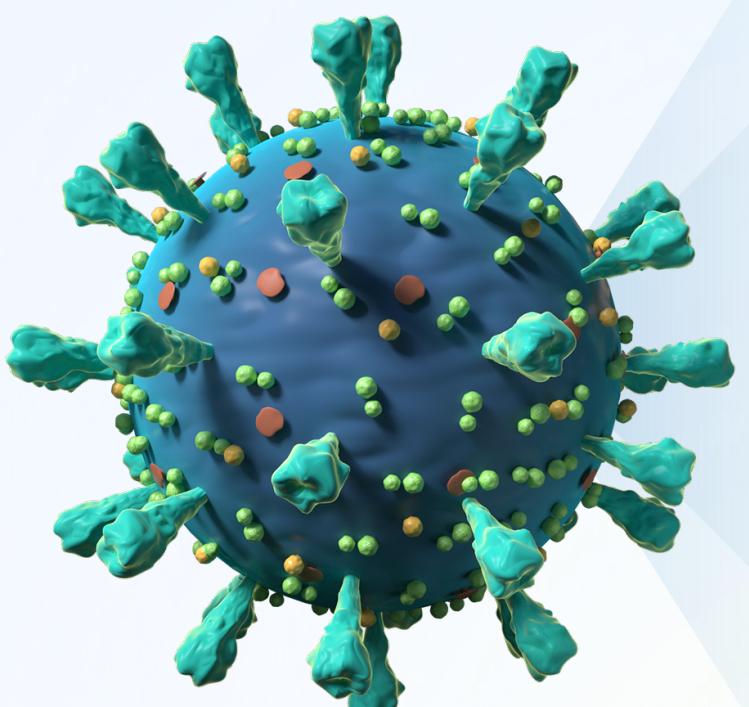
<b><i>The first work of SARS-CoV-2 RBD mutation simulation</i></b>										
Module	Generation				Screening			Feedback		
Details	Massive sequence processing	Pretraining	Fine-tuning	Inference	ACE2 binding	Antibody escape	Transmissibility	Comprehensive analysis	Recursively fine-tuning	
<b>Our work</b>	✓	✓	✓	✓	✓	✓	✓	✓	✓	
Nicola De Maio et al, 2022	✓	X	X	✓	X	X	X	✓	X	
Vladimir Shchur et al, 2022	✓	X	X	✓	X	X	X	✓	X	
Jiahua Chen et al, 2020	X	✓	X	X	✓	X	X	X	X	
Joseph M.Taft et al, 2021	X	✓	X	X	✓	✓	X	✓	X	

Alpha	BF.7
Beta	BQ.1
Gamma	BF.14
Delta	BA.2.75.2
Omicron BA.5	BA.4.6

**100% success rate for the prediction of VOCs, except Omicron**

**Most of VUMs are predicted accurately**





Background  
& Highlight

**Technical challenges  
& Innovations**

Science results  
& Key performance

Summary  
& Prospect

# High performance prediction of high-dimensional mutational space

Number of Mutations = (amino acid types)<sup>sequence length</sup> =  $20^{201} = 3.213876 \times 10^{261}$

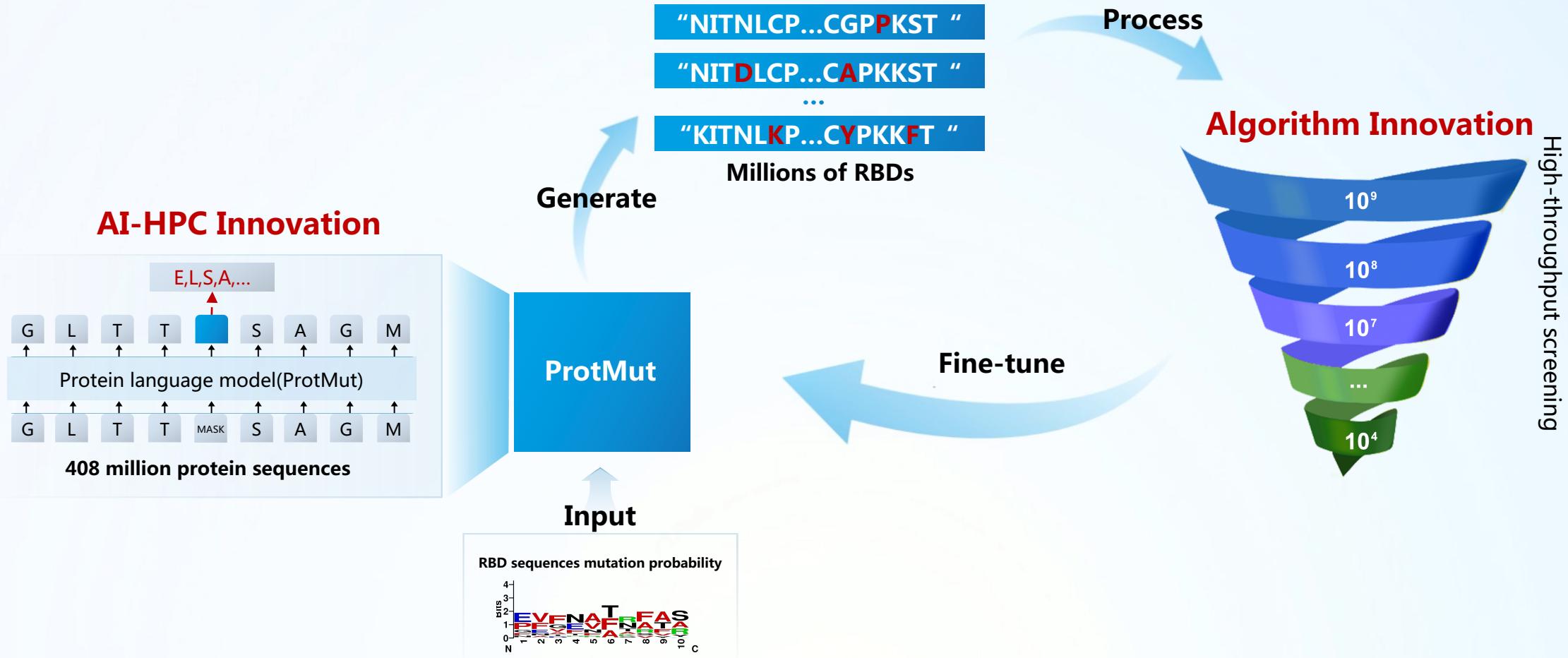
The length of RBD domain is **201** (331-531)

There are **20** amino acids

**We need high-throughput variants generation and screen!**

# Innovation 1: mutation simulation

Simulate the natural process



# Innovation 1: mutation simulation

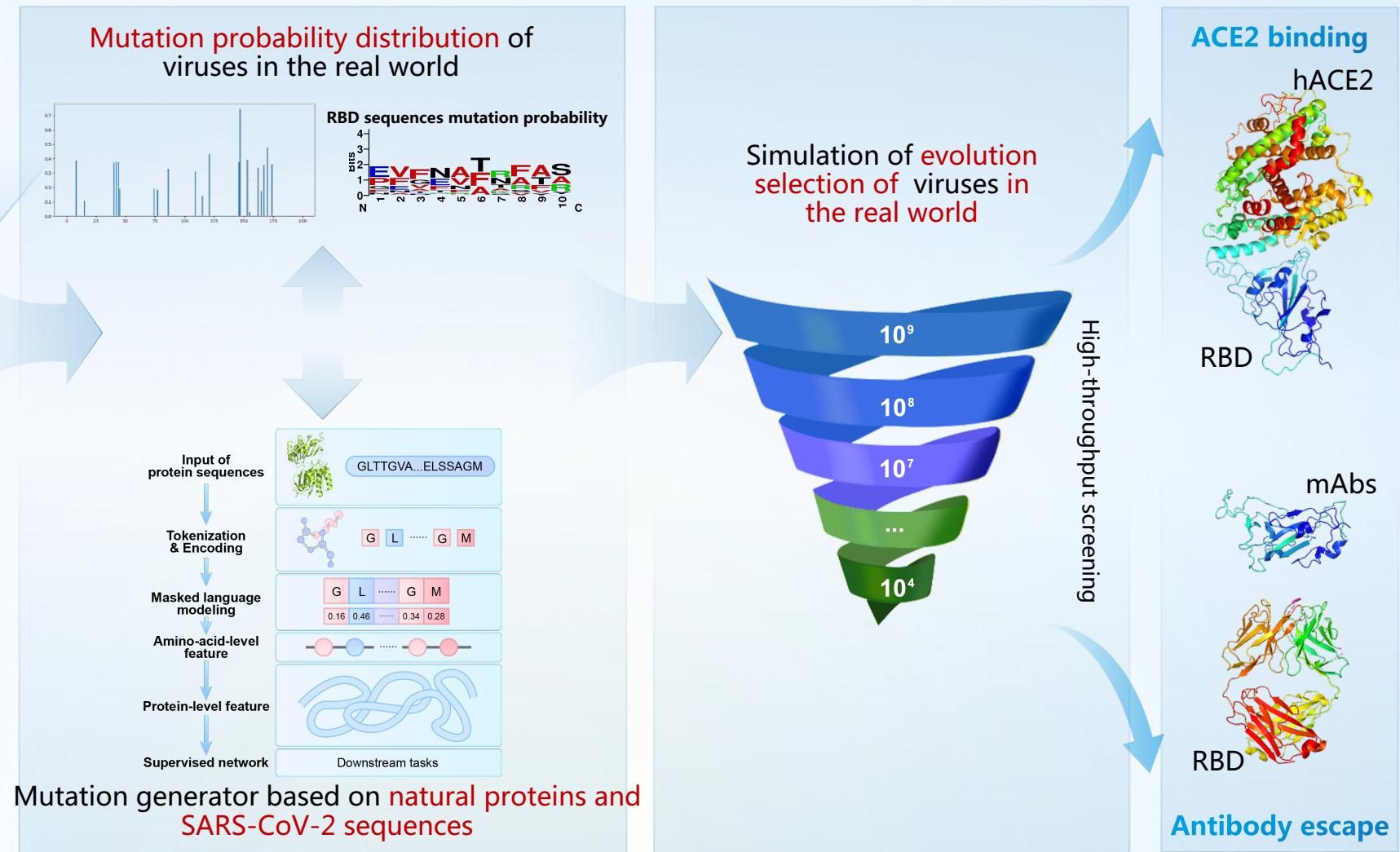
Reduction of  
mutational space

Learning of virus RBD  
sequence

Screening

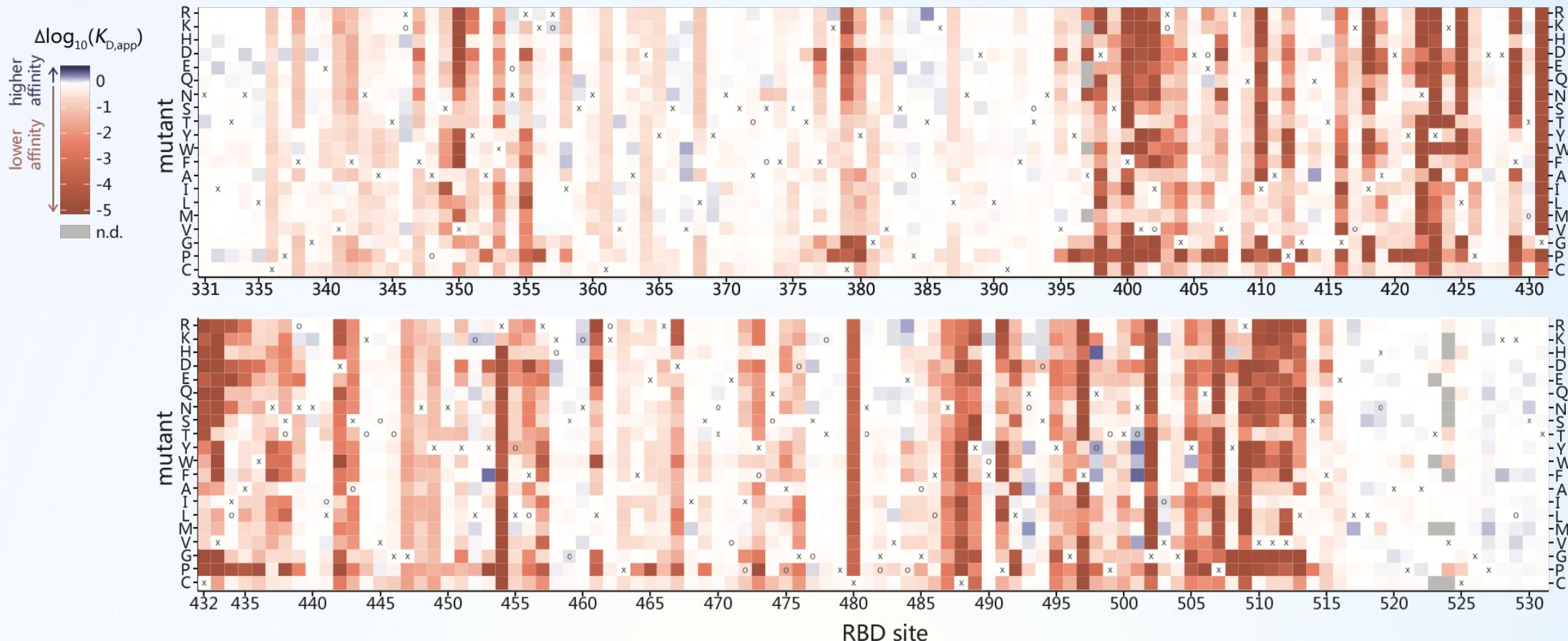
Recursively  
fine-tuning

...  
reduced  
mutational  
space



# Innovation 2: science-driven prediction

## Mutation effects on binding

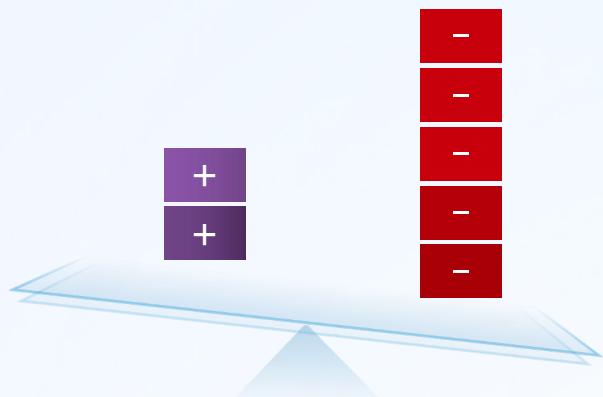


**Deep mutational scanning**

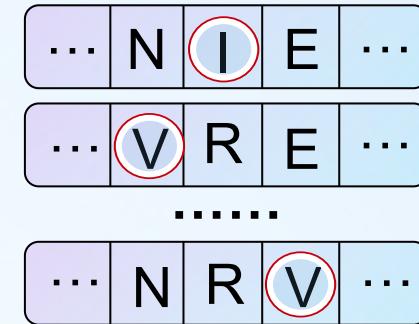
Starr T N, Greaney A J, Hilton S K, et al. Cell, 2020, 182(5): 1295-1310. e20.

# Innovation 2: science-driven prediction

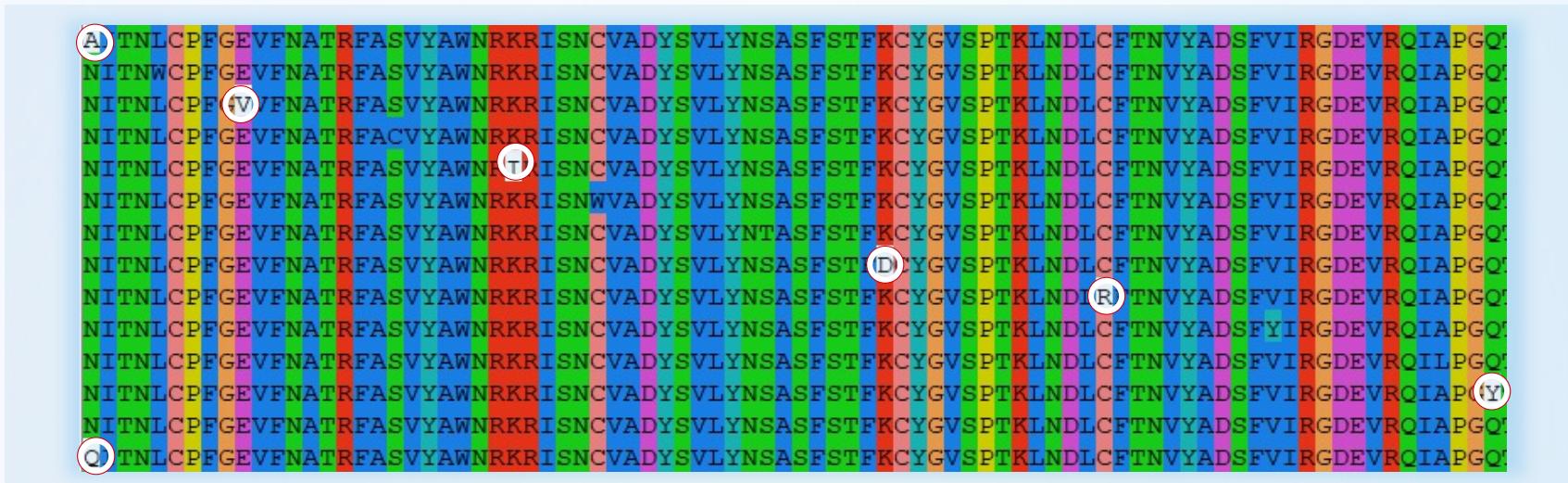
Data characteristics 1:  
serious category imbalance



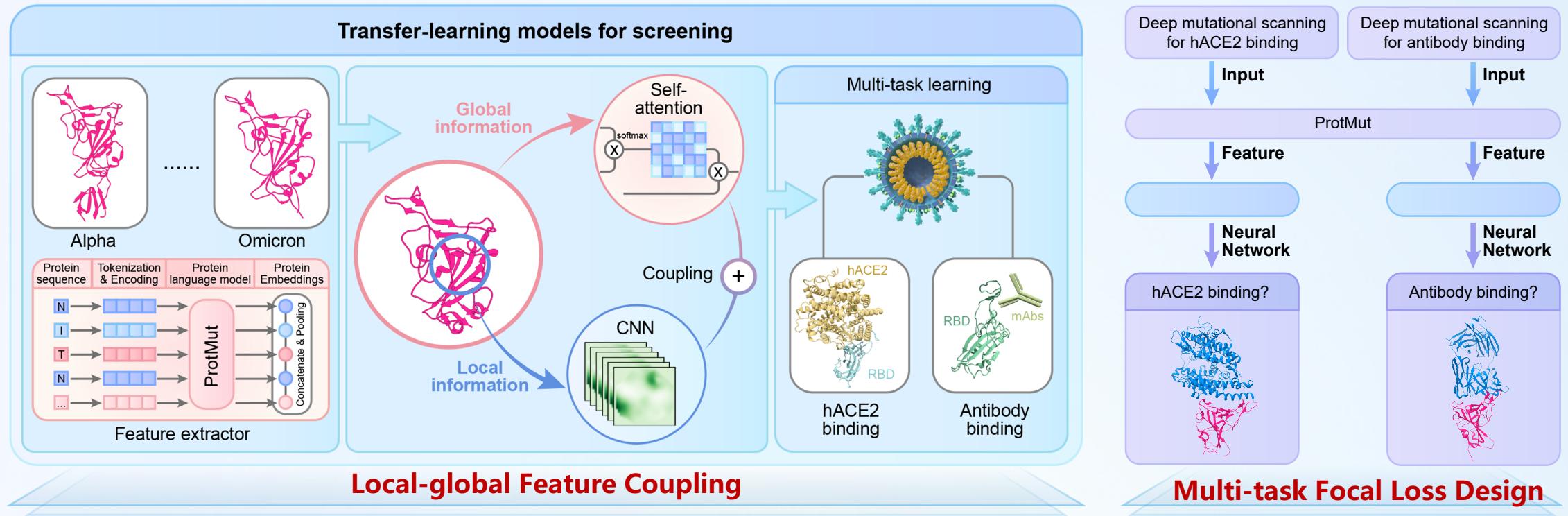
Data characteristics 2:  
very similar between sequences



## Sequences of DMS

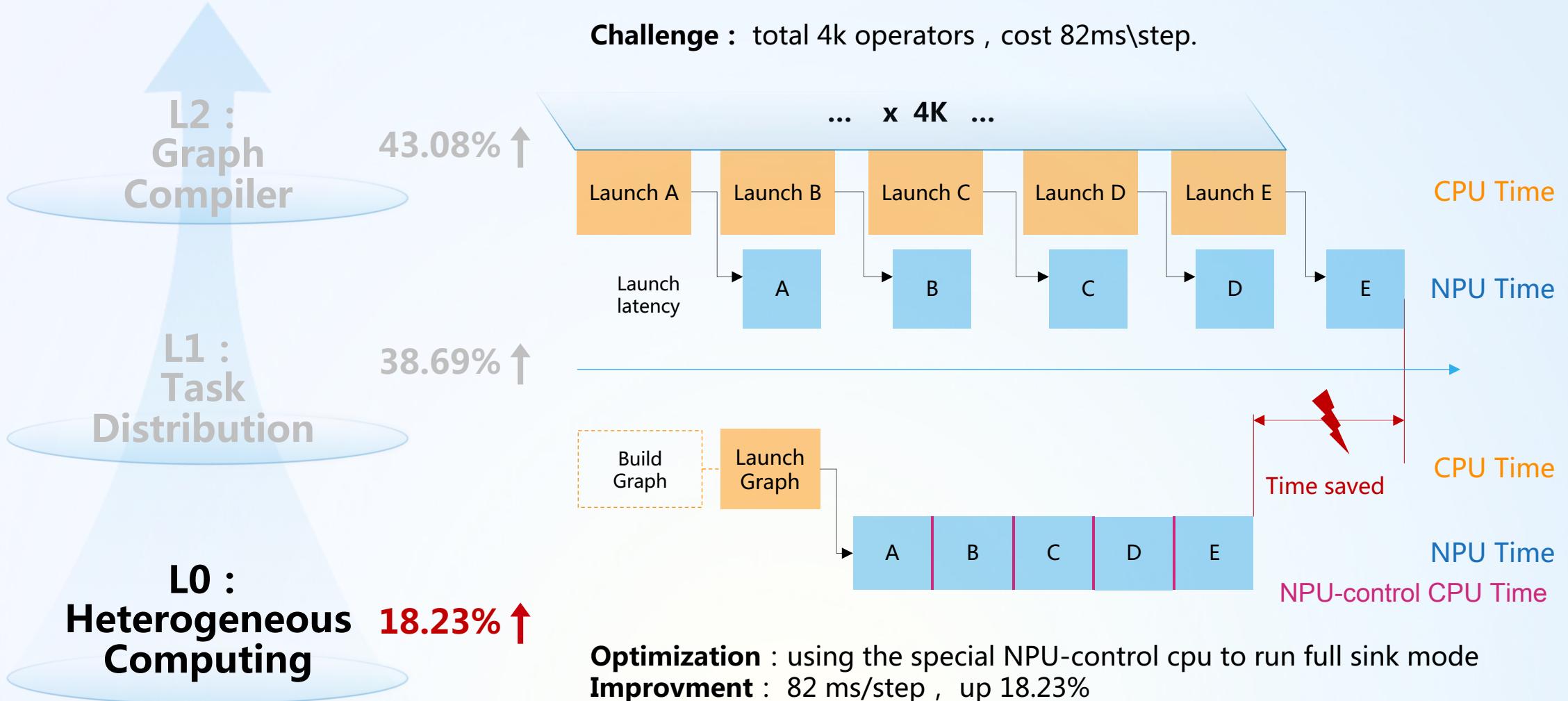


# Innovation 2: science-driven prediction

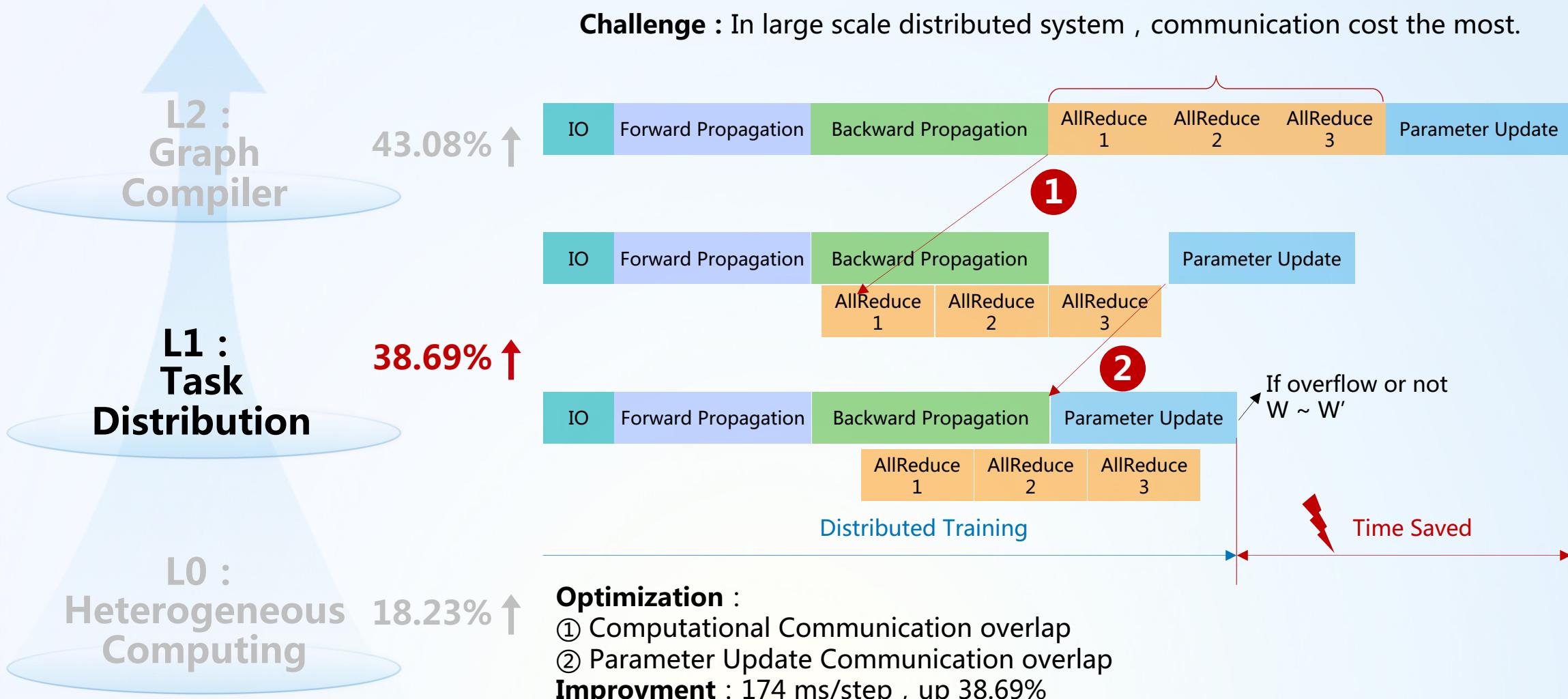


Variants	RBD-31	RBD-32	RBD-33	RBD-34	RBD-312	RBD-36	.....	RBD-71
<b>Consistent with wet-lab experiments?</b>	True	True	True	True	True	True	True	True

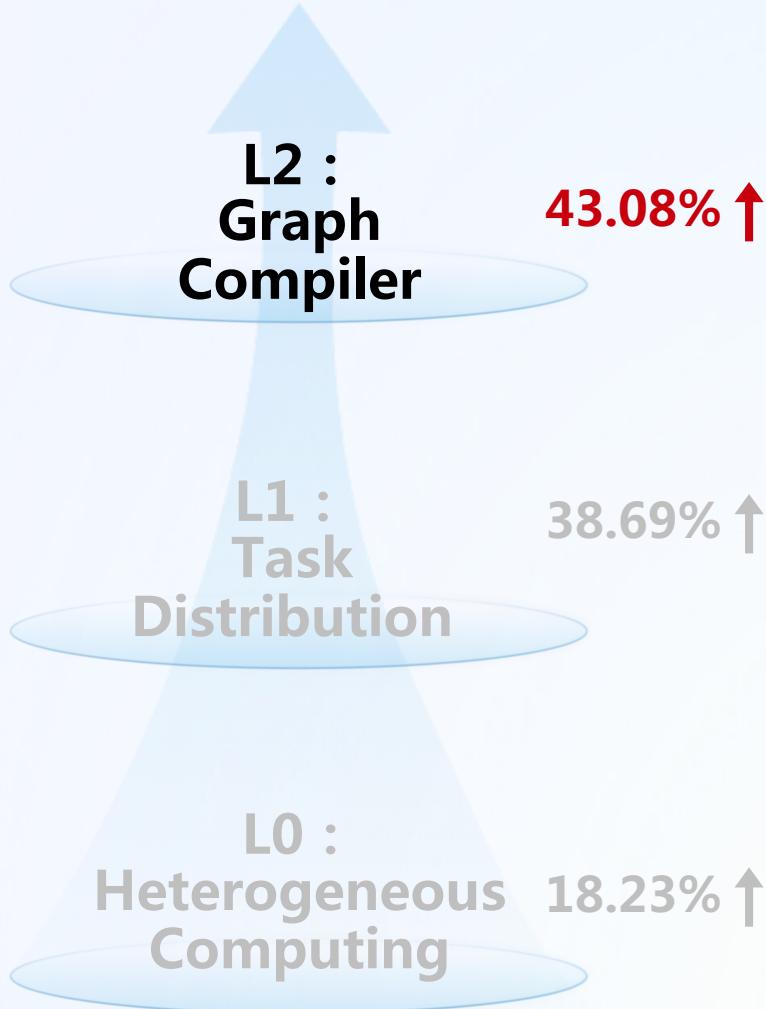
# Innovation 3: multi-level optimization strategy



# Innovation 3: multi-level optimization strategy

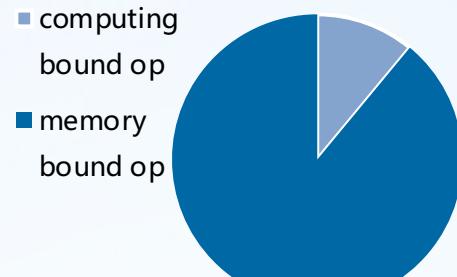


# Innovation 3: multi-level optimization strategy

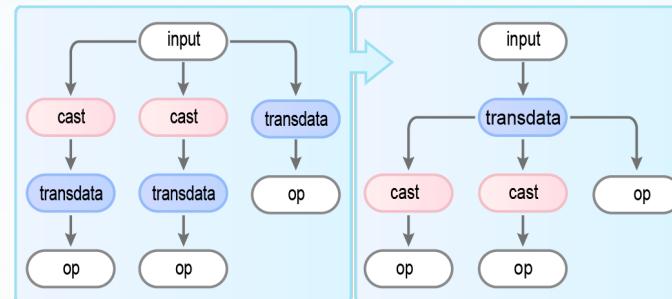


**Challenge :** large nums of memory bound operators lead to extra time overhead

Cube task OP nums distribution



Operator fusion



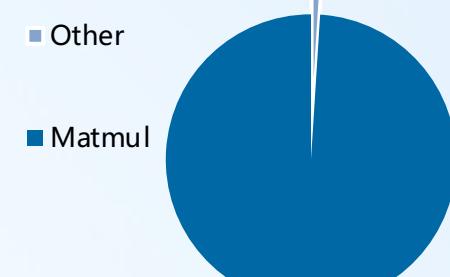
**Optimization :**

- layerNorm fusion
- matmul +add
- allreduce+gradients

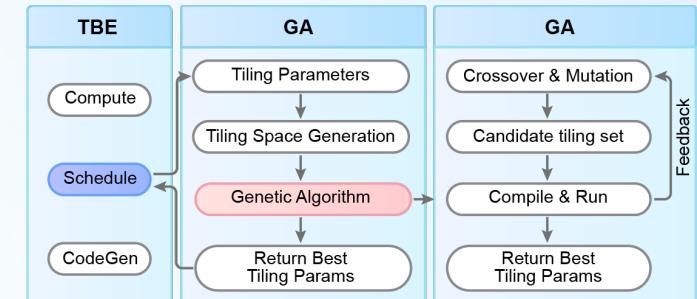
**Improvement :** up 30%

**Challenge :** Matmul computation >99% and time consumption>40%

Matmul ratio of computation

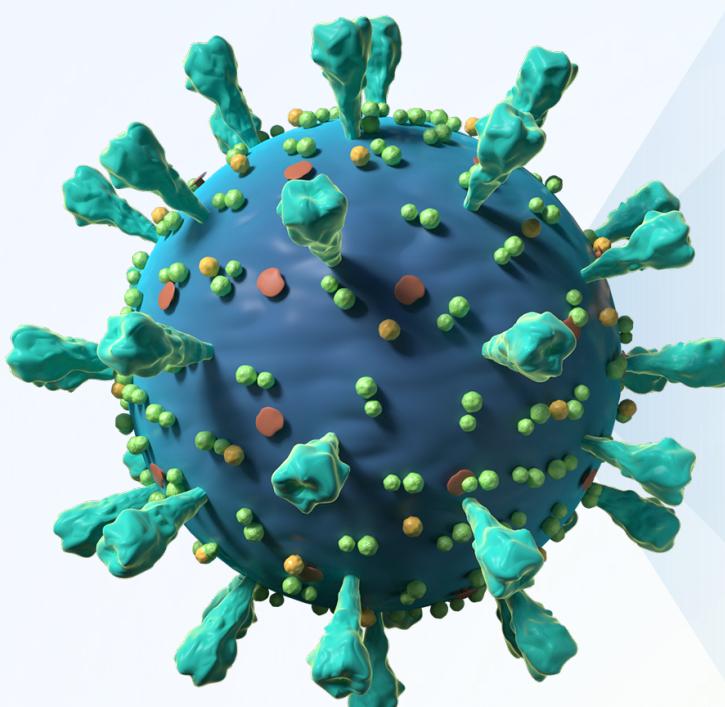


Operator auto-tuning



**Optimization :** auto-tuning , using GA\RL algorithm

**Improvement :** up 13%



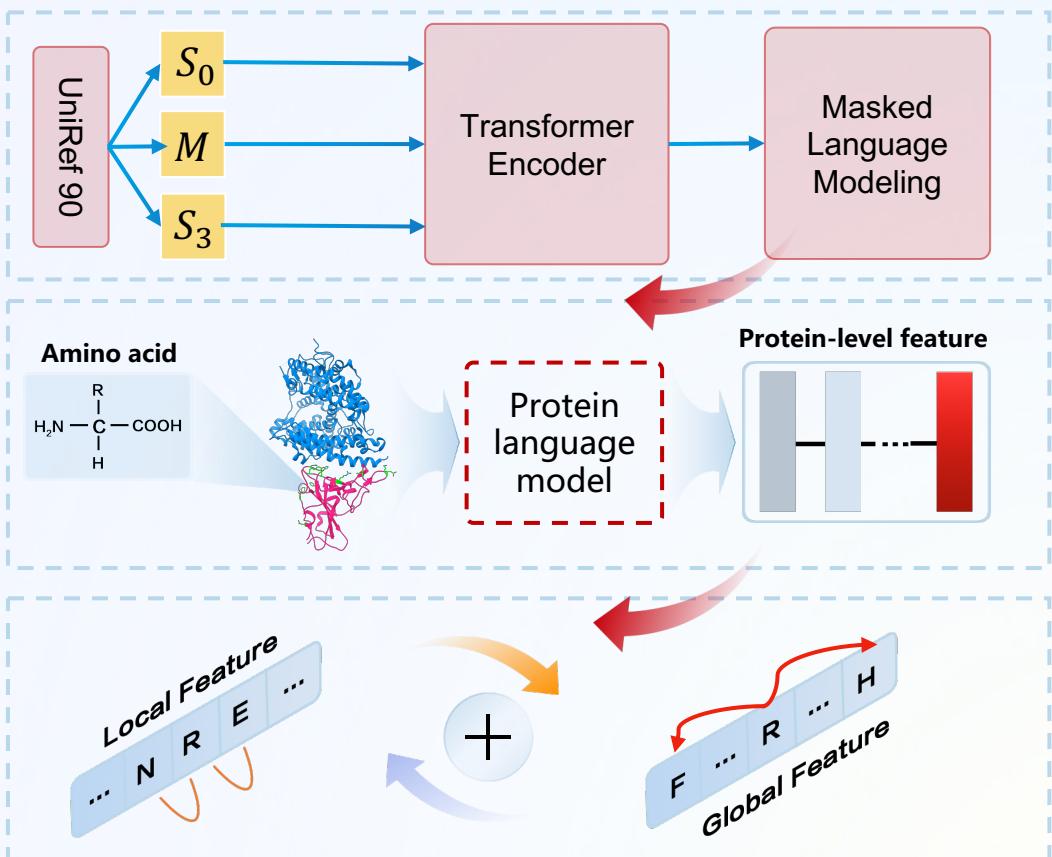
Background  
& Highlight

Technical challenges  
& Innovations

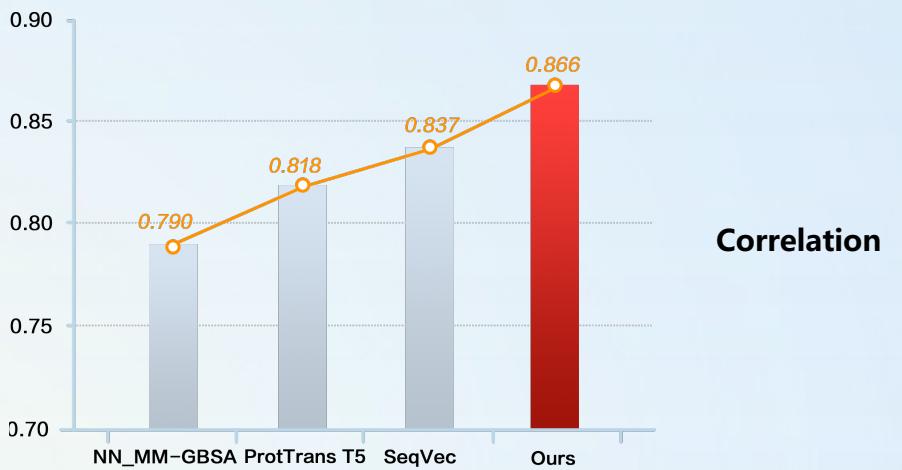
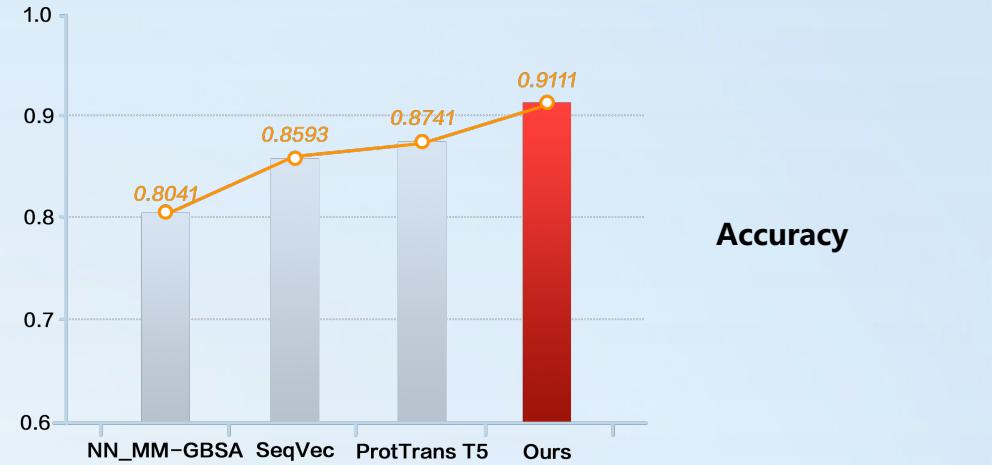
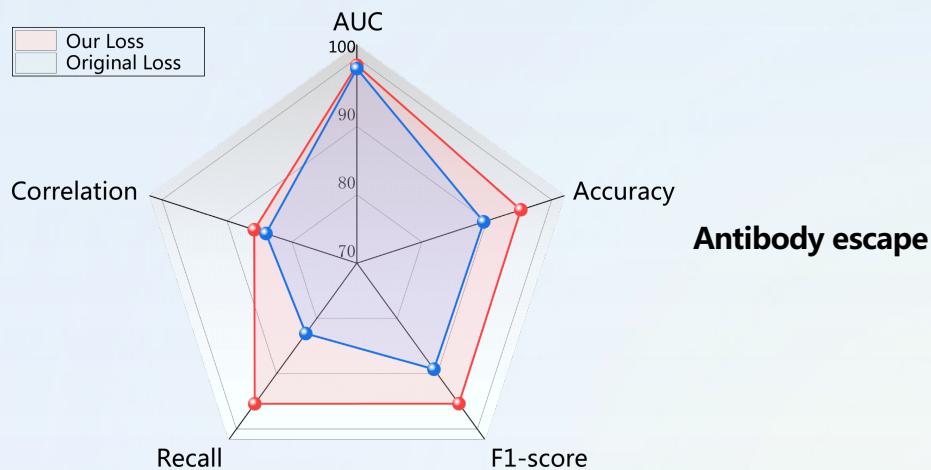
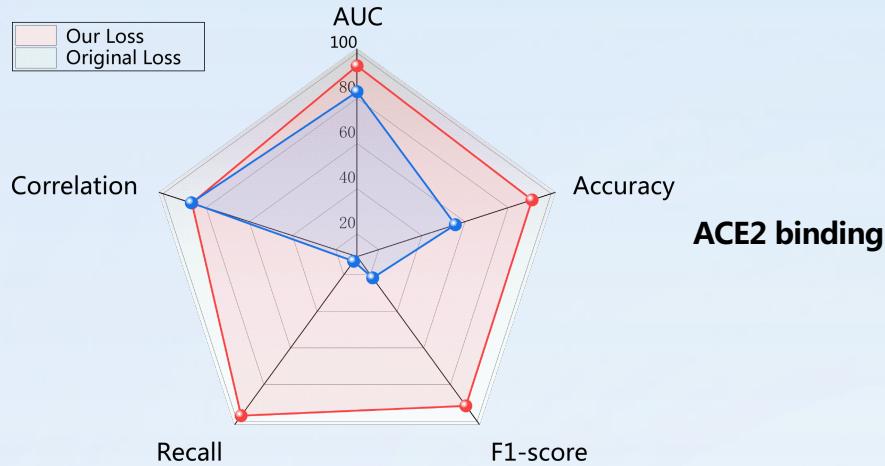
**Science results  
& Key performance**

Summary  
& Prospect

# Reveal the latent virus RBD sequence information

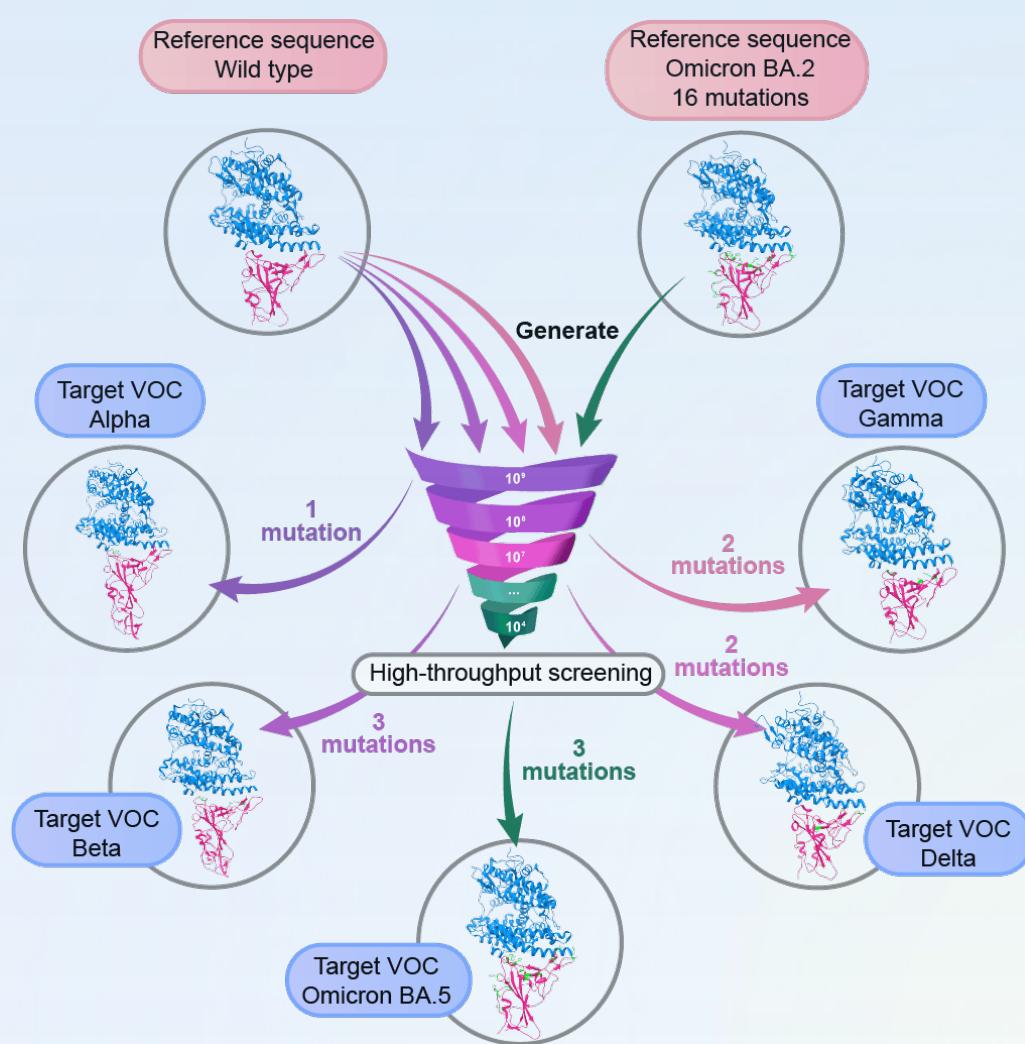


# Algorithm achieved SOTA



# Prediction of VOCs and VUMs

## Two stage validation of VOCs



## Validation of VUMs



# High throughput generating & screening

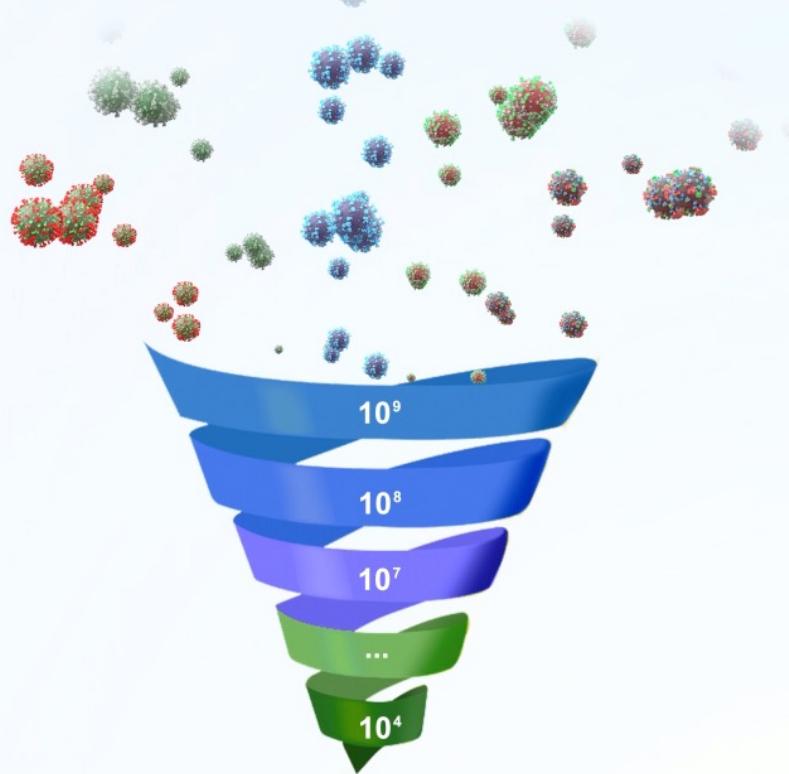


Table1:High throughput screening of various variants

Variant	1st screening*	2st screening**
Alpha	39.8%	2.0%
Beta	13.3%	51.3%
Gamma	45.2%	33.8%
Delta	46.7%	19.1%
Omicron BA.5	90.4%	80.2%

\*Proportion after hACE2 binding screening

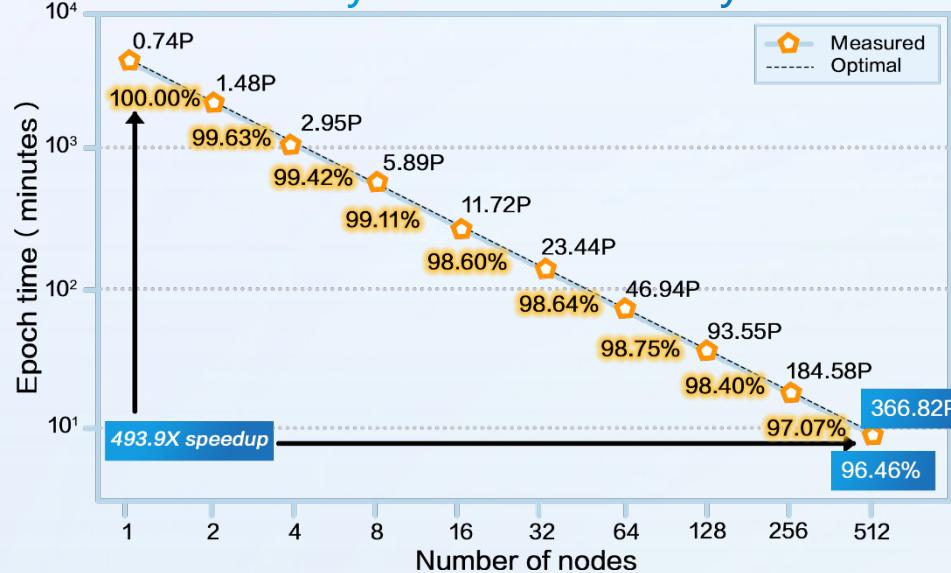
\*\*Proportion after antibody binding screening

## Results on full-scale system

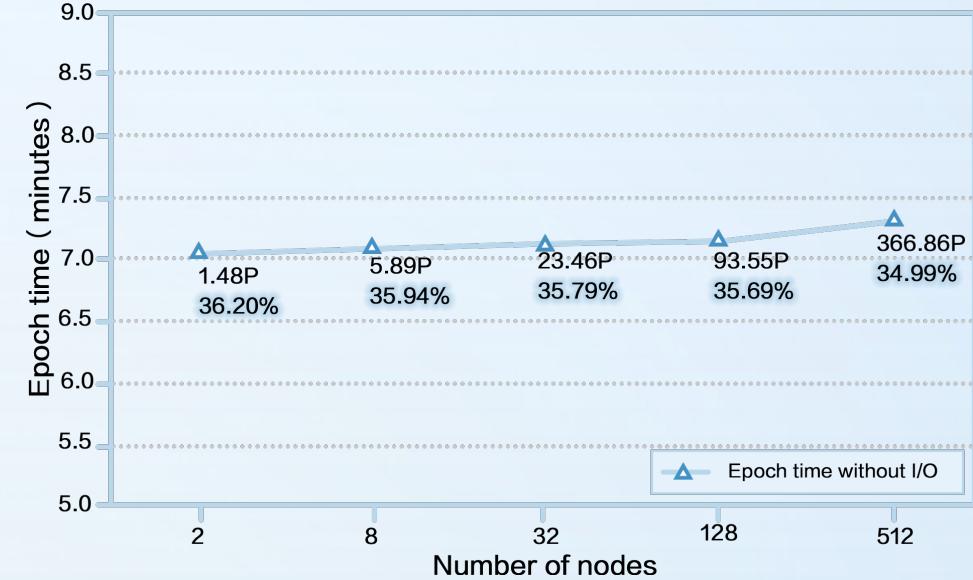
VOCs	Number of generation	Time of generation & screening (h)
Alpha	$10^6$	0.00024
Beta	$10^{11}$	24
Gamma	$10^{10}$	2.4
Delta	$10^8$	0.024
Omicron BA.5	$10^9$	0.24

# Strong scaling & Weak scaling

Nearly Perfect Scalability



Stable Utilization



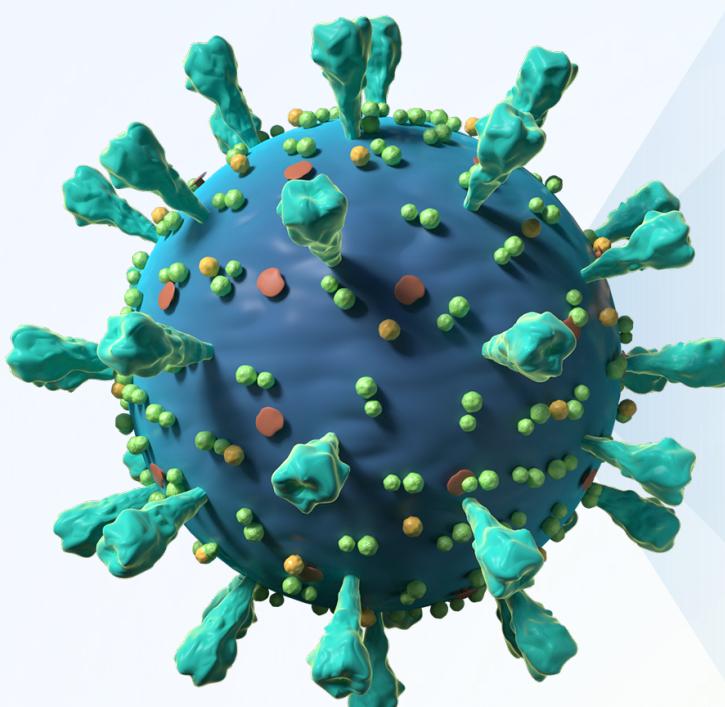
- 4096 Neural Processing Units (NPUs) and 2048 Kunpeng 920 Processors, provides 1 EFLOPS.
- 512 Nodes , each node has 4 Processors and 8 NPUs.
- Each NPU with 32GB memory provides 256 TFLOPS.
- Every Processor is connected via 100G RoCE.



## Key performance

- **63.4 zettaflops\***
- **366.8 PFLOPS** on 4096 NPUs
- **96.5% scalability**
- **493.9X speedup**
- **34.9% utilization (training)**

\* Measured in mixed-precision



## Background & Highlight

Technical challenges  
& Innovations

Science results  
& Key performance

**Summary  
& Prospect**

# 63.4 zettaflops in exchange for almost 5 months

Scientific Results

HPC Performance

**First simulation of SARS-CoV-2 RBD**

**Achieving over 19.8 zettaflops**

**Algorithm achieved SOTA**

**366.8 PFLOPS on 4096 NPUs**

**Almost all VOCs are predicted accurately**

**96.5% scalability**

**Most of VUMs are predicted accurately**

**493.9X speedup**

**The period between outbreaks  $\approx$  5 months**

***In silico* simulation time on full-scale system  $\approx$  2 days**