



Biophysical principles predict fitness of SARS-CoV-2 variants

Dianzhuo Wang^{a,b}, Marian Huot^{a,c}, Vaibhav Mohanty^{a,d,e}, and Eugene I. Shakhnovich^{a,1}

Edited by Andrej Sali, University of California San Francisco, San Francisco, CA; received August 22, 2023; accepted April 19, 2024

SARS-CoV-2 employs its spike protein's receptor binding domain (RBD) to enter host cells. The RBD is constantly subjected to immune responses, while requiring efficient binding to host cell receptors for successful infection. However, our understanding of how RBD's biophysical properties contribute to SARS-CoV-2's epidemiological fitness remains largely incomplete. Through a comprehensive approach, comprising large-scale sequence analysis of SARS-CoV-2 variants and the identification of a fitness function based on binding thermodynamics, we unravel the relationship between the biophysical properties of RBD variants and their contribution to viral fitness. We developed a biophysical model that uses statistical mechanics to map the molecular phenotype space, characterized by dissociation constants of RBD to ACE2, LY-CoV016, LY-CoV555, REGN10987, and S309, onto an epistatic fitness landscape. We validate our findings through experimentally measured and machine learning (ML) estimated binding affinities, coupled with infectivity data derived from population-level sequencing. Our analysis reveals that this model effectively predicts the fitness of novel RBD variants and can account for the epistatic interactions among mutations, including explaining the later reversal of Q493R. Our study sheds light on the impact of specific mutations on viral fitness and delivers a tool for predicting the future epidemiological trajectory of previously unseen or emerging low-frequency variants. These insights offer not only greater understanding of viral evolution but also potentially aid in guiding public health decisions in the battle against COVID-19 and future pandemics.

viral evolution | receptor binding domain | fitness landscape | SARS-CoV-2 | antibody

Since its emergence, the SARS-CoV-2 virus has undergone continuous genetic changes, giving rise to variants with increased transmissibility such as Alpha, Delta, and the recent Omicron. Each has contributed to significant surges in global COVID-19 cases. These genetic alterations in the viral genome have a profound impact on the structure and function of viral proteins, causing consequential changes in viral fitness (defined as the capacity of the virus to infect). Variants of concern (VoCs), such as Omicron BA.1, possess specific mutations in the spike protein. They have been linked to enhanced transmissibility (1, 2), augmented binding to host cell receptors, and heightened resistance to antibody neutralization (3, 4). Understanding the relationship between these mutations and viral fitness requires investigating their influence on molecular properties of affected proteins. Key viral proteins, like the receptor binding domain (RBD) of the spike protein, play a critical role in facilitating viral entry into host cells by binding to angiotensin-converting enzyme 2 (ACE2) (5), a functional receptor on cell surfaces. Furthermore, RBD serves as primary targets for the most potent SARS-CoV-2-neutralizing antibodies (6) and is subject to evolutionary pressure from the human immune system. Therefore, mutations on the RBD have been shown to be highly correlated with increases of fitness.

On the experimental side, Starr et al. (7) systematically scanned through every amino acid substitution in the isolated RBD to determine the mutation effect on RBD folding and ACE2 binding and showed a substantial number of mutations are well tolerated or could even enhance ACE2 binding. In more recent research, Moulana et al. (8, 9) conducted a thorough examination of the binding affinity across all combinations of the 15 RBD mutations found in the BA.1 variant of SARS-CoV-2 in comparison to the original Wuhan Hu-1 strain. This exploration covered a total of 32,768 genotypes and involved testing against four monoclonal antibodies (LY-CoV016, LY-CoV555, REGN10987, and S309) as well as the ACE2 receptor. Additionally, global initiatives that promote data sharing, such as the Global Initiative on Sharing All Influenza Data (GISAID) (10), provide us with the ability to derive viral fitness based on prevalence data.

Prior studies have derived a quantitative correlation between fitness and molecular properties. For instance, Cheron et al. (11) and Rotem et al. (12) devised a theoretical framework to assess fitness of RNA viruses and validated it using experimental and

Significance

This research presents a biophysical model that maps the molecular properties of SARS-CoV-2's receptor binding domain into an epistatic fitness landscape. By linking the binding affinities of the virus to its epidemic fitness, we offer a powerful tool for understanding and predicting the emergence and success of new viral variants. Our model, validated with real-world data and informed by theoretical insights, provides a foundation for interpreting the evolutionary trajectory of past pandemics and predicting those of the future. The adaptability of this biophysical model extends to the key proteins of other viruses as well, signifying its potential in guiding public health interventions, and advancing our understanding of viral evolution.

Author affiliations: ^aDepartment of Chemistry and Chemical Biology, Harvard University, Cambridge, MA 02138; ^bJohn A. Paulson School of Engineering and Applied Sciences, Harvard University, Cambridge, MA 02138; ^cÉcole Polytechnique, Institut Polytechnique de Paris, Palaiseau 91128, France; ^dHarvard/MIT MD-PHD Program, Harvard Medical School, Boston, MA 02115 and Massachusetts Institute of Technology, Cambridge, MA 02139; and ^eProgram for Health Sciences and Technology, Harvard Medical School, Boston, MA 02115 and Massachusetts Institute of Technology, Cambridge, MA 02139

Author contributions: D.W., M.H., V.M., and E.I.S. designed research; D.W., M.H., and V.M. performed research; D.W., M.H., and V.M. analyzed data; and D.W., M.H., V.M., and E.I.S. wrote the paper.

The authors declare no competing interest.

This article is a PNAS Direct Submission.

Copyright © 2024 the Author(s). Published by PNAS. This open access article is distributed under Creative Commons Attribution License 4.0 (CC BY).

¹To whom correspondence may be addressed. Email: shakhnovich@chemistry.harvard.edu.

This article contains supporting information online at <https://www.pnas.org/lookup/suppl/doi:10.1073/pnas.2314518121/-DCSupplemental>.

Published May 31, 2024.

computational methods. Central to their premise is that fitness comes from the proportion of viral capsid proteins in a folded state free of antibodies, with state-occupation probability determined from Boltzmann distribution. SpikePro, a computational model, uses the spike protein's amino acid sequence and structure to estimate SARS-CoV-2 fitness. The model considers the stability of the spike protein, its binding affinity with ACE2, and the potential for immune evasion (13). While it has shown effectiveness in identifying dominant viral strains, it is worth noting that this is an empirical model whose foundation is not grounded in biophysical principles. Furthermore, experimental verification for the model's predictions has been scarce, which highlights the need for more rigorous, physics-based models.

The central aim of our study is thus to bridge the gap between viral fitness and biophysical properties of the RBD. We concentrate specifically on how emerging mutations affect both fitness and the binding energies of the RBD to ACE2 and antibodies. By doing so, we aim to develop a robust methodology based on statistical mechanics to forecast RBD fitness anchored in its molecular properties.

This study establishes a biophysical link between binding affinities and relative fitness for SARS-CoV-2 mutants. To that end, we have constructed a genotype-to-fitness mapping for the SARS-CoV-2 RBD under the constraints of successful cellular entry via ACE2 and influence of neutralizing antibodies. This mapping equips us with a predictive tool for assessing fitness of emerging SARS-CoV-2 variants.

Results

The Model. Our RBD fitness function is based on thermodynamics of protein folding and binding, as described in Fig. 1. For each mutant *mut*, the relative fitness compared to the wild type *wt*

(Wuhan-Hu-1), defined as $\frac{F_{\text{mut}}}{F_{\text{wt}}}$, follows the same relationship as the absolute fitness. To avoid confusion, we will simply refer to this relative measure as the “fitness” *F*.

We consider the contribution of fitness (infectivity) from RBD to the virus, denoted *F*, as proportional to the fraction of RBD that are folded and free from antibodies. We also consider the experimental observation that RBD on the spike protein could adopt two distinct conformations: “up” and “down”, with only “up” RBD exposing the receptor-binding motif (14–17). As a result, we adopt a multistate microscopic configuration model for the RBD that includes the following states: unfolded; folded in both up and down states and free; folded in up and down states, and bound to ACE2; and finally, folded, in up and down states, and bound to one of four distinct antibodies. These configurations correspond to respective free energies G_u , G_f^{\uparrow} , G_f^{\downarrow} , G_{bA}^{\uparrow} , G_{bA}^{\downarrow} , G_{ai}^{\uparrow} , G_{ai}^{\downarrow} where *i* indexes antibodies.

In our model, the RBD can only be bound to one antibody at a time, or ACE2 and to no antibody, or be free from ACE2 and antibodies. We then assume that an RBD can exist at thermodynamic equilibrium over these states at some finite temperature T_s , denoted by $\beta = 1/(k_B T_s)$, where k_B is the Boltzmann constant. Following Rotem et al. (12), we then propose that fitness *F* is proportional to the probability of RBD being found in the folded state free from antibodies or in the folded, bound to ACE2 state. This approach is grounded in the understanding that while binding to the ACE2 receptor is crucial for initiating infection within an individual, the transmission of the virus between individuals is predominantly driven by “free viruses”—those not yet bound to any receptors (18). These free viruses, contained in the host’s body and airborne droplets, represent a potential state for initiating further infections in new hosts.

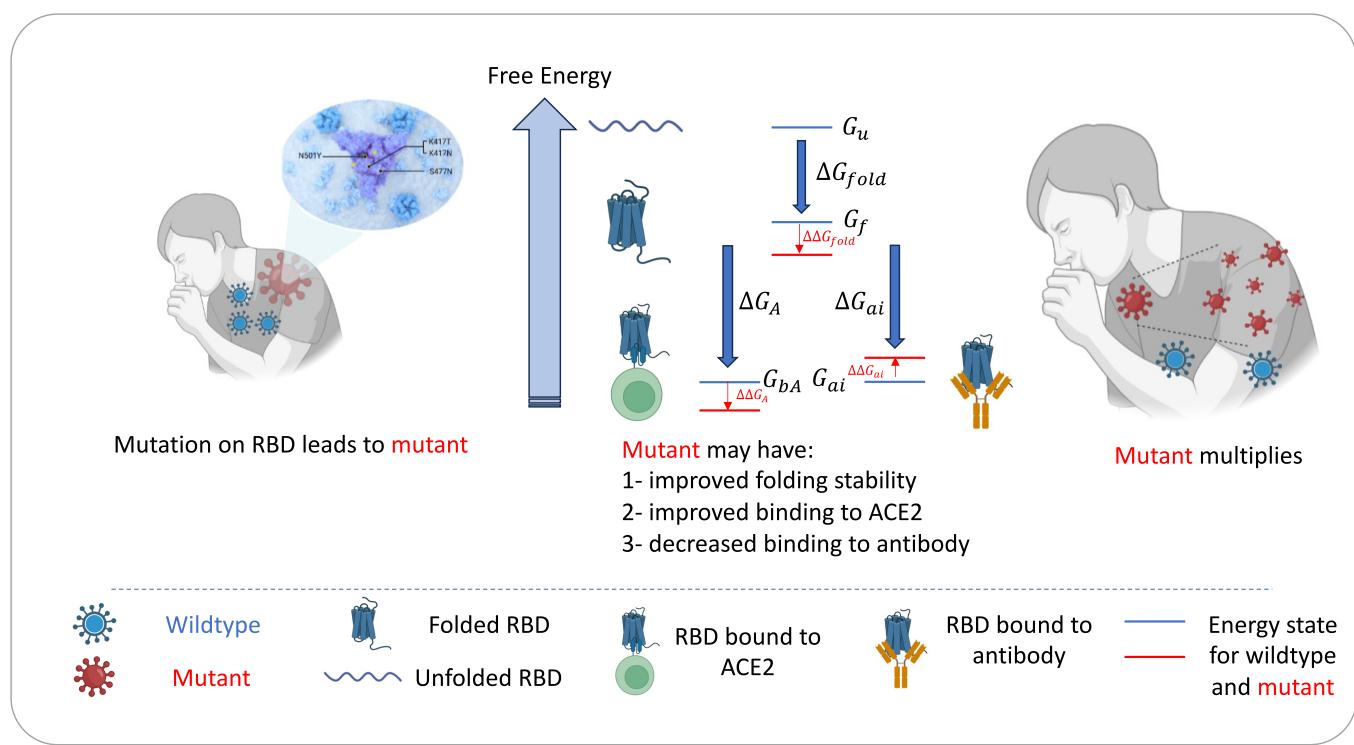


Fig. 1. Model illustration. Link between fitness and thermodynamics of protein folding and binding. High fitness variants may exhibit improved stability in the folded state or in the ACE2-bound state to facilitate cellular entry or have the capacity to destabilize bound-to-antibody states, thereby enabling evasion.

We define the partition function for the up RBD as $\Xi^{\uparrow} = Ce^{-\beta G_{ba}^{\uparrow}} + \sum_i C_i e^{-\beta G_{ai}^{\uparrow}} + e^{-\beta G_f^{\uparrow}}$, where $C = \frac{[ACE2]}{C_0}$ and $C_i = \frac{[Ab_i]}{m_i C_0}$. Here, [...] represents concentration, standard reference concentration C_0 allows us to express C and C_i in dimensionless form, and m_i accounts for the quantity of antibodies required to neutralize the virus. Similarly, the partition function for the down RBD, Ξ^{\downarrow} , can be defined analogously, with all instances of \uparrow replaced by \downarrow . Fitness F can be expressed as

$$F = \lambda \frac{Ce^{-\beta G_{ba}^{\uparrow}} + e^{-\beta G_f^{\uparrow}} + Ce^{-\beta G_{ba}^{\downarrow}} + e^{-\beta G_f^{\downarrow}}}{\Xi^{\uparrow} + \Xi^{\downarrow} + e^{-\beta G_u}}, \quad [1]$$

where λ is a scaling factor. Considering the experimental challenge of measuring free energy for the down state RBD, we need to further simplify this equation.

Since the down RBD does not bind with ACE2, we have $Ce^{-\beta G_{ba}^{\downarrow}} = 0$. We further assume that for RBD variants, the population ratio of down RBD vs. up RBD, is given by $\frac{e^{-\beta G_f^{\downarrow}}}{e^{-\beta G_f^{\uparrow}}} = k$. For a given antibody i , we express the difference in binding free energy between the up and down states of the RBD as $G_{ai}^{\downarrow} = \epsilon_i + G_{ai}^{\uparrow}$.

Defining free energy differences between states as $\Delta G_{fold}^{\uparrow} = G_u - G_f^{\uparrow}$, $\Delta G_A^{\uparrow} = G_{ba}^{\uparrow} - G_f^{\uparrow}$ and $\Delta G_{ai}^{\uparrow} = G_{ai}^{\uparrow} - G_f^{\uparrow}$, and given that naturally occurred RBDs are stable—with the unfolded states having significantly higher free energy than the folded states (that is, $e^{-\beta \Delta G_{fold}^{\uparrow}} \ll 1$; see *Materials and Methods*)—we can simplify the model as follows:

$$F = \lambda \frac{\tilde{C} e^{-\beta \Delta G_A^{\uparrow}} + 1}{\tilde{C} e^{-\beta \Delta G_A^{\uparrow}} + \sum_i \tilde{C}_i e^{-\beta \Delta G_{ai}^{\uparrow}} + 1}, \quad [2]$$

where $\tilde{C} = \frac{C}{k+1}$ and $\tilde{C}_i = \frac{(1+e^{-\beta \epsilon_i})C_i}{k+1}$. In the yeast display experiments, the dissociation constant, K_D , is measured on isolated RBD with no binding sites obstructed. Therefore we take K_D as a proxy for K_D^{\uparrow} , resulting in $\Delta G^{\uparrow} \propto \ln(K_D)$. Consequently, we can express fitness as a logistic function of the logarithm of the dissociation constants K_{DA} and K_{Dai} for ACE2 and antibodies. Our study is centered on inferring F by fitting a scaling parameter λ and effective molecular concentrations in population \tilde{C} and \tilde{C}_i to the biophysical model (*Materials and Methods*).

Fit of Biophysical Model to Fitness Obtained from Population Data.

We separated our data into training and testing sets (see *Materials and Methods* for details) and then calibrated our biophysical model using observed viral variants (19) from the population study and incorporated experimental measurements of binding affinities (8, 9) as an input.

Training of the model involved adjusting six parameters (λ , \tilde{C} , and \tilde{C}_i for i ranging from 1 to 4) to achieve the highest correlation between model prediction and fitness inferred from population data in the training set.

Given the lack of precise information for k and ϵ_i for each variant, and to identify regression coefficients that are invariant across RBD variants, the k and $e^{-\beta \epsilon_i}$ are assumed to be constants across RBD variants. The validity of these assumptions will be further discussed in *Discussion* and *SI Appendix*. Furthermore,

energy scale T was fixed to 1.6, so that fitted concentrations (on the order ~ 1 to ~ 100 nM) agree with values of antibody concentrations observed in human serum (between 1 and ~ 60 nM in severe symptoms) (20). For the calculation of the effective concentration, we estimate that m is in the range of 10 to 100 (12) (see *Materials and Methods* and *SI Appendix* for details about fitting).

To demonstrate the predictive power of our biophysical model, we trained the model on 2% of the observed variants (22 points), achieving a highly satisfactory fit (average $R^2 = 0.97$) on the training set as shown in Fig. 2A. This result was further corroborated by the predictive power demonstrated on the test set with around 1,000 variants (average $R^2 = 0.91$), thus confirming the absence of overfitting (Fig. 2B). Our model's performance reflects the biophysical understanding that a combination of strong ACE2 binding and reduced antibody binding from the four chosen antibodies could provide the virus with a fitness strategic advantage, enhancing its ability to spread in the population.

To further explore the behavior of our model, we fixed all but one of the dissociation constants at their mean value across mutants and studied the variation of inferred relative fitness as a function of the unfixed constant. Interestingly, we observed that variants carrying a combination of Omicron mutations consistently fell into the upper plateau or linear segment of our biophysical model (Fig. 2 C–G). This suggests that natural selection did not favor combinations of mutations that would lead to high antibody binding. It should be noted that in each curve on Fig. 2 C–G, four out of five dissociation constants are fixed at the mean values. Thus, the inferred fitness for each data point in these one-dimensional cross-sections of the fitness landscape does not represent its actual fitness.

Predicting Fitness for Variants between wt and BA.1. Our model constitutes an effective tool for forecasting fitness contribution of RBD to the virus, given molecular properties of their RBD. We first examine the predictive power for variants between *wt* and BA.1.

In Fig. 3A, the model is trained on the experimental dataset that excludes the G446S mutation and is subsequently tested on variants containing this specific mutation. Remarkably, our model even succeeds in predicting fitness of variants bearing the G446S mutation, a demanding task considering that this mutation causes escape from the REGN10987 antibody (9). The model's ability of predicting the impact of the G446S mutation on fitness, despite the complexity of correlating complete immune escape from REGN10987 to fitness, highlights the potential of our model to predict effects of unseen mutations.

This point is further highlighted in Fig. 3B. For each mutation, labeled as m , our model was systematically trained on all variants that exclude mutation m and then used to predict fitness of all existing variants carrying mutation m in combination with other possible mutations. This methodology, implemented across all 14 mutations, yielded a strong R^2 coefficient of over 0.8. The efficacy of the model is demonstrated by its ability to consistently and accurately predict the fitness of previously unobserved mutation between Wuhan-Hu-1 and Omicron BA.1.

To further demonstrate the predictive capacity of our model, we assessed its ability to forecast fitness of future variants. The biophysical model was fitted on the dataset comprising all registered variants prior to May 1, 2021, which included only 42 data points. This trained model was then deployed to predict fitness of subsequent variants that emerged between May 1, 2021,

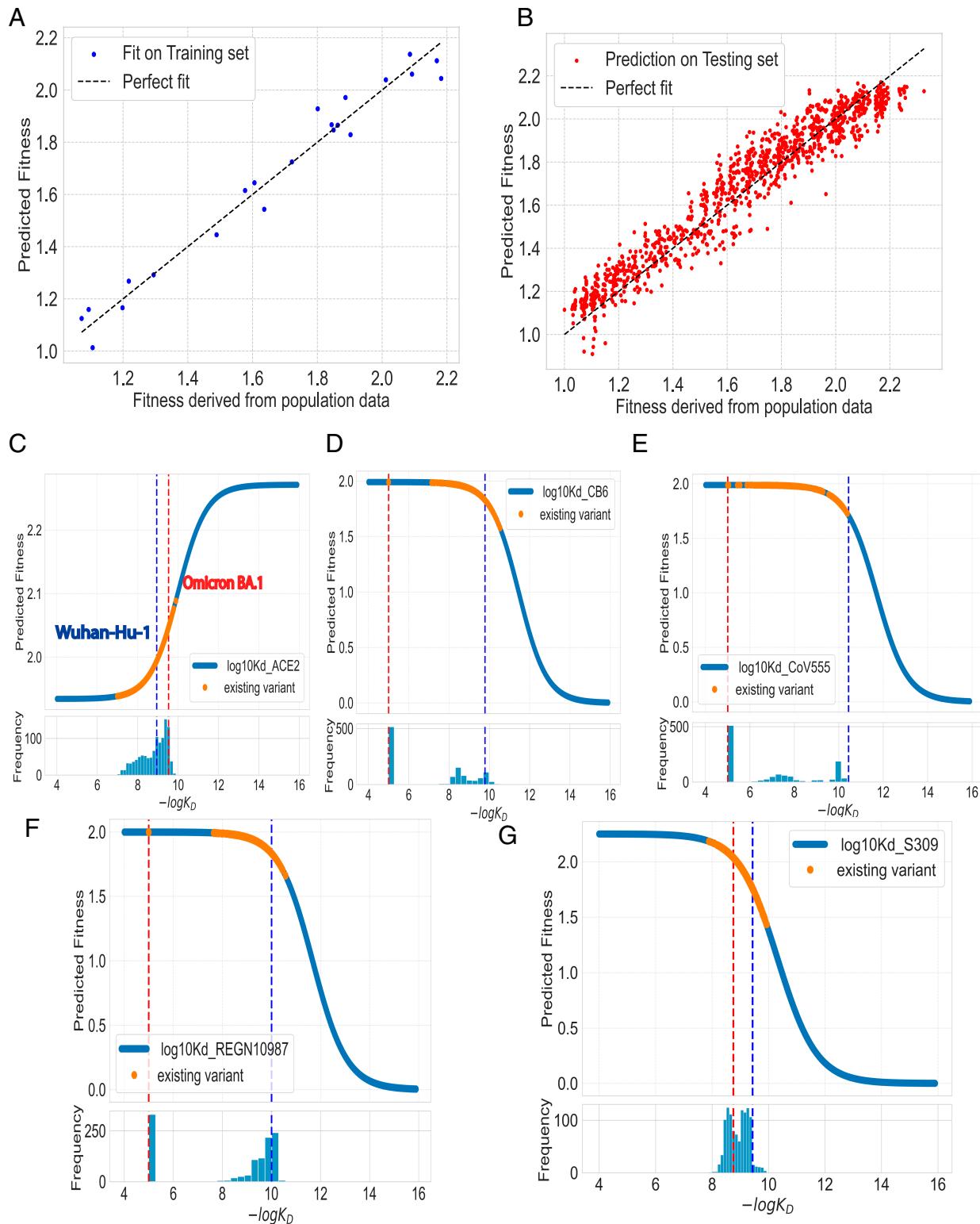


Fig. 2. Biophysical model analysis. (A) Fit of the predictive model on the training set. Each dot represents a variant in the training dataset, plotted against the fitness derived from population data on the x-axis and the predicted fitness by the model on the y-axis ($R^2 = 0.97$). (B) Model's performance on the testing set. The model's predictions align well with the fitness observed in population, as reflected by $R^2 = 0.94$ suggesting that the model maintains a strong predictive power on unseen data. (C–G) The dependence of fitness function on the logarithm of each dissociation constant. Existing variants are depicted as yellow regions on the curve, and a histogram showing the distribution of the dissociation constants is provided beneath each plot. Red and blue dashed lines represent wt and BA.1 respectively.

and May 1, 2022, a period in which 843 unknown variants arose. Importantly, our model exhibited a notable predictive power, accurately predicting the leap in fitness induced by the Omicron BA.1 variant and its neighboring variants, as depicted in Fig. 3C, with an R^2 value of 0.76. This underscores the model's efficiency in predicting fitness of unseen variants, even when trained on relatively sparse data from the early stages of the pandemic.

Predicting Fitness for Variants without Experimental K_D . Moreover, while our model has been trained on mutations spanning from the Wuhan-Hu-1 to the Omicron BA.1, it is not restricted to this specific spectrum. The model exhibits the capability to capture the effects of mutations outside the scope of Wuhan to Omicron, provided the relevant RBD's biophysical properties are available, either through experimental data or simulations.

To estimate the fitness of variants beyond the experimental dataset of 2^{15} sequences, it is essential first to determine the K_D values. We employed a supervised learning approach, utilizing transformer embedding and neural networks (see *Materials and Methods* for details), proven effective in estimating the binding affinity of unseen RBD as demonstrated in Han et al. (21). It is important to note that while the method we employed may not represent the most state-of-the-art K_D estimator available, its application here aims to illustrate how our biophysical fitness model can be synergistically integrated with either simulation or ML-based K_D estimators.

Our ML model is trained using 20,000 variants between Wuhan and Omicron BA.1 and their corresponding K_D values. The remaining 12,565 variants were then utilized for validation. On the validation set, we achieved R^2 of 0.89, 0.98, 0.84, 0.75, and 0.79, respectively, for ACE2 and four antibodies, as demonstrated in *SI Appendix*, Fig. S1.

Then the ML model was applied to all RBD sequences observed in GISAID for which we could calculate population infectivity but lacked K_D values for the biophysical model. We estimate these K_D values using our ML model, then feed them into the biophysical model and compared against actual fitness metrics derived from population data as shown in Fig. 3D. Importantly, this pipeline allows us to predict fitness for new variants, based solely on their sequences. Despite being trained on variants between Wuhan and Omicron BA.1, the model could correctly predict the fitness of variants with combinations of unseen mutations, achieving a Spearman correlation of 0.88. A selection of these variants is identified and labeled in Fig. 3D for reference. Please note that for Fig. 3D and subsequent panels E and F, variants between Wuhan-Hu-1 and Omicron BA.1, for which we already have experimental K_D data, are excluded from the analysis.

In Fig. 3E, we assess the model's capacity to predict RBD fitness for variants that emerged in 2020, 2021, and 2022. Our model not only accurately predicted the fitness of various variants but also consistently identified the top variants within each of these time frames. However, in 2023, a noticeable divergence emerged between the predicted and actual fitness, particularly for the Omicron subvariants BA.2, BA.4, and BA.5. While our model successfully recognized them as the most infectious variants of the year, it was unable to discern the fitness differences among these subvariants. We hypothesize that this limitation arises because the evolution of these variants is not predominantly driven by the four antibodies studied in this paper, given that BA.1 had already demonstrated escape from three out of the four antibodies. We then studied the model's capacity to infer

the fitness trend during the pandemic in Fig. 3F. Although our model's predictions start to diverge from the actual fitness metrics from early 2022 onward due to the underestimation of the most-fit variant, it maintains alignment with the general trend post-2022, showcasing its continued capability in accurately predicting fitness for other variants that occurred during this period.

Epistasis. After fitting our biophysical model to all existing variants, we extrapolated fitness across all 32,768 possible mutation combinations which we have experimental K_D data. Notably, we observed a fitness threshold for the RBD, beyond which additional mutations cease to enhance fitness (22). Essentially, once the virus achieves a certain fitness level, characterized by high immune escape and robust binding to the cell receptor, it becomes increasingly challenging for further mutations to improve this balance. As a result, fitness begins to plateau, as illustrated in Fig. 4A. This phenomenon, which was not predicted by the work of Obermeyer et al. (19) or other nonepistatic models, can be attributed to two key factors. First, as suggested by Moulana et al. (9), mutations that enhance antibody escape tend to reduce the virus's affinity for ACE2 receptors. This indicates a trade-off in viral evolution between immune evasion and the ability to infect host cells, a factor inherently accounted for in K_D measurements and integrated into the biophysical model. Second, the logistic function used in our biophysical model naturally leads to a plateau in fitness.

To calculate the epistatic coefficients, we utilized a linear model described in *Materials and Methods*. We applied this model to fitness values inferred from our biophysical model across all 32,768 possible combinations of mutations. The performance of the model was evaluated (expressed as R^2) for different orders of epistasis on a withheld test dataset, constituting 10% of the total data. A first-order model (comprising 16 coefficients) produced an R^2 value of 0.958, while a second-order model (comprising 121 coefficients) achieved an R^2 of 0.985. An F-test was conducted, yielding a p-value of 10^{-16} . This confirms the robustness of the second-order model despite its increased parameter count over the first-order epistatic model. The high correlation values and satisfactory representation of the data suggest that a second-order epistatic model sufficiently captures the key dynamics, thereby alleviating the necessity for higher-order epistatic models. Following the training of the second-order model across the complete dataset to derive final coefficients, we observed that most single mutations (reflected on the diagonal of the matrix in Fig. 4C) have a positive impact on fitness. An interesting aspect of our findings is the behavior of the Y505H mutation. While it displayed a negative first-order coefficient, its interactions with other mutations are positive, which effectively offsets the first-order negative impact.

Furthermore, we observed a unique dynamic with the Q493R mutation. Alone, this mutation is beneficial, contributing to the virus's escape from the LY-CoV555 antibody. However, its co-occurrence with other mutations reverses this benefit by reducing the binding affinity to ACE2 (23). Particularly noteworthy are adjacent mutations in the crystal structure such as Q493R-E484A, Q493R-Q498R, and Q493R-K417N. These proximal mutation pairs demonstrate significant deleterious second-order effects, as illustrated in Fig. 4B. Intriguingly, the evolutionary trajectory of Q493R seems responsive to these complex interactions. A reversal of the Q493R mutation was observed in subsequent lineages, including BA.4, BA.5, BA.2.75,

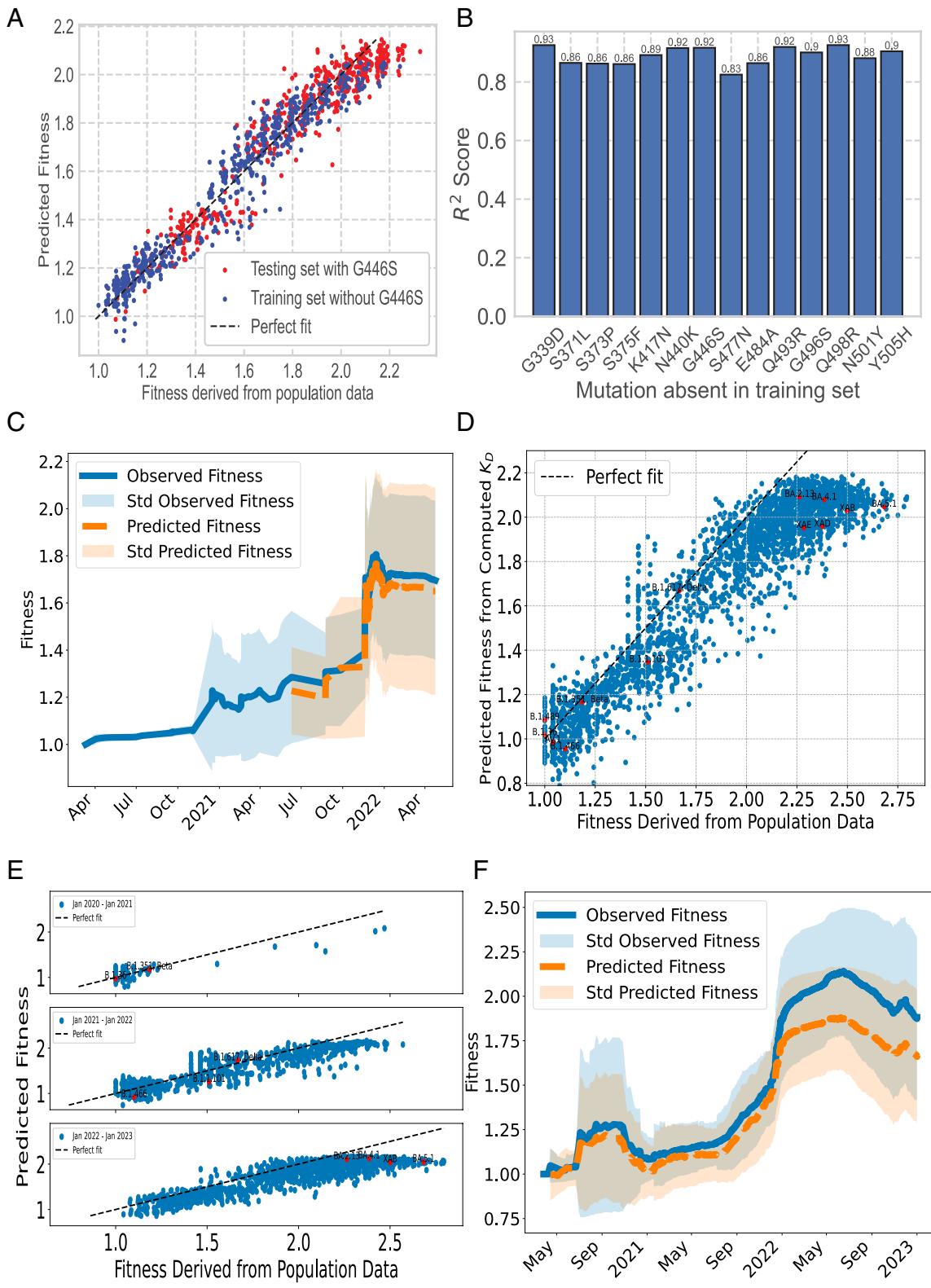


Fig. 3. Assessing predictive power of the model using experimental K_D (A–C) and computed K_D (D–F). (A) Fitness prediction for variants carrying the G446S mutation which was not included in the training set yields $R^2 = 0.92$. (B) R^2 derived from a model trained on variants excluding a specific mutation, then used to predict fitness of variants exhibiting that mutation. (C) Predictions of fitness compared with actual fitness trend for variants between Wuhan-Hu-1 and Omicron BA.1. Variants observed before May 2021 are used as training set for the model. The model uses experimental K_D . (D) Predicted fitness from biophysical model against actual fitness derived from population data. The model uses K_D acquired from ML and is fit on Wuhan-Omicron set. Selected variants are highlighted. (E) Predicted fitness compared with actual fitness for variants observed in 2020, 2021, and 2022. (F) Predictions of fitness using ML derived K_D over three-month rolling windows, compared with the actual fitness trends for variants between 2021 and 2023.

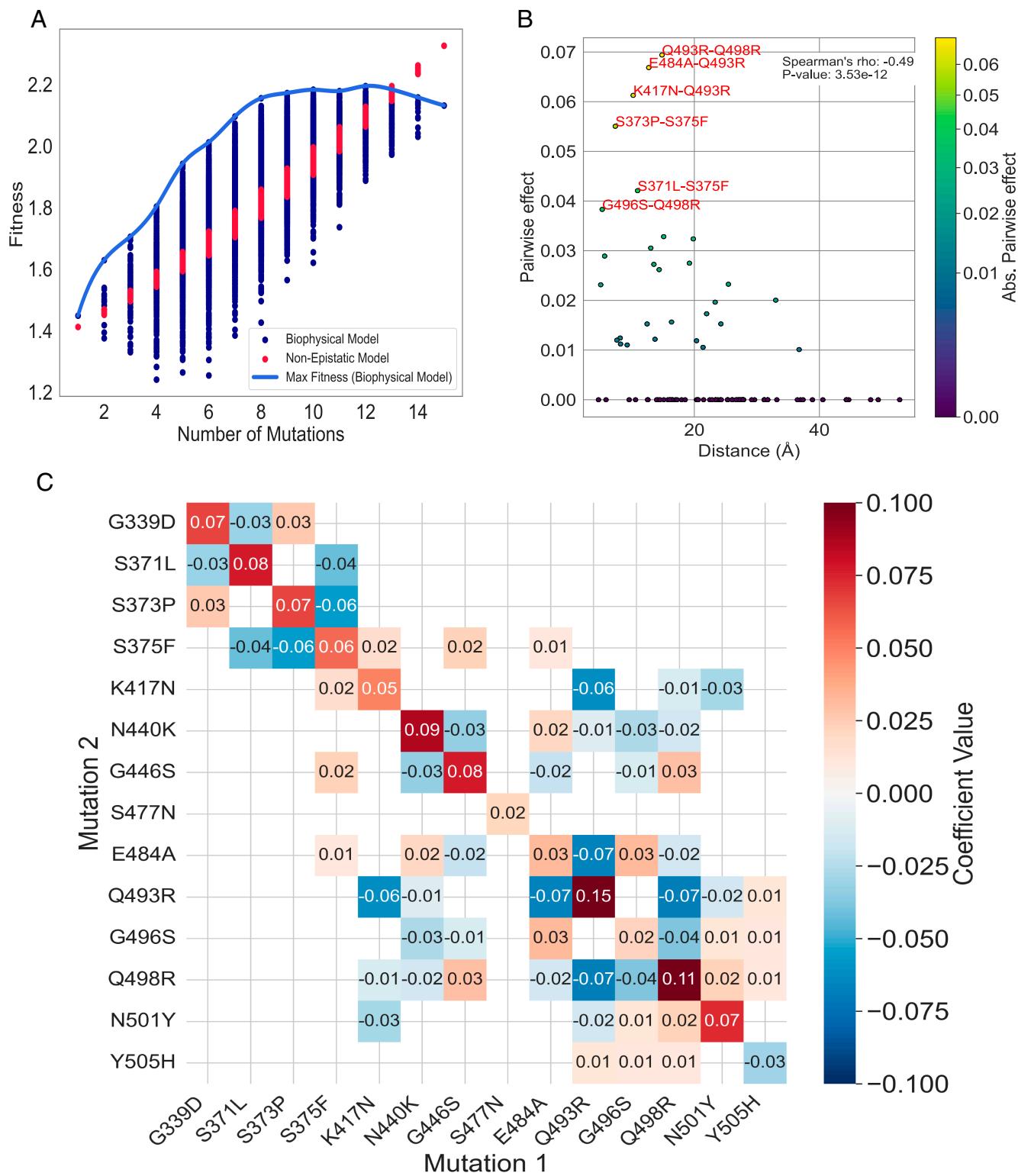


Fig. 4. Epistasis analysis. (A) Predicted fitness values with the nonepistatic model of Obermeyer et al. and with our epistatic biophysical model plotted against the genome's mutation count, for all mutation combinations with mutation T478K. We define "Max Fitness" as the maximum fitness prediction from our biophysical model. Max Fitness curve begins to plateau with a higher mutation count, demonstrating a diminishing returns effect in epistatic. (B) Pairwise (second order) interaction coefficients against the spatial distances between the corresponding residues, with mutations colored in accordance with the absolute value of their pairwise coefficient. (C) Coefficients of epistasis: Diagonal coefficients denote first-order interactions, whereas off-diagonal coefficients represent second-order interactions. Coefficients smaller than 0.01 have been masked for clarity.

BQ.1, and XBB. This reversal suggests an evolution dynamics driven by the trade-offs between immune escape and ACE2 binding.

Overall, these findings highlight the complex dynamics of viral evolution, where multiple mutations interact nonlinearly to enhance or reduce viral fitness.

Discussion

Although the complexity of fitness landscapes is undeniable, recent research indicates that in certain biologically and clinically significant systems, such as evolution of bacterial resistance against antibiotic (24), evolution of viral resistance against antiviral treatments (11, 25), as well as norovirus evolution against a neutralizing antibody (12), these fitness landscapes can be systematically and quantitatively delineated. Our research builds upon these findings and reveals that the fitness landscape of the SARS-CoV-2 RBD, undergoing evolution against neutralizing antibodies, can be systematically described through its biophysical properties. In our study, we found the strength of binding to the neutralizing antibody, as well as to ACE2, plays crucial roles in determining RBD fitness. Importantly, the significance of these biophysical parameters is not confined to SARS-CoV-2 alone. Similar traits, namely antibody binding affinity and protein folding stability, have been crucial in influencing the evolution of influenza viruses (25–27). This observation suggests that our model may have broader applications, potentially extending to other viruses beyond SARS-CoV-2.

To forecast the fitness of emerging novel variants with our model, acquiring the dissociation constants for the corresponding mutated RBD is imperative. While experimental data sets the gold standard, obtaining them early in a pandemic may pose challenges. Fortunately, computational methodologies have demonstrated efficacy in predicting dissociation constants for unobserved mutations using molecular dynamics (MD) and ML. For instance, Lacam et al. (28) achieved exceptional accuracy by employing a framework that integrates MD and potential-of-mean-force calculations. Their study accurately determined the binding free energy of the RBD for four prevalent variants and the wild type, when complexed with ACE2 or antibodies S2E12 and H11-D4. In parallel, Sergeeva et al. (29) offered an effective strategy to determine the impact of interfacial mutations on the binding affinities between RBD and ACE2, using free energy perturbation. Williams et al. (30) constructed a multilayer neural network, using biophysical parameters as inputs to predict binding affinities of SARS-CoV-2 antibodies with various VoCs. Similarly, Chen et al. developed the NN-MM-GBSA model (31) using MD and neural network to predict binding affinity between SARS-CoV-2 spike RBD variants and ACE2, reaching a correlation coefficient of 0.73 on prediction of dissociation for single variants in the work Starr et al. (7).

In this study, for RBD variants without experimental measured K_D , we estimated binding affinities using a computationally efficient alternative by employing a neural network that takes in transformer embedding for sequences and outputs K_D . In this way, we could predict mutational effects not only for combinations of mutations outside of the experimental dataset but also for completely unseen mutations. By incorporating these predicted dissociation constants into our biophysical model, we demonstrate that this streamlined computational approach yields accurate predictions for VoCs both at the early and later stages of the pandemic.

During the early stages of a pandemic, data on viral fitness or infectivity is typically limited. Unlike K_D measurements, which can be derived from wet lab experiments or simulations, comprehensive fitness data are usually only accessible after a variant has extensively spread and been subjected to population-level sequencing. Therefore, the ability of our model to train and predict effectively with minimal data points is particularly valuable in the context of pandemic preparedness. Considering the RBD's susceptibility to mutations, our model has the

potential to be a powerful tool in understanding and predicting the fitness of emerging variants.

Despite our biophysical model being trained on nonepistatic population data from Obermeyer et al. (19), the model allows us to generate an epistatic map from genotype to fitness using k_D . In our model, we observe a tendency for fitness to plateau in the face of the increasing number of mutations relative to the wild type. This phenomenon of diminishing returns, manifested as a fitness plateau, has been extensively studied in the existing literature (22, 32–35).

Our results also indicate that epistasis constrains evolution. While many of the mutations we investigated are beneficial, specific combinations of mutations could be deleterious. Some mutations require the concurrent occurrence of stabilizing mutations to counterbalance their adverse consequences. This observation is consistent with findings of Gong et al. (25) and Rodrigues et al. (24) underscoring the critical role of stabilizing mutations in fixation of subsequent destabilizing mutations that could hold adaptive value.

Our results emphasize the importance of accounting for the interactive effects of multiple mutations in viral evolution modeling and prediction. It is particularly noteworthy that our model could explain the reversal of Q493R from a viral fitness perspective. Furthermore, in our model, second-order pairwise effects between mutations tend to weaken as their separation in the crystal structure increases. These observations demonstrate that our model effectively captures the essential aspects of viral fitness, including the epistatic effects that drive viral evolution.

A fundamental assumption of our model is that evolution of the SARS-CoV-2 RBD is predominantly driven by its capacity to bind the ACE2 receptor and evade antibodies. Our model, which is based on this assumption, demonstrates the ability to predict the fitness effect of all mutations in RBD, except for T478K. T478K is an interesting mutation as shown by Moulana et al. (8, 9) that it had a negligible effect on dissociation constants, despite its strong contribution to variant infectivity. The exact reason for T478K's high fitness contribution, despite no apparent change in RBD's biophysical parameters, remains unclear. A prevailing hypothesis is its frequent co-occurrence with the D614G mutation (36). D614G, located on the spike protein but outside of the RBD, has been shown to enhance viral replication and infectivity (37, 38). Our model, focusing on the RBD, does not account for effects of mutations like D614G on the spike protein. This limitation leads to our observation that the RBD containing T478K consistently shows increased fitness, regardless of the mutational background in RBD. Consequently, we segregated the training data based on the presence of the T478K mutation. This approach helps us account for the fitness increase associated with T478K, despite our model's focus on the RBD.

Studies have established that RBD could adopt either up or down conformations on the spike protein (14–17). Despite this, many investigations have measured binding affinities using isolated RBDs, neglecting the complex dynamics between the spike protein and RBD. In the main text and *SI Appendix*, we demonstrated that incorporating both conformational states into our model achieves a logistic function analogous to that derived from considering the RBD solely in its up state, with effective molecular concentration renormalized as $\tilde{C} = \frac{C}{k+1}$ and $\tilde{C}_i = \frac{(1+e^{-\beta e_i})C_i}{k+1}$. This provides a theoretical foundation for fitting the logistic regression with dissociation constants measured on isolated RBDs, given two assumptions. One of the assumptions is that k is roughly a constant across different

variants. This is experimentally verified as Omicron BA.1 Spike preferentially adopts the one-RBD-up conformation (39) similar to the wild type (14–16). However, we acknowledge that this assumption may become less reliable as mutations accumulate. For instance, Cryo-EM studies indicate a shift toward more 3-RBD-down configurations in the Omicron BA.2 variant (40), potentially explaining the observed underestimation of BA.2's fitness in our model. Another assumption is that $e^{-\beta\epsilon_i}$ is also a constant across variants. In the context of REGN10987 and S309 antibodies, which bind outside of the receptor-binding motif, we estimate $e^{-\beta\epsilon_3} \approx e^{-\beta\epsilon_4} \approx 1$. For LY-CoV016 and LY-CoV555, whose binding sites overlap with the receptor-binding motif $0 < e^{-\beta\epsilon_1} < 1$ and $0 < e^{-\beta\epsilon_2} < 1$.

On the other hand, extrapolating our fitness predictions beyond RBD to encompass the entire viral sequence, presents a more formidable challenge. Our model's predictive power on RBD fitness comes from the fact that binding affinity to ACE2 and antibodies are major evolutionary forces that drive RBD's evolution. It is likely that mutations outside the RBD will have functional and structural effects that extend beyond alterations to dissociation constants with cell receptors and antibodies, complicating the predictions.

Furthermore, our model currently considers antibody binding to ACE2 and four monoclonal antibodies: LY-CoV016, LY-CoV555, REGN10987, and S309, an oversimplification given the complexity of human immune responses. While our model already exhibits accurate predictions and could easily be extended to other antibodies if data are available, the inclusion of additional factors such as other antibodies, vaccination effects, the replicative capacity within the infected cell (41), transmission dynamics (42), and drug resistance (43) could potentially enhance its predictive power and realism.

Materials and Methods

RBD Fitness Data Analysis. We acquired the fitness ratio of each RBD compared to the wild type from the work of Obermeyer et al. (19). In their study, fitness label is obtained by modeling the relative growth rate of SARS-CoV-2 lineages using a hierarchical Bayesian regression model. The model combines individual mutations and clusters genetically similar genomes to estimate the incremental effect of amino acid changes on growth rate within each lineage, which enables the model to share statistical strength among similar lineages. Specifically, the proportion of lineages is modeled as a multinomial distribution whose probability parameter is a multivariate logistic function softmax ($\alpha + tb/\tau$). For each lineage, the slopes b are linearly regressed against the presence of each possible amino acid substitution $X_m \in \{0, 1\}$ as $b = \sum_m b_m X_m$. These linear coefficients b_m can be directly interpreted as the effect of a mutation m on a lineage's fitness. This model assumes each single point mutation independently linearly contributes to change in fitness. Authors reported that fitting a similar model of both single and pair mutations leads to no pairwise mutations stronger than the top 100 single mutations.

This enabled us to estimate fitness F_{mut} of each existing RBD mutant, compared to wild type. Using $b = \log(\frac{F_{mut}}{F_{wild}})$, we get

$$\frac{F_{mut}}{F_{wild}} = \exp \left(\sum_{m \in \text{RBD}} b_m X_m \right).$$

RBD Binding Affinity. We acquired the binding affinity data from the work of Moulana et al. (8, 9). In their study, they systematically examined the interactions between all possible combinations of 15 RBD mutations (totaling 32,768 genotypes) and ACE2, as well as four monoclonal antibodies (LY-CoV016, LY-CoV555, REGN10987, and S309) via Tite-seq measurement. In situations where binding affinities in their dataset were too low to measure accurately,

we have chosen to substitute these with a fixed value of 5. We stress that this choice of value is not expected to influence our study's outcomes. As indicated in Fig. 2 C–G, the antibody escape largely resides in the upper plateau region of the fitness curve; thus, this preset value for variants with immune escape to antibodies does not have a substantial impact on the logistic regression results. Furthermore, we eliminated approximately 100 genotypes from the analysis that did not have measured ACE2 binding affinity.

Effect of Mutations on RBD Stability. In the presented results, we excluded the unfolded state from the model, on the assumption that this variable shows minimal variation across different variants and hence, would not significantly influence fitness. To verify this assumption, we employed DDGUN, an untrained, high-throughput tool (44), to compute $\Delta\Delta G_{fold}$, the variance in folding free energy difference, for mutants relative to the wild-type: The maximum variation was under 2 kcal/mol. Recalling that $\Delta G_{fold} \approx -10$ kcal/mol (45), we deduce that most mutations do not significantly destabilize RBD.

This could be indicative of the universally efficient folding of RBDs seen in nature. The selection pressure acting on the RBD primarily focuses on binding to ACE2 and immune evasion, and therefore the mutations are predominantly on the protein surface and do not significantly affect the protein's stability.

Filtering RBD Sequences from GISAID. All 15,371,428 spike sequences on GISAID (10) as of April 14, 2023 were downloaded and aligned, following the approach in the work of Starr et al. (46). Sequences from nonhuman origins and with lengths outside [1260, 1276] were removed. They were then aligned via mafft (47) and sequences containing unicode errors, gap, or ambiguous characters were removed. Overall, we retained 11,976,984 submissions represented by 25,725 unique RBD sequences. RBD amino acid mutations were enumerated compared to the reference Wuhan-Hu-1 SARS-CoV-2 RBD sequence (Genbank MN908947, residues N331-T531).

We then remove all RBD sequences that do not match any of the possible intermediates between Wuhan Hu-1 and Omicron BA.1. To do this, we allow all possible combinations of 15 mutations (G339D, S371L, S373P, S375F, K417N, N440K, G446S, S477N, T478K, E484A, Q493R, G496S, Q498R, N501Y, and Y505H) between Wuhan and Omicron BA.1, which lead to $2^{15} = 32,768$ possible combinations. We calculated the number of occurrences of each RBD sequence as well as the time of its first occurrence, which we approximated by taking 5% quantile of time data for each RBD sequence. From this analysis, we obtained 1,121 unique observed RBDs. These RBDs correspond to 2.5 million sequences out of the 12 million sequences we initially screened from the GISAID database.

Fitting the Model with Logistic Regression. For the purpose of logistic regression analysis, we utilized the intersecting data obtained from *Materials and Methods* RBD Fitness and RBD Binding Affinity. This accumulated dataset comprises 1,118 unique RBDs observed in the GISAID database (48), for which the K_D values have been experimentally determined. We further partitioned this data based on the presence or absence of the T478K mutation within the sequence, resulting in two distinct subsets (*SI Appendix* and *Discussion*).

The ratio between system temperature T_S (corresponding to body temperature of the host) and experimental temperature T_E [corresponding to the temperature of experiments conducted in work from Moulana et al. (8, 9), leading to $\Delta G = RT_E \ln(K_D)$] was treated as a hyperparameter T (simply referred as "energy scale"), which is a parameter whose value is chosen before the fit is done.

We calculated unknown parameters λ , \tilde{C} , and \tilde{C}_i by fitting:

$$F = \lambda \frac{\tilde{C} e^{-\ln(K_D)/T} + 1}{\tilde{C} e^{-\ln(K_D)/T} + \sum_i \tilde{C}_i e^{-\ln(K_{D,i})/T} + 1}, \quad [3]$$

where energy scale T is fixed to 1.6. We emphasize that the hyperparameter can be chosen with knowledge of the training set alone (and thus does not invalidate prediction capabilities of the model) and that the model behavior is only slightly affected by changes in this hyperparameter (see *SI Appendix* for model performance at different energy scales).

Model fitting was performed with nonlinear least square regression (*scipy.optimize* package) on a randomly selected training set. The model is then evaluated on the remaining testing set to prove absence of overfitting. Additionally, to mitigate the effects of randomness, we implemented 10-fold cross-validation wherever feasible.

Fitness Prediction Using ML Estimated K_D . Considering the potential unavailability of experimental measures for dissociation constants early in a pandemic, we advocate for the application of our methodology to K_D estimates derived from a computationally inexpensive deep learning pipeline. Notably, Hie et al. successfully predicted viral escape using a machine learning technique designed for natural language processing (49). Similarly, Han et al. introduced an online platform based on deep learning models, incorporating transformers, for rapid prediction of binding affinity between RBD mutants and ACE2 (50).

Building on these developments, our study showcases a framework that combines a biophysical model with a K_D predictor utilizing protein sequence embedding and a neural network. This integrated system efficiently predicts dissociation constants for emerging variants, and the results are directly fed into the biophysical model. This method enables prompt and effective prediction of viral evolution in response to novel mutations.

ESM-1v (51) is a Transformer-based language model specifically trained on a diverse dataset of 98M protein sequences. This pretrained model inputs a given protein sequence and outputs a vector representation, or an "embedding," of that sequence. This embedding consists of the evolutionary information of the protein sequence, which is pivotal in enabling the machine learning model to predict the effects of various combinations of unseen mutations as demonstrated in these papers (21, 49, 51).

To elaborate the embedding mathematically, consider a protein sequence of length L , we describe it as a sequence of tokens $\mathbf{x} \stackrel{\text{def}}{=} (x_1, \dots, x_L)$. For the RBD, we have $L = 201$. During the forward pass of ESM-1v, we obtain the hidden representations from the final layer, denoted as $(\mathbf{h}_1, \dots, \mathbf{h}_L)$, with each \mathbf{h}_i being a vector in \mathbb{R}^{1080} . To generate a comprehensive representation of the entire sequence, we apply mean pooling to these vectors, resulting in a single sequence representation $\mathbf{z} = f_{\text{esm}}(\mathbf{x}) = \frac{1}{L} \sum_{i=1}^L \mathbf{h}_i$.

After acquiring the embedding, we trained a neural network that takes the embedding as input and output the variant binding affinity with ACE2 and four antibodies. This network was trained using K_D data from 20,000 variants and subsequently validated on a separate set comprising the remaining 12,565 variants.

After training, the neural network was deployed to estimate the K_D for existing variants with known population fitness, encompassing a total of 5,460 variants. Excluding the variants between Wuhan and Omicron BA.1 (for which we already have experimental binding affinities), we focused on evaluating the predictive power of the remaining 3,334 variants. These estimated K_D values were subsequently integrated into the biophysical model, as described in *Fitting the Model with Logistic Regression*.

Epistasis Analysis. Epistasis describes how mutation interactions can affect fitness F . If there is no epistasis then fitness can be described as linear

combinations of the presence of each mutation $X_m \in \{0, 1\}$, leading to a first-order epistatic model:

$$F = \sum_i c_i X_i. \quad [4]$$

If we consider pairwise epistatic interactions between mutated sites, we get a second-order epistatic model:

$$F = \sum_i c_i X_i + \sum_{i < j} c_{ij} X_i X_j, \quad [5]$$

c_i are considered as "first-order" epistatic coefficient, while c_{ij} are "second-order" epistatic coefficients as they illustrate the nonlinear epistatic interaction between mutated sites i and j .

To make sure the linear model does not overfit and can generalize on unseen data, we implemented a 10-fold cross-validation strategy (dataset split: 90%/10%) and identified a linear model involving first- and second-order coefficients as described in Eq. 5 gives a better representation of data than a first-order model ($R^2 = 0.958$ vs. $R^2 = 0.985$ on test set). An F-test was conducted to assess the statistical significance of the improvement in R^2 when moving from the 16-parameter first-order model to the 121-parameter second-order model. Our analysis resulted in an F-statistic of 540, with an extremely low P -value of 10^{-16} . The very low P -value rejects the null hypothesis that the additional parameters in the second-order model do not contribute to an improved fit.

To analyze the relationship between second-order epistatic coefficients c_{ij} and the distance between mutated sites, we calculated the latter as the Euclidean distance between the average position (computed as the mean of positions of all nonhydrogen atoms in the amino acid) of each mutated site.

Data, Materials, and Software Availability. The code for our analyses is available on GitHub at <https://github.com/Dianzhuo-Wang/COVID19-Biophysical-Model> (52). The fitness data used in this study are available in ref. 19 at <https://github.com/broadinstitute/pyro-cov> (53). The binding affinity data used in this work is described in refs. 8 and 9 and can be accessed at https://github.com/desai-lab/compensatory_epistasis_omicron (54) and https://github.com/desai-lab/omicron_ab_landscape (55).

ACKNOWLEDGMENTS. This work is supported by NIH R35GM139571 (to E.I.S.) V.M. acknowledges support from NIGMS T32GM144273, a Hertz Foundation Fellowship, and a Paul & Daisy Soros Fellowship. The content is solely the responsibility of the authors and does not necessarily represent the official views of the National Institute of General Medical Sciences or the National Institutes of Health. We gratefully acknowledge all data contributors, i.e., the Authors and their Originating laboratories responsible for obtaining the specimens, and their submitting laboratories for generating the genetic sequence and metadata and sharing via the GISAID Initiative, on which this research is based. We would like to thank Vaibhav Upadhyay, Krishna Mallela, Zechen Zhang, and Junlang Liu for useful discussions.

1. X. Deng *et al.*, Transmission, infectivity, and neutralization of a spike L452R SARS-CoV-2 variant. *Cell* **184**, 3426–3437.e8 (2021).
2. K. Leung, M. H. Shum, G. M. Leung, T. T. Lam, J. T. Wu, Early transmissibility assessment of the N501Y mutant strains of SARS-CoV-2 in the United Kingdom, October to November 2020. *Eurosurveillance* **26**, 2002106 (2021).
3. S. Celk *et al.*, Escape of SARS-CoV-2 501Y.V2 from neutralization by convalescent plasma. *Nature* **593**, 142–146 (2021).
4. V. Upadhyay, A. Lucas, S. Panja, R. Miyauchi, K. M. Mallela, Receptor binding, immune escape, and protein stability direct the natural selection of SARS-CoV-2 variants. *J. Biol. Chem.* **297**, 101208 (2021).
5. M. Letko, A. Marzi, V. Munster, Functional assessment of cell entry and receptor usage for SARS-CoV-2 and other lineage B betacoronaviruses. *Nat. Microbiol.* **5**, 562–569 (2020).
6. B. Ju *et al.*, Human neutralizing antibodies elicited by SARS-CoV-2 infection. *Nature* **584**, 115–119 (2020).
7. T. N. Starr *et al.*, Deep mutational scanning of SARS-CoV-2 receptor binding domain reveals constraints on folding and ACE2 binding. *Cell* **182**, 1295–1310.e20 (2020).
8. A. Moulana *et al.*, Compensatory epistasis maintains ACE2 affinity in SARS-CoV-2 omicron BA.1. *Nat. Commun.* **13**, 7011 (2022).
9. A. Moulana *et al.*, The landscape of antibody binding affinity in SARS-CoV-2 omicron BA.1 evolution. *eLife* **12**, e83442 (2023).
10. S. Elbe, G. Buckland-Merrett, Data, disease and diplomacy: GISAID's innovative contribution to global health: Data, disease and diplomacy. *Glob. Chall.* **1**, 33–46 (2017).
11. N. Chéron, A. W. R. Serohijos, J. M. Choi, E. I. Shakhnovich, Evolutionary dynamics of viral escape under antibodies stress: A biophysical model. *Protein Sci.* **25**, 1332–1340 (2016).
12. A. Rotem *et al.*, Evolution on the biophysical fitness landscape of an RNA virus. *Mol. Biol. Evol.* **35**, 2390–2400 (2018).
13. F. Pucci, M. Rooman, Prediction and evolution of the molecular fitness of SARS-CoV-2 variants: Introducing SpikePro. *Viruses* **13**, 935 (2021).
14. R. Yan *et al.*, Structural basis for the different states of the spike protein of SARS-CoV-2 in complex with ACE2. *Cell Res.* **31**, 717–719 (2021).
15. A. C. Walls *et al.*, Structure, function, and antigenicity of the SARS-CoV-2 spike glycoprotein. *Cell* **181**, 281–292.e6 (2020).
16. D. Wrapp *et al.*, Cryo-EM structure of the 2019-nCoV spike in the prefusion conformation. *Science* **367**, 1260–1263 (2020).
17. Y. Chen *et al.*, Broadly neutralizing antibodies to SARS-CoV-2 and other human coronaviruses. *Nat. Rev. Immunol.* **23**, 189–199 (2023).

18. C. C. Wang *et al.*, Airborne transmission of respiratory viruses. *Science* **373**, eabd9149 (2021).
19. F. Obermeyer *et al.*, Analysis of 6.4 million SARS-CoV-2 genomes identifies mutations associated with fitness. *Science* **376**, 1327–1332 (2022).
20. K. A. de Jong, H. Rosing, M. Vermunt, A. D. Huitema, J. H. Beijnen, Quantification of anti-SARS-CoV-2 antibodies in human serum with LC-QTOF-MS. *J. Pharm. Biomed. Anal.* **205**, 114319 (2021).
21. W. Han *et al.*, Predicting the antigenic evolution of SARS-CoV-2 with deep learning. *Nat. Commun.* **14**, 3478 (2023).
22. C. S. Wylie, E. I. Shakhnovich, A biophysical protein folding model accounts for most mutational fitness effects in viruses. *Proc. Natl. Acad. Sci. U.S.A.* **108**, 9916–9921 (2011).
23. Q. Wang *et al.*, Antibody evasion by SARS-CoV-2 omicron subvariants BA.2.12.1, BA.4 and BA.5. *Nature* **608**, 603–608 (2022).
24. J. V. Rodrigues *et al.*, Biophysical principles predict fitness landscapes of drug resistance. *Proc. Natl. Acad. Sci. U.S.A.* **113**, E1470–E1478 (2016).
25. L. I. Gong, M. A. Suchard, J. D. Bloom, Stability-mediated epistasis constrains the evolution of an influenza protein. *elife* **2**, e00631 (2013).
26. J. M. Fonville *et al.*, Antibody landscapes after influenza virus infection or vaccination. *Science* **346**, 996–1000 (2014).
27. E. Y. Klein *et al.*, Stability of the influenza virus hemagglutinin protein correlates with evolutionary dynamics. *mSphere* **3**, e00554–17 (2018).
28. E. G. Coderc, M. de Lacam, H. Chen Blazhynska, J. C. Gumbart, C. Chipot, When the dust has settled: Calculation of binding affinities from first principles for SARS-CoV-2 variants with quantitative accuracy. *J. Chem. Theory Comput.* **18**, 5890–5900 (2022).
29. A. P. Sergeeva *et al.*, Free energy perturbation calculations of mutation effects on SARS-CoV-2 RBD:ACE2 binding affinity. *J. Mol. Biol.* **435**, 168187 (2023).
30. A. H. Williams, C. G. Zhan, Fast prediction of binding affinities of SARS-CoV-2 spike protein and its mutants with antibodies through intermolecular interaction modeling-based machine learning. *J. Phys. Chem. B* **126**, 5194–5206 (2022).
31. C. Chen *et al.*, Computational prediction of the effect of amino acid changes on the binding affinity between SARS-CoV-2 spike RBD and human ACE2. *Proc. Natl. Acad. Sci. U.S.A.* **118**, e2106480118 (2021).
32. X. Wei, J. Zhang, Patterns and mechanisms of diminishing returns from beneficial mutations. *Mol. Biol. Evol.* **36**, 1008–1021 (2019).
33. H. H. Chou, H. C. Chiu, N. F. Delaney, D. Segrè, C. J. Marx, Diminishing returns epistasis among beneficial mutations decelerates adaptation. *Science* **332**, 1190–1192 (2011).
34. S. Kryazhimskiy, D. P. Rice, E. R. Jerison, M. M. Desai, Global epistasis makes adaptation predictable despite sequence-level stochasticity. *Science* **344**, 1519–1522 (2014).
35. A. W. R. Serohijos, E. I. Shakhnovich, Contribution of selection for protein folding stability in shaping the patterns of polymorphisms in coding regions. *Mol. Biol. Evol.* **31**, 165–176 (2014).
36. S. Di Giacomo, D. Mercatelli, A. Rakhimov, F. M. Giorgi, Preliminary report on severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) spike mutation T478K. *J. Med. Virol.* **93**, 5638–5643 (2021).
37. J. A. Plante *et al.*, Spike mutation D614G alters SARS-CoV-2 fitness. *Nature* **592**, 116–121 (2021).
38. Y. J. Hou *et al.*, SARS-CoV-2 D614G variant exhibits efficient replication ex vivo and transmission in vivo. *Science* **370**, 1464–1468 (2020).
39. Z. Zhao *et al.*, Omicron SARS-CoV-2 mutations stabilize spike up-RBD conformation and lead to a non-RBM-binding monoclonal antibody escape. *Nat. Commun.* **13**, 4958 (2022).
40. V. Stalls *et al.*, Cryo-EM structures of SARS-CoV-2 omicron BA.2 spike. *Cell Rep.* **39**, 111009 (2022).
41. F. Touret *et al.*, Replicative fitness of a SARS-CoV-2 20I/501Y.V1 variant from lineage B.1.1.7 in human reconstituted bronchial epithelium. *mBio* **12**, e0085021 (2021).
42. R. Burioni, E. J. Topol, Has SARS-CoV-2 reached peak fitness? *Nat. Med.* **27**, 1323–1324 (2021).
43. N. Matange, S. Hegde, S. Bodkhe, Adaptation through lifestyle switching sculpts the fitness landscape of evolving populations: Implications for the selection of drug-resistant bacteria at low drug pressures. *Genetics* **211**, 1029–1044 (2019).
44. L. Montanucci, E. Capriotti, Y. Frank, N. Ben-Tal, P. Fariselli, DDGun: An untrained method for the prediction of protein stability changes upon single and multiple point variations. *BMC Bioinformatics* **20**, 335 (2019).
45. N. Tokuriki, F. Stricher, J. Schymkowitz, L. Serrano, D. S. Tawfik, The stability effects of protein mutations appear to be universally distributed. *J. Mol. Biol.* **369**, 1318–1332 (2007).
46. T. N. Starr *et al.*, Prospective mapping of viral mutations that escape antibodies used to treat COVID-19. *Science* **371**, 850–854 (2021).
47. K. Katoh, D. M. Standley, MAFFT multiple sequence alignment software version 7: Improvements in performance and usability. *Mol. Biol. Evol.* **30**, 772–780 (2013).
48. S. Elbe, G. Buckland-Merrett, Data, disease and diplomacy: GISaid's innovative contribution to global health. *Glob. Chall.* **1**, 33–46 (2017).
49. B. Hie, E. D. Zhong, B. Berger, B. Bryson, Learning the language of viral evolution and escape. *Science* **371**, 284–288 (2021).
50. J. Han *et al.*, D3AI-Spike: A deep learning platform for predicting binding affinity between SARS-CoV-2 spike receptor binding domain with multiple amino acid mutations and human angiotensin-converting enzyme 2. *Comput. Biol. Med.* **151**, 106212 (2022).
51. A. Rives *et al.*, Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences. *Proc. Natl. Acad. Sci. U.S.A.* **118**, e2016239118 (2019).
52. D. Wang *et al.*, COVID-19 Biophysical Model. GitHub. <https://github.com/Dianzhuo-Wang/COVID19-Biophysical-Model>. Deposited 16 May 2024.
53. F. Obermeyer *et al.*, Pyro-Cov. GitHub. <https://github.com/broadinstitute/pyro-cov>. Accessed 1 May 2023.
54. A. Moulana *et al.*, Compensatory Epistasis Omicron. GitHub. https://github.com/desai-lab/compensatory_epistasis_omicron. Accessed 1 May 2023.
55. A. Moulana *et al.*, Omicron Antibody Landscape. GitHub. https://github.com/desai-lab/omicron_ab_landscape. Accessed 1 May 2023.