

1 Vaccination shapes evolutionary trajectories of SARS-CoV-2

2 Matthijs Meijers^a, Denis Ruchnewitz^a, Marta Luksza^b, Michael Lässig^{a,*}

3 ^a Institute for Biological Physics, University of Cologne, Zülpicherstr. 77, 50937 Köln, Germany

4 ^b Tisch Cancer Institute, Departments of Oncological Sciences and Genetics and Genomic Sciences, Icahn

5 School of Medicine at Mount Sinai, New York, NY, USA

6 * To whom correspondence should be addressed. Email: mlaessig@uni-koeln.de

7 Abstract

8 The large-scale evolution of the SARS-CoV-2 virus has been marked by rapid turnover of genetic
9 clades. New variants show intrinsic changes, notably increased transmissibility, as well as anti-
10 genic changes that reduce the cross-immunity induced by previous infections or vaccinations^{1–4}.
11 How this functional variation shapes the global evolutionary dynamics has remained unclear.
12 Here we show that selection induced by vaccination impacts on the recent antigenic evolution
13 of SARS-CoV-2; other relevant forces include intrinsic selection and antigenic selection induced
14 by previous infections. We obtain these results from a fitness model with intrinsic and antigenic
15 fitness components. To infer model parameters, we combine time-resolved sequence data⁵, epi-
16 demiological records^{6,7}, and cross-neutralisation assays^{8–10}. This model accurately captures the
17 large-scale evolutionary dynamics of SARS-CoV-2 in multiple geographical regions. In partic-
18 ular, it quantifies how recent vaccinations and infections affect the speed of frequency shifts
19 between viral variants. Our results show that timely neutralisation data can be harvested to
20 identify hotspots of antigenic selection and to predict the impact of vaccination on viral evolu-
21 tion.

22 Introduction

23 Two classes of molecular adaptation have been observed in the evolution of SARS-CoV-2
24 to date. Multiple mutations carry intrinsic changes of viral functions, such as increasing the
25 binding affinity to human receptors¹, the efficiency of cell entry^{2,3}, or the stability of viral
26 proteins^{11,12}. Other mutations, referred to as antigenic changes, decrease the neutralizing ac-
27 tivity of human antibodies^{4,8–10}, thereby reducing the immune protection against secondary
28 infections^{13,14}. The strains that inherit a given mutation define a clade of the evolving viral
29 population. Several of these molecular changes had drastic evolutionary and epidemiological
30 impact, inducing global turnover of viral clades and concurrent waves of the pandemic. Over
31 the last two years, three genetic variants and their associated clades successively gained global
32 prevalence: Alpha (α) from March to June in 2021, Delta (δ) from June to December in 2021
33 and Omicron (ω) in 2022. These were named Variants of Concern (VOCs) by the World Health
34 Organization¹⁵; other VOCs gained temporary regional prevalence. Several studies reported fit-
35 ness advantages of VOCs inferred from epidemiological trajectories and comparative functional
36 studies^{3,16–19}. Importantly, however, the evolutionary impact of antigenic changes is time-
37 dependent, because it depends on previously acquired population immunity: a larger amount of
38 previous infections or vaccinations increases the global fitness advantage of an antigenic escape
39 mutation. Specifically, multi-strain epidemiological models and simulations suggest that vacci-
40 nations can favour the emergence of escape variants^{20–23} and influence the turnover of circulating
41 clades^{24,25}; effects of this kind have been reported for some clades of human influenza²⁶. In
42 the case of SARS-CoV-2, pandemic infection and massive vaccination programs, with a global
43 count of 4.5 billion vaccinations and >200 million confirmed cases in 2021⁶, have built up partial
44 population immunity, but its feedback on viral evolution has not been quantified. This leads to

45 the central question of this paper: what is the impact of vaccination and infection rates on the
46 turnover of SARS-CoV-2 clades? To address this question, we infer a data-driven fitness model
47 for SARS-CoV-2 variants with distinct components of intrinsic fitness and antigenic fitness by
48 vaccination and infection.

49 Results

50 **Trajectories and speed of clade turnover** As a first step, we map the evolutionary trajectories
51 of the three global clade shifts in the last two years. To track circulating clades, we analyse
52 a set of >5M quality-controlled SARS-CoV-2 sequences obtained from the GISAID database⁵.
53 We assign these sequences to genetic clades using a standard set of amino acid changes²⁷; then
54 we infer time-dependent clade frequencies from strain counts smoothed over a period of ~30
55 days (Methods). To obtain accurate, time-resolved data, we record frequency trajectories at the
56 level of regions (countries and US states). Including all regions satisfying uniform criteria of
57 data availability (Methods), we obtain frequency trajectories of the $1 - \alpha$, the $\alpha - \delta$, and the $\delta - o$
58 shift for 11, 16, and 14 regions, respectively. Here, 1 denotes the set of clades circulating prior
59 to α , including the wild type (wt) and the early 614G mutation in the spike protein. Fig. 1a
60 shows trajectories of the ancestral and the invading clade for the $\alpha - \delta$ and $\delta - o$ shifts in Italy;
61 trajectories for all regions of this study are reported in Fig. S1 and S2.

62 Assuming that large-scale frequency shifts of viral clades are adaptive processes, we can
63 infer the underlying selective force from frequency trajectories. Specifically, the fitness difference
64 (selection coefficient) between invading and ancestral clades takes the form

$$\hat{s}(t) = \frac{d}{dt} \log \frac{x_{\text{inv}}(t)}{x_{\text{anc}}(t)}, \quad (1)$$

65 where $x_{\text{inv}}(t)$ and $x_{\text{anc}}(t)$ are the corresponding frequencies (here and below, empirical selection
66 coefficients inferred from frequency trajectories are marked by a hat). We note that this relation
67 is independent of other co-circulating clades (Methods). In Fig. 1b, we show time-resolved,
68 regional selection coefficients of the invading and ancestral clade for the $\alpha - \delta$ and $\delta - o$ shifts.
69 These data reveal two opposing trends: During the $\alpha - \delta$ shift, selection increases with time
70 in 16 of 16 regions. Conversely, selection driving the $\delta - o$ shift decreases with time in 12 of
71 14 regions. Compared to a reference of time-independent selection, the $\alpha - \delta$ shift runs at an
72 accelerating speed, the $\delta - o$ shift at a decelerating speed. The time dependence of selection
73 is statistically significant ($P < 10^{-15}$ for $\alpha - \delta$, $P < 10^{-5}$ for $\delta - o$; two-sided Wald test). In
74 contrast, the earlier $1 - \alpha$ shift does not show a significant signal of time-dependent selection
75 ($P > 0.01$, Fig. S3). In what follows, we will relate this pattern to feedback of vaccination on
76 viral evolution.

77 **Cross-immunity trajectories** Cross-immunity induced by a primary infection against subse-
78 quent infections by related pathogens is routinely tested by neutralisation assays, which measure
79 the minimum antiserum concentration required to neutralise the second antigen. Relative, in-
80 verse concentrations are reported as serum dilution titers; here we use logarithmic titer values,
81 T (with base 2). For SARS-CoV-2, recent work^{8–10,28,29} has established a matrix of titers, T_i^k ,
82 measuring neutralisation of variant i in immune channel k (Fig. 2a, Table S1). Here and below,
83 immune channels label primary challenges inducing specific antisera, including infections by dif-
84 ferent variants ($k = \alpha, \delta, o, \dots$), as well as primary and booster vaccinations ($k = \text{vac}, \text{bst}$; titers
85 shown here are for mRNA vaccines). Together, these data provide a first, coarse-grained cross-
86 immunity landscape of SARS-CoV-2. Infection-induced cross-immunity titers are maximal when
87 primary infection and secondary challenge are by the same variant (Fig. 2a). Similarly, titers in-
88 duced by primary vaccination, T_i^{vac} , are maximal against strains from the ancestral clade, which
89 contains the strain used for vaccination³⁰. Differences of neutralisation titers, $\Delta T_{ij}^k = T_i^k - T_j^k$,

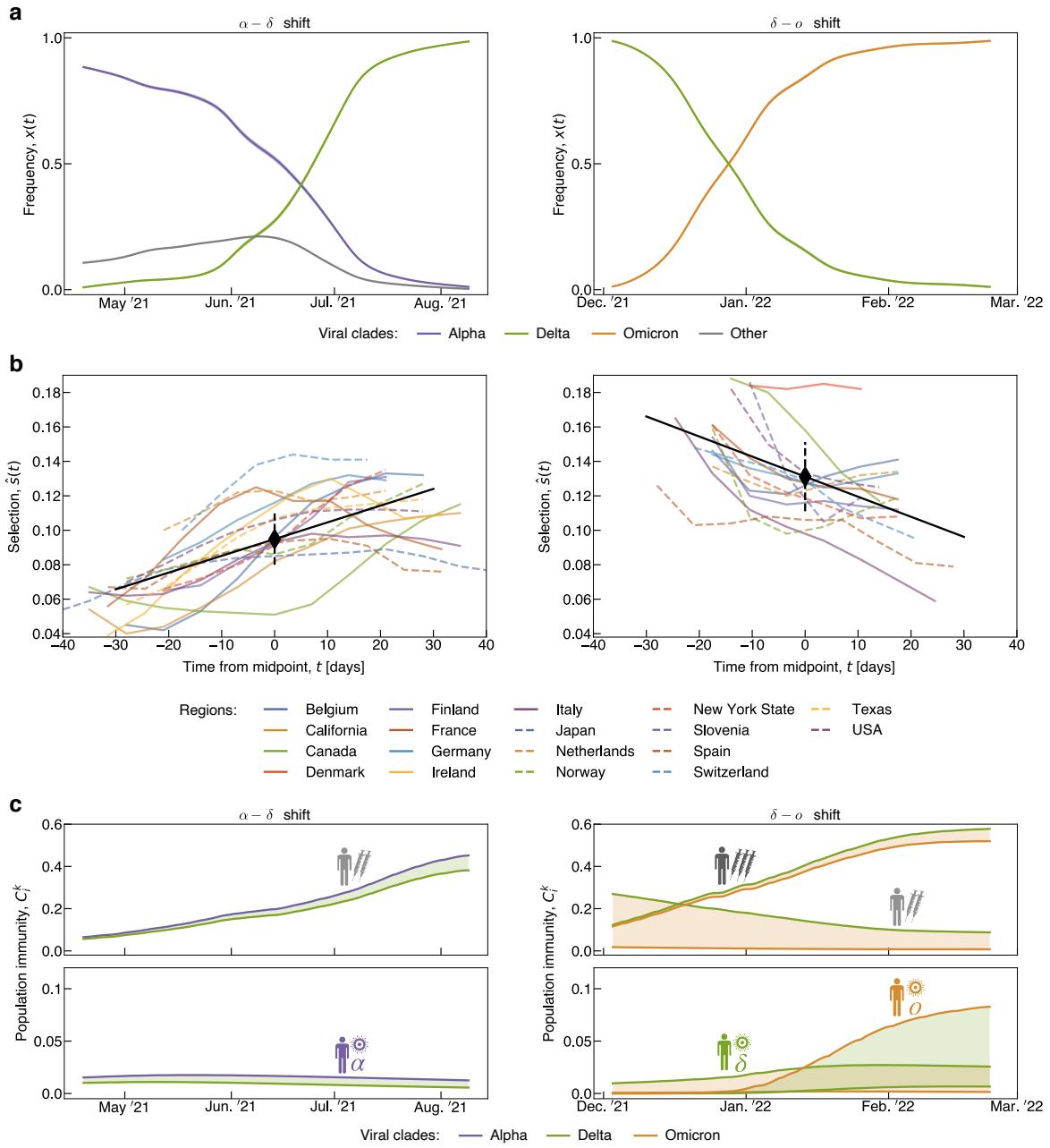


Fig. 1: Evolutionary, epidemiological, and immune tracking of SARS-CoV-2. Time-dependent trajectories are shown for the clade shift from α to δ (left column) and from δ to o (right column). (a) Observed frequency trajectories of relevant clades, $x_i(t)$, for the clade shifts in Italy. (b) Empirical selection coefficient (fitness difference) between invading and ancestral clade, $\hat{s}(t)$, for all regions. Selection trajectories are derived from the frequency trajectories of (a) and plotted against time counted from the midpoint. Summary statistics: cross-region linear regression (black solid line), cross-region average (black diamond), and rms cross-region variation of selection (black dashed line). (c) Population immunity functions of the ancestral and invading variant, $C_{\text{anc}}^k(t)$ and $C_{\text{inv}}^k(t)$, in relevant immune channels k for Italy (coloured lines). Cross-immunity differences in a given channel, $C_{\text{inv}}^k(t) - C_{\text{anc}}^k(t)$, are highlighted by shading (colours indicate which variant receives a fitness advantage). See Figs. S1–S3 for tracking of all shifts in all regions and reporting of rms statistical errors.

measure differences in functional antibody binding between strains of different variants; evolved titer reductions are also referred to as antigenic advance. Notably, each of the global clade shifts observed to date, $1 - \alpha$, $\alpha - \delta$, and $\delta - o$, has decreased neutralisation by vaccination, i.e., generated antigenic advance, $\Delta T_{1\alpha}^{\text{vac}}, \Delta T_{\alpha\delta}^{\text{vac}}, \Delta T_{\delta o}^{\text{vac}} > 0$ (Fig. 2b). Moreover, the in-vivo concentration of neutralising antibodies decays exponentially with time after immunisation^{31,32}. This

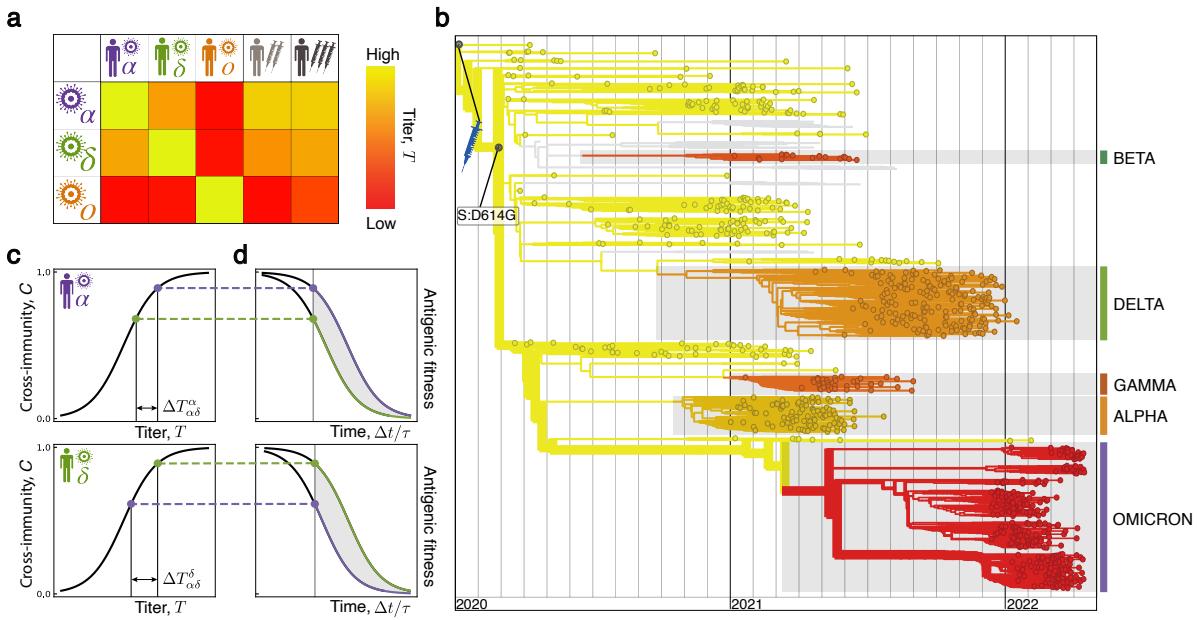


Fig. 2: Cross-neutralisation and antigenic fitness. (a) Neutralisation titers, T_i^k , of human antisera induced by different primary challenges (columns: infection by strains from clade α , δ , o , mRNA primary and booster vaccination) assayed against different test strains (rows: strains from clade α , δ , o); see Table S1. (b) Strain tree of SARS-CoV-2 with lineages i colored by vaccine neutralisation titers, T_i^{vac} . WHO Variants of Concern are marked by bars. The ancestral clade (1) has the highest neutralisation (yellow); the successive clade shifts $1 - \alpha$, $\alpha - \delta$, and $\delta - o$ decrease neutralisation, i.e., induce antigenic advance (see text). (c) The cross-immunity, c_i^k , induced by a primary immunisation in channel k (top: infection by α , bottom: infection by δ) against a secondary infection of clade i (blue dot: α , red dot: δ) is a Hill function of the neutralising titer, T_i^k (ref. [13, 14], Methods). Cross-immunity decreases with increasing antigenic advance ΔT_{ik}^k (bars). (d) Cross-immunity decays with time after primary immunisation, Δt (in units of the characteristic decay time τ ; ref. [31, 32]). According to the fitness model, cross-immunity induces a proportional antigenic fitness cost; the resulting time-dependent selection coefficient (fitness difference) between clades is marked by shading.

95 translates into a linear titer reduction, $T_i^k(\Delta t) = T_i^k - \Delta t/\tau$, with an estimated half life $\tau \sim 65$
96 days (Methods).

97 Importantly, recent work for SARS-CoV-2 has also shown that neutralisation titers predict
98 the cross-immunity c_i^k , defined as the relative drop of secondary infections in human cohorts.
99 Specifically, $c_i^k = H(T_i^k)$ is a Hill function^{13, 14} (Fig. 2c, details are given in Methods), consistent
100 with the underlying biophysics of antibody-antigen binding and with results for other viral
101 pathogens^{33–36}. The post-immunisation decay of antibody concentration induces a decay of
102 cross-immunity, $c_i^k(\Delta t) = H(T_i^k - \Delta t/\tau)$. Together, cross-immunity depends in a predictable,
103 nonlinear way on neutralisation titer and on time since primary immunisation. Fig. 2 shows two
104 examples of this pattern. Primary infection by an α strain induces a high cross-immunity against
105 other α strains and a reduced cross-immunity against δ strains ($c_\alpha^\alpha > c_\delta^\alpha$) (Fig. 2c, top). Both
106 factors decrease by antibody decay; their difference has a maximum at an intermediate time
107 since primary infection (Fig. 2d, top). Infection by a δ strain induces cross-immunity factors of
108 opposite ranking ($c_\delta^\alpha < c_\delta^\delta$) and similar decay (Fig. 2cd, bottom).

109 To track population immunity over time, we combine these cross-immunity factors with
110 infection and vaccination data. In each region, we record cumulative fractions of immunised
111 individuals, $y_k(t)$, in each channel k (clade-specific infections, primary and booster vaccinations;
112 see Figs. S1 and S2). Their derivatives $\dot{y}_k(t)$ are the rates of new immunisations in channel k .
113 Clade-specific infection data are obtained by multiplying the total rate of new infections reported
114 in each region with the simultaneous viral clade frequencies $x_k(t)$ (Fig. 1a). In the regions
115 included in our analysis, vaccination has been predominantly by mRNA vaccines (Methods).
116 By weighting with the time-dependent cross-immunity factors $c_i^k(\Delta t)$, we infer the population

117 cross-immunity against clade i by immunisation in channel k ,

$$C_i^k(t) = \int^t c_i^k(t-t') \dot{y}_k(t') dt'. \quad (2)$$

118 In Fig. 1c, we plot the cross-immunity trajectories relevant for the $\alpha - \delta$ and $\delta - o$ shifts in Italy; 119 trajectories for all regions are reported in Figs. S1 and S2. The $\alpha - \delta$ shift shows sizeable and 120 increasing immunity induced by primary vaccination, while infection-induced immunity remains 121 small. During the $\delta - o$ shift, immunity by primary vaccination declines, while booster- and 122 infection-induced immunity components increase. During the earlier $1 - \alpha$ shift, population 123 cross-immunity is still small in most regions. We conclude that the joint dynamics of new 124 immunisations and antibody decay can produce complex and opposing cross-immunity patterns.

125 **Inference of intrinsic and antigenic selection** To quantify the feedback of cross-immunity 126 on viral evolution, we use a minimal, computable fitness model,

$$f_i(t) = f_i^0 - \sum_k \gamma_k C_i^k(t), \quad (3)$$

127 where $f_i(t)$ is the absolute fitness, or epidemic growth rate, of a viral strain. Fitness is propor- 128 tional to the log of the effective reproductive number, $f_i(t) = \tau_0^{-1} \log R_i(t)$, where τ_0 denotes 129 the infectious period (Methods). Here, we write fitness as the sum of a time-independent in- 130 trinsic component, f_i^0 , and of time-dependent antigenic components, $f_i^k(t) = -\gamma_k C_i^k(t)$ (Meth- 131 ods). Each component is proportional to the corresponding cross-immunity factor $C_i^k(t)$ with a 132 weight factor γ_k for each immune channel k . Hence, selection is generated by cross-immunity 133 differences between competing strains (shading in Fig. 1c and Fig. 2c). This type of fitness 134 model has been established for predictive evolutionary analysis of human influenza^{24,37,38} and 135 is grounded in multi-strain epidemiological models³⁹. The minimal fitness model does not ac- 136 count for differences in cross-immunity between human hosts (for example, through differences 137 in immunodominance⁴⁰) and for correlations between multiple prior infections (antigenic sin⁴¹).

138 For SARS-CoV-2, we compute fitness at the level of variants, neglecting fitness differences be- 139 tween strains within a clade. Similarly, we evaluate cross-immunity at the level of variant-specific 140 prior infection and of primary and booster vaccination, using the trajectories $C_i^k(t)$ calculated 141 above (Fig. 1c). To compare model and data, we compute the fitness difference between invading 142 and ancestral strain for each regional trajectory: $s(t) = f_{\text{inv}}(t) - f_{\text{anc}}(t) = s_0 + s_{\text{ag}}(t)$, where in- 143 trinsic selection, s_0 , and antigenic selection, $s_{\text{ag}}(t) = \sum_k s_k(t)$, are given by equation (3). Then 144 we decompose the model-based selection trajectories into mean and change, $s(t) = \langle s \rangle + \Delta s(t)$ 145 (brackets denote time averages over the trajectory for a given region). The empirical trajectories 146 $\hat{s}(t)$ are decomposed in the same way (Fig. 1b). Cross-region selection differences, measured by 147 the rms deviation of $\langle \hat{s} \rangle$, reflect inhomogeneous conditions of contact limitations, surveillance, 148 geography, and population structure that are not included in the minimal model. In a given 149 region, however, variants compete under more homogeneous conditions. Therefore, we infer 150 antigenic selection from the regional selection change, $\Delta \hat{s}(t)$. We use a minimal model with just 151 3 antigenic parameters: a uniform γ_{vac} for vaccination and boosting (downweighted by a factor 152 a in the $\delta - o$ shift to account for double infections⁴²) and a uniform $\gamma_k = b \gamma_{\text{vac}}$ for all infection 153 channels k (upweighted by a factor b to correct for relative underreporting; Methods). We infer 154 maximum-likelihood (ML) values of these parameters by calibrating computed and empirical 155 trajectories, $\Delta s(t)$ and $\Delta \hat{s}(t)$, for the $\alpha - \delta$ and $\delta - o$ shifts. The intrinsic selection coefficients, 156 s_0 , are then obtained as the time-independent part of selection. Details of the inference proce- 157 dure are given in Methods; ML model parameters and selection coefficients for all clade shifts 158 are reported in Table S2 and S3.

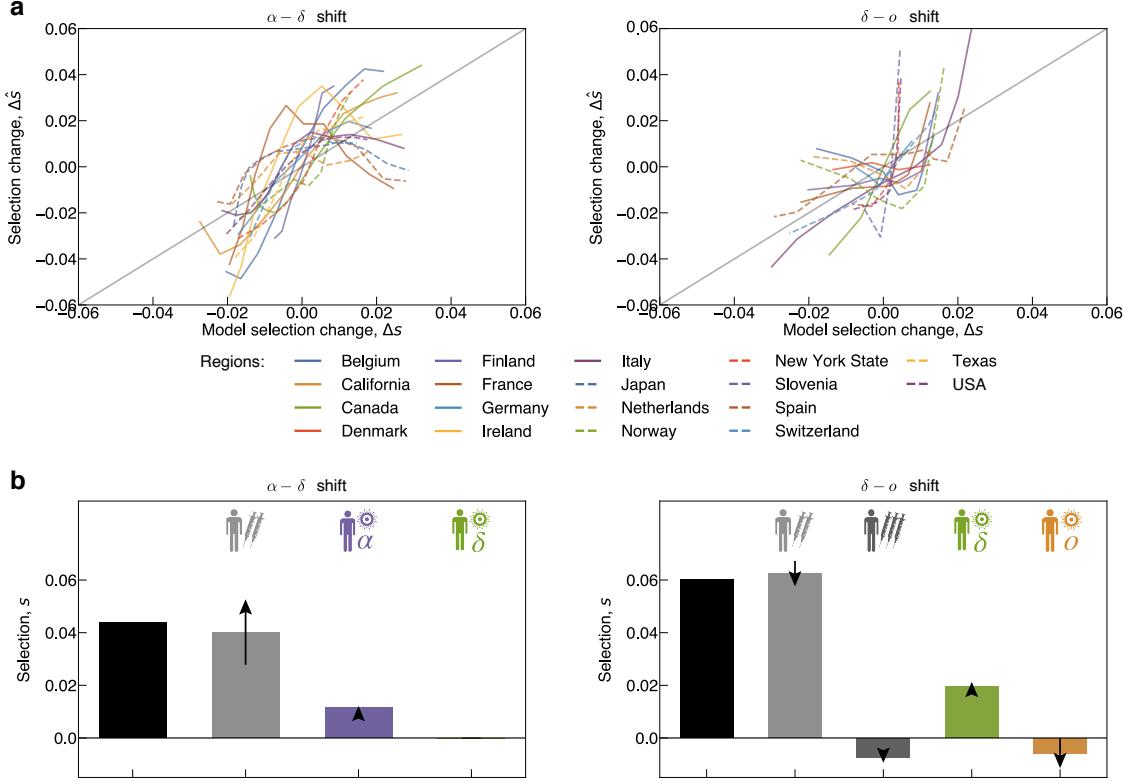


Fig. 3: Antigenic and intrinsic selection drive SARS-CoV-2 evolution. We compare empirical selection trajectories and model predictions for the clade shift from α to δ (left column) and from δ to o (right column). **(a)** Empirical selection change, $\Delta\hat{s}$, obtained from the selection trajectories of Fig. 1a are plotted against model predictions, Δs , for all regions. Rms statistical errors are reported in Figs. S1–S2. **(b)** Breakdown of fitness model components. Intrinsic selection coefficients (black) and antigenic selection coefficients in marked immune channels (coloured), as inferred from the ML fitness model (bars: region- and time-averaged value for each crossover; arrows: region-averaged rms temporal change, $\langle(\Delta s)^2\rangle^{1/2}$, with marked direction; confidence intervals are given in Table S3).

159 In Fig. 3a, we plot Δs from the ML fitness model against the corresponding empirical
 160 selection change, $\Delta\hat{s}$, obtained from the trajectories of Fig. 1b. We obtain a remarkable data
 161 compression: for most regions, the antigenic fitness computed from equation (3) reproduces
 162 the empirical fitness changes (intrinsic selection drops out of this comparison). This can be
 163 further quantified: the covariance between data and ML model, $\langle\Delta s\Delta\hat{s}\rangle$, explains $\sim 50\%$ of
 164 the empirical variance of selection, $\langle(\Delta s)^2\rangle$; this level of covariance is found on average and
 165 in most individual regions. A detailed comparison of data and model trajectories, $\Delta\hat{s}(t)$ and
 166 $\Delta s(t)$, for all regions is shown in Figs. S1 and S2. As a control, the model predicts only small
 167 selection change for the $1 - \alpha$ shift, consistent with the weak time dependence of the empirical
 168 selection trajectories (Fig. S3). We conclude that time-dependent cross-immunity explains the
 169 time-dependence of selection governing SARS-CoV-2 variant shifts.

170 **Impact of vaccination and infection on evolution** From the ML fitness model, we obtain a
 171 breakdown of intrinsic and antigenic selection components relevant for each clade shift. Intrinsic
 172 selection is strong and positive in all three major clade shifts, with average selection coefficients
 173 $s_0 = 0.05 - 0.08$, consistent with strong functional differences observed between the α , δ , and o
 174 variants^{3,43} (Fig. 3b, Table S3). Antigenic selection becomes equally strong in the $\alpha - \delta$ and $\delta - o$
 175 shifts. Its two main components, vaccination- and infection-induced selection, are statistically
 176 significant parts of the fitness model, partial models with only one component have a strongly

177 reduced posterior likelihood (differences in model complexity are accounted for by a Bayesian
178 information criterion; see Methods and Table S2).

179 Vaccination induces cross-immunity differences between variants, resulting in positive anti-
180 genic selection of average strength $s_{\text{vac}} = 0.04$ in the $\alpha - \delta$ shift and $s_{\text{vac}} = 0.06$ in the
181 $\delta - o$ shift (Fig. 3b, Table S3). These selection coefficients quantify the evolutionary impact
182 of primary SARS-CoV-2 vaccination: they measure the relative increase in effective repro-
183 duction number of the invading variant by partial escape from vaccination-induced immunity
184 ($\tau_0 s_{\text{vac}} = R_{\text{inv}}/R_{\text{anc}} - 1$). Vaccination-induced antigenicity also explains the observed time-
185 dependence of selection (Fig. 1b, Figs. S1 and S2): s_{vac} increases during the $\alpha - \delta$ shift be-
186 cause of increasing vaccination levels, but decreases during the $\delta - o$ shift because vaccination-
187 induced immunity fades. In both shifts, primary vaccination generates the dominant compo-
188 nents of antigenic selection (Fig. 3b). Booster vaccinations have increased breadth; they induce
189 higher neutralisation $T_{\delta}^{\text{bst}}, T_o^{\text{bst}}$ and reduced antigenic advance $\Delta T_{\alpha\delta}^{\text{bst}}, \Delta T_{\delta o}^{\text{bst}}$ compared to pri-
190 mary vaccinations^{9,10,44,45} (Fig. 2a, Table S1). Hence, booster vaccinations generate higher
191 cross-immunity but weaker selection for antigenic escape (Fig. 1c, Fig. 4ab). The net effect
192 of boosters in the $\delta - o$ shift is opposite to that of primary vaccinations: we infer a negative
193 selection coefficient $s_{\text{bst}} = -0.01$. This is because boosters remove cross-immunity differences
194 and antigenic selection generated by the preceding primary vaccination (Fig. 3b).

195 Infection-induced antigenic selection increased in net strength from 0.01 in the $\alpha - \delta$ shift
196 to 0.03 in the $\delta - o$ shift. Notably, it always contains components of opposite sign: primary
197 infections by the ancestral clade generate positive selection, while infections by the invading clade
198 generate negative selection. This frequency-dependent negative feedback acts to prolong the
199 coexistence of ancestral and invading clade. Together, antigenic selection can produce complex
200 but computable patterns of time dependence.

201 These results require careful interpretation. They show that vaccination and previous infec-
202 tions induced sizeable antigenic selection on circulating SARS-CoV-2 variants and modulated
203 the speed of successive clade shifts. However, antigenic selection did not cause or prevent any
204 of these shifts, because intrinsic functional changes generated sizeable fitness advantages of the
205 invading variants independently of population immunity. The breakdown of selection given in
206 Fig. 3b applies to the set of regions accessible to our analysis; the relative weights of vaccination-
207 and infection-induced selection components are expected to be different in other regions. The
208 availability of comparable data precludes a fully global model-based analysis. An additional,
209 model-free inference of selection in regions with low vaccination coverage is given in Methods.
210 Most importantly, the fitness model and our data analysis do not predict any simple relation
211 between vaccination coverage and speed of evolution. This is because cross-immunity channels
212 are correlated: fewer vaccinations lead to more infections, generating buildup of cross-immunity
213 in other channels and complex long-term effects.

214 **Fitness trajectories and selection hotspots** The ML fitness model can be applied to the
215 long-term turnover of viral clades up to date, including recent frequency changes between the
216 variants BA.1, BA.2, and BA.4/5 within the o clade (we use shorthands $o1$, $o2$, and $o45$). First,
217 we look at two building blocks of antigenic fitness: antigenic landscapes and immune weights.
218 We define antigenic landscapes for each immune channel k by plotting all cross-immunity factors
219 c_i^k against their corresponding titers T_i^k (using a fixed time delay $\Delta t = \tau$ to account for antibody
220 decay in an approximate way; Methods). These landscapes visualise antigenic drift, that is, the
221 partial escape from population immunity by gradual evolutionary steps^{46,47} (Fig. 4a-c; arrows
222 mark sizeable steps between successive variants). In this picture, the time-dependence of cross-
223 immunity is captured by immune weight functions $Q_k(t)$, which measure recent infections or
224 vaccinations in channel k , again over a time window of order τ (Fig. 4d, Methods). Next, we
225 juxtapose these immune trajectories to long-term trajectories showing the antigenic selection
226 between successive variants, $s_{\text{ag}}(t)$ (Fig. 4e), and the fitness gap of each variant, $\delta f_i(t) = f_i(t) -$

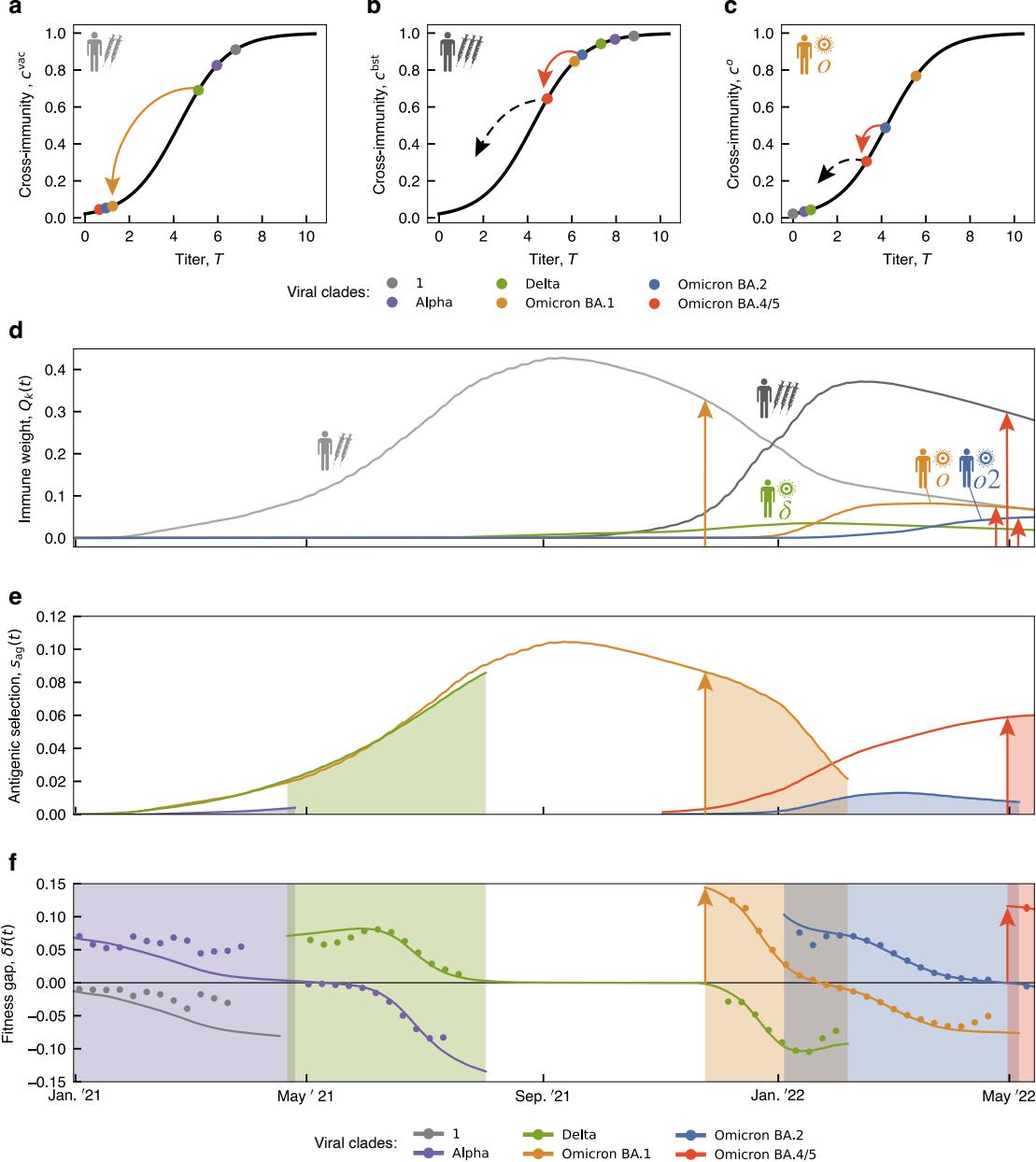


Fig. 4: Antigenic landscapes and immune weights generate selection hotspots. (a) Antigenic landscapes. Cross-immunity factors of major and recent clades (colored dots) are plotted against neutralisation titers in different immune channels: (a) primary vaccination, (b) booster vaccination, (c) α -induced infection. (d) Long-term immune weight trajectories of different channels, $Q_k(t)$. (e, f) Long-term fitness trajectories of major and recent clades. Periods of clade shifts are highlighted by shading. (e) Antigenic selection between successive variants, $s_{\text{ag}}(t)$. Model-based trajectories for each variant pair are shown up to the end of the corresponding clade shift. (f) Time-dependent fitness gap, $\delta f(t)$ (Methods). Model-based trajectories for each variant i (lines) are shown from in the time interval where $0.01 \leq x_i(t) \leq 0.99$; empirical selection is marked by dots. Selection hotspots: when sizeable cross-immunity drop on the flank of an antigenic landscape (arrows in a-c) coincides with large immune weights (arrows in d), the fitness model predicts time windows of strong selection for antigenic escape (arrows in e, f mark clade shifts starting in a selection hotspot). Immune weight and fitness trajectories are averaged over regions (see Fig. S4 for regional trajectories).

227 $\bar{f}(t)$ (Fig. 4f); these trajectories are computed from equation (3). Fitness gaps are shifted by the
 228 mean population fitness $\bar{f}(t) = \sum_i x_i(t)f_i(t)$ and include the intrinsic component (Methods).
 229 As expected from the analysis above, the ML model is in quantitative agreement with empirical
 230 selection (dots in Fig. 4f). The trajectories of Fig. 4ef are averaged over 14 regions (for regional

231 trajectories, see Fig. S4).

232 Together, the trajectories of Fig. 4 show a pattern of selection hotspots. The fitness model
233 predicts time windows of strong antigenic selection when antigenic advance on the flank of a Hill
234 landscape generates sizeable cross-immunity loss and coincides with high immune weight. We
235 now trace this pattern through successive clade shifts. At early stages of evolution, until spring
236 2021, all immune weights were small (Fig. 4d). Hence, intrinsic selection governed the $1 - \alpha$
237 shift; consistently, the antigenic advance (neutralisation titer drop) $\Delta T_{1\alpha}^{\text{vac}}$ was small (Fig. 4a).
238 Between spring 2021 and spring 2022, primary and booster vaccination generated the dominant
239 immune weights and induced directional antigenic selection for escape from the vaccine strain,
240 while infection-induced immunity remained relatively small (Fig. 4d). The clade shifts $\alpha - \delta$ and
241 $\delta - o$ carried increasing antigenic advance $\Delta T_{\alpha\delta}^{\text{vac}}$ and $\Delta T_{\delta o}^{\text{vac}}$, and smaller advance with respect to
242 boosting (Fig. 4ab). These changes mark the onset of antigenic drift. The fitness model identifies
243 a first clear hotspot of antigenic selection driving the $\delta - o$ shift. At the start of this shift, large
244 cross-immunity change and large immune weight coincided in the primary vaccination channel
245 (orange arrows in Fig. 4ade). This is consistent with the observed dynamics: $\delta - o$ was faster
246 than the previous shifts and under exceptionally high initial selection, $\hat{s} = 0.15$ (orange arrow in
247 Fig. 4f, Fig. 1b). The following shift, $o1 - o2$, involved subclades with similar neutralisation by
248 primary and booster vaccination. This shift is inferred to be governed predominantly by intrinsic
249 selection (Table S1 and S3); consistently, selection is only weakly time-dependent (Fig. S3).

250 The most recent viral-immune co-evolution shows two important novelties. In spring 2022,
251 infection immunity increased, while vaccination-induced immunity decreased; both components
252 are reaching comparable weights (Fig. 4d). Whether vaccination remains at sizeable immune
253 weight will depend on availability and acceptance of vaccines in the future. These immune weight
254 changes also mark the onset of immune drift, that is, the response of population immunity to
255 antigenic drift of the viral population. Recently, population immunity has shifted its center of
256 mass from wt towards o ; components cognate to each of these clades have reached comparable
257 weight (red arrows in Fig. 4d).

258 At this point, our fitness model predicts the next selection hotspot for novel variants car-
259 rying antigenic advance away from vaccination *and* from o infections. The recent antigenic
260 evolution within the o clade follows this scenario: the emerging variants $o4$ and $o5$ (BA.4 and
261 BA.5) combine antigenic advance in three channels^{45,48} (red arrows in Fig. 4bce). Consistently,
262 these variants show fast initial growth with empirical selection $\hat{s} = 0.12$ (red arrow in Fig. 4f).
263 Moreover, near-future mutations carrying antigenic advance in the same channels are predicted
264 to be in the same hotspot (dashed arrows in Fig. 4bc).

265 If these emerging variants develop into major clade shifts, they will further increase the o
266 immune weight factors. On the other hand, population immunity against earlier variants could
267 be maintained by backboosting⁴⁹ of o infections or by bivalent vaccines with a wt component.
268 The future evolutionary trajectories of new variants will also depend on their mutual antigenic
269 relations, which have not yet been assayed comprehensively to date. Together, the recent evolu-
270 tionary dynamics signals the unfolding of antigenic complexity towards coexistence of multiple
271 antigenic variants and immune classes.

272 Discussion

273 Here we have established a data-driven, multi-component fitness model for the evolution of
274 SARS-CoV-2. By applying this model to recent evolutionary trajectories in multiple regions,
275 we have quantified intrinsic and antigenic selection driving the genetic and functional evolution
276 of the virus. In particular, primary vaccination impacted on the speed of global clade shifts in
277 2021. Booster vaccination generated higher cross-protection, but weaker selection for antigenic
278 escape in the same period (Fig. 3). These results underscore that vaccine breadth is important
279 for constraining antigenic escape evolution. More broadly, they highlight the need to integrate

280 evolutionary feedback into vaccine design.

281 In the recent evolution of SARS-CoV-2, two general trends are revealed by our analysis.
282 Antigenic selection has increased in strength and has broadened its target: primary infection
283 by distinct viral variants has generated an increasing number of antigenic selection components
284 (Fig. 3b, Fig. 4). These trends mark the transition from initial, post-zoonotic adaptation of
285 the virus to evolution to an endemic state, where antigenic evolution continues to be fuelled by
286 the buildup of population immunity to circulating viral variants. A plausible end point of this
287 transition becomes clear by comparison with influenza, a long-term endemic virus in humans.
288 In influenza, the viral escape from population immunity follows a specific mode of antigenic
289 drift, where multiple variants with different cross-immunity profiles compete for prevalence⁵⁰.
290 This mode is marked by continuous, adaptive clade turnover with characteristic time scales
291 of several months, which is substantially slower than the recent prevalence shifts of SARS-
292 CoV-2 variants (Fig. 1a). In contrast, non-antigenic mutations in influenza proteins are under
293 broad negative selection; observed changes often compensate the deleterious collateral effects
294 of antigenic evolution on conserved molecular traits (including protein stability and receptor
295 binding)^{50–52}. If SARS-CoV-2 reaches a similar endemic state, antigenic evolution is expected
296 to slow down and most intrinsic changes, e.g., in binding affinity to human receptors, will
297 become compensatory. Recent findings of compensatory evolution leading to Omicron support
298 this scenario⁵³.

299 The expected transition of SARS-CoV-2 to gradual, multi-faceted antigenic evolution will
300 open the possibility to predict the future evolution of the viral population by data-driven fitness
301 models^{24,37,38} and to inform preemptive vaccination strategies⁵⁴. Previous work has estab-
302 lished an important prerequisite of predictions: neutralisation assays of human antisera against
303 viral strains quantify the immune protection of human cohorts against secondary infections^{13,14}.
304 Here, we have shown that this data can be harvested at the population scale, to compute im-
305 mune drift and inform antigenic fitness models. As a first step of short-term predictions, we have
306 identified emerging variants in antigenic selection hotspots, in quantitative agreement with their
307 observed clade growth (Fig. 4). This and future predictions of SARS-CoV-2 evolution require
308 integrated analysis of genome sequences, epidemiological records, and increasingly complex anti-
309 genic data. While sequence and epidemiological data are already collected in large amounts, our
310 analysis calls for world-wide, real-time tracking of antigenic evolution by cross-neutralisation
311 assays. This will be critical for our ability to predict antigenic escape evolution and to integrate
312 such predictions into vaccine design.

313 Methods

314 **Sequence data and primary sequence analysis.** The study is based on sequence data from the GI-
 315 SAID EpiCov database⁵ available until 06-22-2022. For quality control, we truncate the 3' and 5' regions
 316 of sequences and remove sequences that contain more than 5% ambiguous sites or have an incomplete
 317 collection date. We align all sequences against a reference isolate from GenBank⁵⁵ (MN908947), using
 318 MAFFT v7.490⁵⁶. Then we map sequences to Variants of Concern/Interest (VOCs/VOIs), using the set
 319 of identifier amino acid changes given in Outbreak.info²⁷. As a cross-check, we independently infer a
 320 maximum-likelihood (ML) strain tree from quality-controlled sequences under the nucleotide substitution
 321 model GTR+G of IQTree⁵⁷, using the reference isolate hCoV-19/Wuhan/Hu-1/2019 (GISAID-Accession:
 322 EPI ISL 402125) as root. For assessment of the tree topology, we use the ultrafast bootstrap function⁵⁸
 323 with 1000 replicates. Internal nodes are timed by TreeTime⁵⁹ with a fixed clock rate of 8×10^{-4} under
 324 a skyline coalescent tree prior⁶⁰. Consistently, variants are mapped to unique genetic clades (subtrees)
 325 of the ML tree (Fig. 2b).

326 **Frequency trajectories of variants.** For a given variant i , we define the smoothed count $n_i(t) =$
 327 $Z^{-1} \sum_{\nu \in i} \exp[-(t - t_\nu)^4/\delta^4]$, where the sum runs over all sequences ν mapped to variant i , t_ν is the
 328 collection date of sequence ν , and Z is a normalisation constant. We use a smoothing period $\delta =$
 329 33d. The corresponding variant frequency is then defined by normalisation over all co-existing variants,
 330 $x_i(t) = n_i(t) / \sum_i n_i(t)$. These frequency trajectories, evaluated separately for each region of this study,
 331 are shown in Fig. 1 and Figs. S1, S2.

332 **Inference of empirical selection.** In a population of different variants, the absolute fitness of each
 333 variant is defined as the growth rate of its population,

$$f_i(t) = \frac{\dot{N}_i(t)}{N_i(t)} \quad (4)$$

334 ($i = 1, \dots, n$). The absolute fitness is related to its reproductive number, defined as the mean number of
 335 new infections generated by an individual during its infectious period τ_0 ,

$$R_i(t) = \exp[\tau_0 f_i(t)]. \quad (5)$$

336 The fitness difference (selection coefficient) between a given pair of variants, $s_{ij}(t) = f_i(t) - f_j(t)$, is
 337 given by $s_{ij}(t) = (d/dt)(\log N_i(t) - \log N_j(t)) = (d/dt)\log((N_i(t)/N_j(t)))$, independently of the other
 338 co-circulating variants. This relation can also be written in terms of the population frequencies $x_i(t) =$
 339 $N_i(t) / \sum_k N_k(t)$, leading to equation (1) of the main text.

340 For each clade shift and each region included in the study, we infer a trajectory of empirical selection,
 341

$$\hat{s} = (\hat{s}(t_1), \hat{s}(t_2), \dots, \hat{s}(t_n)), \quad (6)$$

342 which records the time-dependent fitness difference between invading and ancestral strain, $\hat{s}(t_i) =$
 343 $f_{\text{inv}}(t_i) - f_{\text{anc}}(t_i)$ ($i = 1, \dots, n$). A hat distinguishes these empirical selection coefficients from their
 344 model-based counterparts introduced below. At each point of the trajectory, we evaluate the selection
 345 gradient of equation (1),

$$\hat{s}(t_i) = \frac{1}{\Delta t} \left[\log \left(\frac{x_{\text{inv}}(t_i + \Delta t/2)}{x_{\text{anc}}(t_i + \Delta t/2)} \right) - \log \left(\frac{x_{\text{inv}}(t_i - \Delta t/2)}{x_{\text{anc}}(t_i - \Delta t/2)} \right) \right] \quad (i = 1, \dots, n), \quad (7)$$

346 using a time window $\Delta t = 30$ d for the $1 - \alpha$ and $\delta - o$ shifts and $\Delta t = 40$ d for the $\alpha - \delta$ shift (which
 347 extends over a longer period). Increasing Δt reduces the statistical error of $\hat{s}(t_i)$ but reduces the time
 348 span covered by a trajectory \hat{s} . We evaluate equation (6) for the maximal time interval such that
 349 $x_{\text{anc}}(t_i \pm \Delta t/2) > 0.01$ and $x_{\text{inv}}(t_i \pm \Delta t/2) > 0.01$ along the entire trajectory. The start point
 350 t_1 is the first day when $x_{\text{inv}}(t - \Delta t/2) > 0.01$. From this point, selection is recorded weekly,
 351 $t_i - t_{i-1} = 7$ d ($i = 2, \dots, n$), and t_n is the last point of this sequence where $x_{\text{anc}}(t + \Delta t/2) > 0.01$.
 352 Single measurements $\hat{s}(t_i)$ are excluded when at least one of the sequence counts $n_{\text{anc}}(t_i \pm \Delta t/2)$
 353 or $n_{\text{inv}}(t_i \pm \Delta t/2)$ is < 10 . Statistical errors for selection trajectories are evaluated by binomial
 354 sampling of counts $n_{\text{anc}}(t)$ and $n_{\text{inv}}(t)$ with a pseudocount of 1. Empirical selection trajectories
 355 are reported in Fig. 1b and Figs. S1-S3.

356 For the subsequent analysis, we grade the complete clade shifts $1 - \alpha$, $\alpha - \delta$, $\delta - o1$, $o1 -$
357 $o2$ by the time dependence of their empirical selection trajectories (Fig. S3). We evaluate
358 two summary statistics: (i) the amount of systematic time-dependent variation of selection,
359 defined as $\text{Var}(s_{\text{lin}})$, averaged over regions, where $s_{\text{lin}}(t)$ is a linear regression to the ensemble
360 of trajectories; (ii) the statistical significance of the linear regression, P (two-sided Wald test).
361 This identifies two shifts with substantial, statistically significant time-dependent variation of
362 selection, $\alpha - \delta$ and $\delta - o$.

363 **Infection and vaccination trajectories.** Daily vaccination and infection rates for individual
364 regions have been obtained from Ourworldindata.org⁶ and from CDC COVID Data Tracker⁷
365 for US states (download date: 06-22-2022). Clade-specific infection rates $y_k(t)$ are computed by
366 multiplying the total daily infection rates reported in each region with the simultaneous viral
367 clade frequencies $x_k(t)$. The resulting cumulative population fractions of infected individuals,
368 $y_k(t)$, together with cumulative population fractions of primary and booster vaccinations, $y_{\text{vac}}(t)$
369 and $y_{\text{bst}}(t)$, are reported in Figs. S1-S2.

370 **Data integration for regional analysis.** This study is based on sequence data and epidemiological data from multiple regions (countries and US states) for parallel analysis. Sequence data is used to infer empirical selection trajectories for individual clade shifts, as defined in equations (6) and (7). Epidemiological records provide input to the antigenic fitness model, equations (2) and (3). Evaluation of the fitness model, which is detailed below, integrates data of both categories and requires stringent criteria of data availability and comparability.

371 To enable this analysis, we choose the set of countries to be included in model inference
372 based on uniform criteria. Additionally, we include 3 US states (New York, Texas, California),
373 each representative of a different geographic region, that satisfy the same criteria. For each
374 clade shift $\text{anc} - \text{inv}$, we require the following: (i) anc and inv are majority variants at times
375 t and $t' > t$ of the clade shift, respectively; i.e., $x_{\text{anc}}(t) > 0.5$ and $x_{\text{inv}}(t) > 0.5$. This criterion
376 excludes regions where other variants are prevalent during the shift $\text{anc} - \text{inv}$ (e.g., Brazil and
377 South Africa have $x_\alpha < 0.5$ throughout the $1 - \alpha$ shift). (ii) anc and inv have a combined,
378 smoothed sequence count $n_{\text{anc}}(t) + n_{\text{inv}}(t) > n_0$ throughout the clade shift. This criterion
379 ensures that the empirical frequencies $x_{\text{anc}}(t)$ and $x_{\text{inv}}(t)$, especially minority frequencies, can
380 be estimated with reasonable statistical errors. We use threshold values $n_0 = 500$ for $1 - \alpha$ and
381 $\alpha - \delta$ and $n_0 = 750$ for $\delta - o$ (reflecting the increased sequence availability). (iii) The empirical
382 selection trajectory \hat{s} contains at least 4 ($1 - \alpha$, $\delta - o$) or 6 ($\alpha - \delta$) measured points $\hat{s}(t_i)$; the
383 threshold values reflect the relative duration of shifts. This criterion ensures a sufficient signal-
384 to-noise ratio for inference of temporal variation along the trajectory. (iv) In the $\delta - o$ shift, the
385 cumulative fraction of o infections exceeds a threshold value, $y_o > 0.01$. The o variant, which is
386 characterised by many less severe cases, is likely to be particularly affected by underreporting.
387 This criterion excludes regions with very low o count (y_o is less than $\sim 20\%$ of the remaining
388 regions) and ensures that cross-immunity trajectories, as given by equation (2), can be evaluated
389 across regions with sufficient consistency. (v) Vaccinations have been predominantly by mRNA
390 vaccines and epidemiological records in the database⁶ are complete. This criterion ensures that
391 antigenic data for mRNA vaccines can be used uniformly (Table S1). It excludes regions with
392 substantial use of viral vector vaccines (e.g., the UK) and with partial records (e.g., for booster
393 vaccinations in Sweden and Croatia).

394 Based on these criteria, our analysis includes (i) 11 regions for the $1 - \alpha$ shift, (ii) 16 regions
395 for the $\alpha - \delta$ shift, and (iii) 14 regions for the $\delta - o$ shift (Fig. 1, Fig. 3, Figs. S1-S3). Regions
396 analysed for both $\alpha - \delta$ and $\delta - o$ are used for the long-term trajectories (Fig. 4, Fig. S4). Scope
397 and limitations of this set of regions for the inference of selection are described below.

398 **Antigenic data.** Neutralisation assays for SARS-CoV-2 test the potency of antisera induced
399 by a given primary immunisation to neutralise viruses of different variants. Log dilution titers

405 measure the minimum antiserum concentration required for neutralisation,

$$T_i^k = \log_2 \frac{K_0}{K_i^k}, \quad (8)$$

406 relative to a reference concentration K_0 . Hence, \log_2 titer differences, or neutralisation fold
 407 changes, $\Delta T_{ij}^k \equiv \Delta T_i^k - T_j^k$, measure differences in antigenicity between variants, $\Delta T_{ij}^k =$
 408 $\log_2(K_j^k/K_i^k)$. We note that these differences are specific to each primary challenge (immune
 409 channel) k . For example, the inequality $T_{\alpha\delta}^{\text{bst}} < T_{\alpha\delta}^{\text{vac}}$ reflects the increased breadth of booster
 410 vaccinations compared to primary vaccinations. In contrast, uni-valued antigenic distances be-
 411 tween variants, d_{ij} , can be computed from the titer matrix (T_i^k) by multi-dimensional scaling
 412 methods^{61,62}. Such distance measures average over inhomogeneities between immune channels.

413 Here we define a matrix of titer drops ΔT_i^k ,

$$\Delta T_i^k = T_*^k - T_i^k \quad (i = \alpha, \delta, o(01), o2, o45; \quad k = \alpha, \delta, o(01), o2, o45, \text{vac}, \text{bst}), \quad (9)$$

414 with respect to a reference for each immune channel, $T_*^k = T_1^k$ ($k = \alpha, \delta, o(01), o2, o45$) and
 415 $T_*^k = T_1^k$ ($k = \text{vac}, \text{bst}$). This procedure eliminates technical differences between assays in
 416 absolute antibody concentration. We assemble this matrix in Table S1, using primary data
 417 from different sources^{3,8,9,28,29,44,45,48,63–78}. We proceed as follows: (i) For matrix elements
 418 with available data, ΔT_i^k is the average of the corresponding primary measurements. This
 419 procedure eliminates technical differences between assays in absolute antibody concentration.
 420 As appropriate for the analysis in our set of regions, all vaccination titers refer to mRNA vaccines.
 421 (ii) If no data are available for ΔT_i^k but the conjugate titer ΔT_k^i has been measured, we use the
 422 approximate substitution $\Delta T_i^k \approx \Delta T_k^i$, as discussed in ref. [79]. (iii) If no data are available for
 423 ΔT_i^k but the titer ΔT_j^k of a closely related clade has been measured, we use the approximate
 424 substitution $\Delta T_i^k \approx \Delta T_j^k$, which should be understood as a lower bound (this applies to the
 425 recent variants $o2$ and $o45$).

426 The matrix of absolute neutralisation titers, T_i^k , is then computed by equation (9), combining
 427 the titer drops ΔT_i^k of Table S1 and the reference titers $T_*^k = 6.5$, ($k = \alpha, \delta, o(01), o2, o45$),
 428 $T_*^{\text{vac}} = 7.8$, $T_*^{\text{bst}} = 9.8$ reported in ref. [30]. A titer difference between vaccination and booster,
 429 $T_*^{\text{bst}} - T_*^{\text{vac}} \approx 2.0$, has been observed in several studies^{9,10,44}. The titers T_i^k enter the cross-
 430 immunity functions c_i^k , $C_i^k(t)$, and \bar{c}_i^k defined below and are shown in Fig. 2ab.

431 The decay of antibody concentration after primary immunisation has been characterised in
 432 recent work^{31,32}. Here we describe this effect by a linear titer reduction with time after primary
 433 challenge,

$$T_i^k(\Delta t) = T_i^k - \frac{\Delta t}{\tau}, \quad (10)$$

434 corresponding to an exponential decay of antibody concentration, with a uniform decay time
 435 90d (i.e., half life $\tau = 65$ d). This is broadly consistent with experimental data; we infer decay
 436 times in the range [60, 170]d from several studies^{8,9,31,32,44}. In addition, we check that varying
 437 τ in this range does not affect our results (in particular, the rank order of variants with respect
 438 to antigenic fitness remains unchanged).

439 **Cross-immunity trajectories.** The cross-immunity factor c_i^k is defined as the relative reduction
 440 in infections by variant i induced by (recent) immunisation in channel k . As shown in recent
 441 work^{13,14}, absolute titers of SARS-CoV-2 neutralisation assays can predict cross-immunity, $c_i^k =$
 442 $H(T_i^k)$ with

$$H(T) = \frac{1}{1 + \exp[-\lambda(T - T_{50})]}. \quad (11)$$

443 This relation has been established in ref. [13] with constants $T_{50} = 4.2$ and $\lambda = 0.9$. The resulting
 444 cross-immunity factors $c_i^k(\Delta t)$ include antibody decay, as given by equation (10). Hence, they

445 depend on the time since primary immunisation,

$$c_i^k(\Delta t) = H(T_i^k - \frac{\Delta t}{\tau}). \quad (12)$$

446 These factors enter the population cross-immunity functions $C_i^k(t)$, equation (2), which become

$$C_i^k(t) = \int^t H(T_i^k - \frac{t-t'}{\tau}) \dot{y}(t') dt'. \quad (13)$$

448 These functions enter all evaluations of the fitness model, equation (3) (Fig. 1c, Fig. 3, Fig. 4ef,
449 Figs. S1, S2, and S4).

450 To display the emergence of selection hotspots, we approximate the cross-immunity functions,
451 equation (13), by time-independent effective factors and immune weights. (i) The effective
452 cross-immunity factors \bar{c}_i^k are obtained from equation (12) at a fixed time delay $\Delta t = \tau$ after
453 primary immunisation, $\bar{c}_i^k = H(T_i^k - 1)$, which accounts for the decay of immune response in
454 an approximate way. These factors define the antigenic landscapes shown in Fig. 4a-c. (ii) The
455 immune weight functions,

$$Q_k(t) = \int_{-\infty}^t H(T_0 - \frac{t-t'}{\tau}) \dot{y}(t') dt', \quad (14)$$

456 measure the effective population fractions of immune individuals in channel k . They account
457 for immune decay from a fixed reference titer $T_0 = 6.5$ and follow the time-dependence of
458 cross-immunity functions $C_i^k(t)$ and selection coefficients $s_k(t)$ in an approximate way (Fig. 4de,
459 Fig. S4). Selection hotspots emerge if large steps on an antigenic landscape, $\bar{c}_{\text{inv}}^k - \bar{c}_{\text{anc}}^k$, coincide
460 with sizeable immune weights $Q_k(t)$ in one or more immune channels (Fig. 4).

461 **Fitness model.** Equation (3) expresses the fitness of viral variants as a sum of intrinsic and
462 antigenic fitness components, $f_i(t) = f_i^0 + \sum_k f_i^k(t)$. Intrinsic fitness, f_0 , integrates contributions
463 from several molecular phenotypes, including protein stability, host receptor binding, and traits
464 related to intra-cellular viral replication. The antigenic components $f_i^k(t)$ describe the impact of
465 antibody binding on viral growth, summed over the immune repertoire components of different
466 channels of primary infection or vaccination. Importantly, the input of this fitness model can
467 be learned by integration of sequence data, epidemiological records, and antigenic assays.

468 The additive form the fitness model neglects epistasis between fitness components. The
469 additivity assumption is justified between intrinsic and antigenic fitness, because these compo-
470 nents are associated to different stages of the viral replication cycle. The additivity of antigenic
471 fitness components rests on the approximation of a well-mixed host population and short infec-
472 tion times. In this approximation, each viral lineage is subject to a dense sequence of random
473 encounters with hosts of different immune channels k , leading to averaging of antigenic fitness
474 effects. Multiple infections in an individual can generate additional immune channels; however,
475 these effects are relatively small over the short periods of SARS-CoV-2 evolution studied in this
476 paper.

477 Our analysis of the fitness model focuses on selection coefficients between co-existing variants,

$$s_{ij}(t) \equiv f_i(t) - f_j(t) = s_{ij}^0 - \sum_k \gamma_k [C_i^k(t) - C_j^k(t)], \quad (15)$$

479 because these can directly be compared with their empirical counterparts $\hat{s}_{ij}(t)$. Of equal impor-
480 tance, selection coefficients within a region decouple from the changes in viral ecology within that
481 region. Specifically, seasonality and contact limitations can generate strongly time-dependent
482 reproductive numbers. However, any modulation of the form $R(t) \rightarrow \alpha(t)R(t)$ leaves the se-
483 lection coefficients s_{ij} invariant, as can be seen from equation (5). Our inference of empirical
484 selection, as described above, is also independent of the underlying infectious period τ_0 , which
485 may itself be under evolutionary pressure and change with time^{80,81}. To keep this independence,
486 we report all selection coefficients in fixed units [1/d].

487 **Inference of fitness model parameters.** The free parameters γ_k ($k = 1, \dots, n$) measure the
488 fitness effect of each cross-immunity component. These parameters calibrate the model to data
489 of complex real populations differing, for example, in population structure (including incidence
490 structure), infection histories, and monitoring of infections. To avoid overfitting, we use a
491 minimal model with just 3 global antigenic parameters: (i) A basic rate $\gamma_{\text{vac}} = \gamma_{\text{bst}}$ translates
492 cross-immunity generated by vaccination into units of selection. (ii) This rate is downweighted
493 to a value $\gamma'_{\text{vac}} = a\gamma_{\text{vac}}$ for the shift $\delta - o$ and later shifts. This can be seen as a heuristic to
494 account for the effect of double infections⁴², which increase cross-immunity and decrease cross-
495 immunity differences between variants. (iii) Cross-immunity in all infection channels is uniformly
496 upweighted, $\gamma_k = b\gamma_{\text{vac}}$, to account for underreporting of infections relative to vaccinations.

497 Our inference proceeds in two steps. First, we train the antigenic fitness model using data
498 from the clade shifts $\alpha - \delta$ and $\delta - o$. These shifts are suitable because they carry sizeable antigenic
499 advance $\Delta T_{\alpha\delta}^{\text{vac}}$ and $\Delta T_{\delta o}^{\text{vac}}$ (Fig. 2a), selection shows a substantial and statistically significant
500 time dependence (Fig. 1b), and population immunity has started to pick up (Fig. 4d). We infer
501 the ML likelihood model by aggregation of log likelihood scores over the sets of regional selection
502 trajectories for the clade shifts $\alpha - \delta$ and $\delta - o$. We use the score function

$$L(\hat{\mathbf{s}}, \mathbf{s}) = - \sum_{i=1}^n \frac{(\Delta\hat{s}(t_i) - \Delta s(t_i))^2}{2\sigma^2(t_i)} \quad (16)$$

503 for a single empirical selection trajectory $\hat{\mathbf{s}}$, equation (6), and its model-based counterpart \mathbf{s} . This
504 score evaluates selection change, $\Delta s(t) = s(t) - \langle s \rangle$, where brackets denote averaging over time.
505 Hence, the fitness model is trained on the time-dependence of selection in each region, in order to
506 avoid the confounding factor of heterogeneity across regions. The expected square deviation is
507 $\sigma^2(t_i) = \sigma_s^2(t_i) + \sigma_0^2$; the first term describes the sampling error of sequence counts, which enters
508 frequency and empirical selection estimates, the second term summarises fluctuations unrelated
509 to sequence counts. The total log likelihood score is the sum $L = \sum L(\mathbf{s}, \hat{\mathbf{s}})$, which runs over
510 both shifts and all included regions. Table S2 lists the ML parameters $\gamma_{\text{vac}}, a, b$ and the ML
511 score L relative to a null model of time-independent selection (see below). The 95% confidence
512 intervals of the inferred parameters are computed by resampling the empirical selection data with
513 fluctuations σ^2 . We note that the ML values $a < 1, b > 1$ are consistent with the interpretation
514 of these parameters as weighting factors accounting for double-infections and underreporting
515 (see above). Second, we infer the intrinsic selection for each shift as the difference between
516 empirical selection and ML antigenic selection, $s_0 = \langle\langle \hat{s} - s_{\text{ag}} \rangle\rangle$, where the double brackets
517 denote averaging over time and regions. The ML antigenic selection coefficients, $\langle\langle s_k \rangle\rangle$, and the
518 intrinsic selection coefficient s_0 between invading and ancestral variant are listed in Table S3;
519 see also Fig. 3b. Confidence intervals are computed by resampling model parameters with their
520 confidence intervals. Consistently, we infer weak antigenic selection for the shifts $1 - \alpha$ and
521 $o1 - o2$, which also show only weak time dependence of selection (Fig. S3).

522 **Significance analysis of the fitness model.** To assess the statistical significance of our inference,
523 we compare four fitness models of the form (3): the full model used in the main text (VI:
524 antigenic selection by vaccination and infection, intrinsic selection), two partial models (V:
525 antigenic selection only by vaccination, intrinsic selection; I: antigenic selection only by infection,
526 intrinsic selection), and a null model (0: intrinsic selection only). We infer conditional ML
527 parameters for each model and we rank models by their ML score difference to the null model,
528 $\Delta L = L - L_0$ (Table S2). An alternative ranking by BIC score⁸², which contains a score
529 penalty for the number of model parameter, leads to the same result. We observe the following:
530 (i) All antigenic fitness models have significantly higher scores than the null model, which shows
531 that the empirical selection data are incompatible with time-independent selection. (ii) The
532 full model has a significantly higher score than any of the other models; both vaccination and
533 infection are significant components of antigenic selection. (iii) Vaccination explains a larger part

534 of the time-dependent data than infection ($\Delta L_V > \Delta L_I$), which is consistent with the ranking of
535 selection coefficients inferred from the full model (Fig. 3b, Table S3). (iv) The score gain of the
536 full model is less than the sum of its parts ($\Delta L_{IV} < \Delta L_V + \Delta L_I$). This can be associated with
537 statistical correlations in the input data for both antigenic model components. For example, the
538 fraction of vaccinated individuals y_{vac} is weakly anti-correlated with the fraction of δ infections,
539 y_δ .

540 **Fitness trajectories.** Long-term fitness trajectories display clade turnover of multiple successive
541 shifts (Fig. 4f, Fig. S4). For each of the variants $\alpha, \delta, o(01), o2, o45$, we plot the time-
542 dependent fitness gap, $\delta f_i(t) = f_i(t) - \bar{f}(t)$, where $\bar{f}(t) = \sum_j x_j(t)f_j(t)$ is the mean population
543 fitness. Like selection coefficients, fitness gaps decouple from ecological factors affecting absolute
544 growth (see the discussion above). Assuming that the fitness difference between ancestral and
545 invading variant, $s(t)$, is dominant during each crossover, we obtain the fitness gap trajectories
546 $\delta f_{\text{anc}}(t) = -s(t)x_{\text{inv}}(t)$ and $\delta f_{\text{inv}}(t) = -s(t)[1 - x_{\text{inv}}(t)]$, as well as their empirical counterparts
547 $\hat{\delta f}_{\text{anc}}(t)$ and $\hat{\delta f}_{\text{inv}}(t)$. For each variant, we patch trajectories from origination to near-fixation
548 (here in the time interval $(t_{0,i}, t_{f,i})$ given by $x_i(t_{0,i}) = 0.01$ and $x_i(t_{f,i}) = 0.99$). The long-term
549 trajectories display selection hotspots and confirm the quantitative agreement between empirical
550 and model-based fitness.

551 **Model-based inference of selection across regions.** As shown by the preceding analysis, we
552 can infer a statistically significant fitness model with few, global parameters from sequence
553 and epidemiological data aggregated over a set of regions and combined with antigenic data.
554 The model describes common time-dependent patterns of selection in these regions and serves
555 two main purposes: to provide a breakdown of selection in intrinsic and antigenic components
556 (Fig. 3) and to display selection hotspots in long-term trajectories (Fig. 4). Our inference
557 procedure rests on stringent criteria for the joint availability of sequence and epidemiological
558 data in each of these regions (as listed above). A number of points support this procedure:
559 (i) The results are robust under variation of the inclusion criteria for regions. In particular, the
560 signal of antigenic selection in data and model is broadly distributed over regions (Figs. S1-S3).
561 Hence, the selection averages reported in Fig. 3b and Table S3 are reproducible in subsampled
562 sets of regions. (ii) Within the set of regions included, the model is applicable beyond the $\alpha - \delta$
563 and $\delta - o$ shifts used for training. The early $1 - \alpha$ shift and the recent $o1 - o2$ shifts serve as
564 controls. In both cases, we infer weak antigenic selection, consistent with weak time dependence
565 of empirical selection (Fig. S4). For the emerging $o2 - o45$ shift, strong antigenic selection is
566 consistent with fast initial growth of the new variants.

567 Our model-based inference of selection excludes a number of regions that do not fulfil the
568 criteria of joint data availability. (i) For the $\alpha - \delta$ shift, several regions are excluded because
569 VOCs other than α were majority variants prior to the shift to δ (for example, Beta in South
570 Africa and Gamma in Brazil). Unlike $\alpha - \delta$, these shifts do not involve antigenic advance in
571 the vaccination channel; i.e., $\Delta T_{\text{anc}\delta}^{\text{vac}} < 0$. However, we lack comprehensive antigenic data on
572 other VOCs as input for the fitness model. (ii) For the $\delta - o$ shift, several regions are excluded
573 because of low reported incidence (e.g., Brazil, California, India, Mexico, Poland, South Korea,
574 Turkey, Texas). Most of these regions show a signal of time-dependent selection consistent with
575 the regions included; however, much lower reported incidence counts prevent reliable immune
576 tracking of infection channels during the $\delta - o$ shift. Variation in reported incidence can, in
577 principle, be incorporated into the fitness model by region-dependent γ_{vac} factors, but this would
578 likely lead to overfitting. We conclude that at current levels of data availability, a comprehensive
579 cross-regional analysis is not feasible.

580 **Model-free inference of selection in regions with low vaccination coverage.** Countries
581 with low vaccination coverage during the $\alpha - \delta$ shift ($y_{\text{vac}} < 0.1$) disqualify for the model-based
582 analysis because α was not a majority variant (India, Malaysia, Russia, Philippines, Indonesia,

583 South Africa, South Korea) or sequence counts are too low for the inference of selection trajec-
584 tories (Australia). For this set of countries, we can still infer a selection coefficient s by fitting a
585 sigmoid function to the frequency trajectory $x_\delta(t)$ (to be interpreted as the growth difference be-
586 tween δ and the average of all other coexisting variants). We find lower region-averaged selection
587 in the set of low-vaccination countries compared to other countries ($\langle\langle \hat{s} \rangle\rangle = 0.08$ vs. $\langle\langle \hat{s} \rangle\rangle = 0.12$).
588 This is qualitatively consistent with our model-based inference of vaccination-induced selec-
589 tion; however, the genetic heterogeneity of this clade shift prevents a systematic breakdown of
590 antigenic selection into immune channels.

591 **Data availability.** The datasets analysed in this study are available in published work.

592 **Code availability.** The code used in this study is available at https://github.com/m-meijers/vaccine_effect

594 **Acknowledgements.** We thank Florian Klein and Kanika Vanshylla for discussions. This
595 work has been partially funded by Deutsche Forschungsgemeinschaft grant CRC 1310 *Predictabil-
596 ity in Evolution*.

- 598 1. Ozono, S. *et al.* SARS-CoV-2 D614G spike mutation increases entry efficiency with enhanced ACE2-binding
599 affinity. *Nature Communications* **12** 848 (2021).
- 600 2. Meng, B. *et al.* Recurrent emergence of SARS-CoV-2 spike deletion H69/V70 and its role in the Alpha
601 variant B.1.1.7. *Cell Reports* **35** 13 (2021).
- 602 3. Mlcochova, P. *et al.* SARS-CoV-2 B.1.617.2 Delta variant replication and immune evasion. *Nature* **599**
603 114-119 (2021).
- 604 4. Harvey, W. T. *et al.* SARS-CoV-2 variants, spike mutations and immune escape *Nature reviews microbiology*
605 **19** 409-424 (2021).
- 606 5. Shu, Y. & McCauley, J. GISAID: Global initiative on sharing all influenza data – from vision to reality
607 *EuroSurveillance* **22** 13 (2017).
- 608 6. Ritchie, H. *et al.* Coronavirus pandemic (covid-19). *Our World in Data* (2020).
609 <Https://ourworldindata.org/coronavirus>.
- 610 7. Centers for Disease Control and Prevention. COVID Data Tracker. Atlanta, GA: US Department of Health
611 and Human Services, CDC. <Https://covid.cdc.gov/covid-data-tracker>. Accessed: 22-06-2022.
- 612 8. Planas, D. *et al.* Reduced sensitivity of SARS-CoV-2 variant Delta to antibody neutralization. *Nature* **596**,
613 276–280 (2021).
- 614 9. Planas, D. *et al.* Considerable escape of SARS-CoV-2 Omicron to antibody neutralization. *Nature* **602**
615 671–675 (2021).
- 616 10. Garcia-Beltran, W. F. *et al.* mRNA-based COVID-19 vaccine boosters induce neutralizing immunity against
617 SARS-CoV-2 Omicron variant. *Cell* **185**, 457–466 (2022).
- 618 11. Wrobel, A. G. *et al.* SARS-CoV-2 and bat RaTG13 spike glycoprotein structures inform on virus evolution
619 and furin-cleavage effects. *Nature Structural and Molecular Biology* **27**, 763–767 (2020).
- 620 12. Zeng, C. *et al.* Neutralization and stability of SARS-CoV-2 Omicron variant. *bioRxiv*: (2021). URL
621 <Https://doi.org/10.1101/2021.12.16.472934>.
- 622 13. Khoury, D. S. *et al.* Neutralizing antibody levels are highly predictive of immune protection from symptomatic
623 SARS-CoV-2 infection. *Nature Medicine* **27**, 1205–1211 (2021).
- 624 14. Feng, S. *et al.* Correlates of protection against symptomatic and asymptomatic SARS-CoV-2 infection.
625 *Nature Medicine* **27**, 2032–2040 (2021).
- 626 15. Tracking SARS-CoV-2 variants. <Https://www.who.int/en/activities/tracking-SARS-CoV-2-variants>. Accessed: 22-06-2022.
- 628 16. Davies, N. G. *et al.* Estimated transmissibility and impact of SARS-CoV-2 lineage B.1.1.7 in England.
629 *Science* **372** (2021).
- 630 17. Dhar, M. S. *et al.* Genomic characterization and epidemiology of an emerging SARS-CoV-2 variant in Delhi,
631 India. *Science* **374**, 995–999 (2021).
- 632 18. Kepler, L., Hamins-Puertolas, M. & Rasmussen, D. A. Decomposing the sources of SARS-CoV-2 fitness
633 variation in the united states. *Virus Evolution* **7** (2021).
- 634 19. Ulrich, L. *et al.* Enhanced fitness of SARS-CoV-2 variant of concern Alpha but not Beta. *Nature* **602**,
635 307–313 (2022).
- 636 20. Grenfell, B. T. *et al.* Unifying the epidemiological and evolutionary dynamics of pathogens. *Science* **303**,
637 327–332 (2004).

- 638 21. Rella, S. A., Kulikova, Y. A., Dermitzakis, E. T. & Kondrashov, F. A. Rates of SARS-CoV-2 transmission
639 and vaccination impact the fate of vaccine-resistant strains. *Scientific Reports* **11** (2021).
- 640 22. Saad-Roy, C. M. *et al.* Epidemiological and evolutionary considerations of SARS-CoV-2 vaccine dosing
641 regimes. *Science* **372**, 363–370 (2021).
- 642 23. Lobinska, G. *et al.* Evolution of resistance to COVID-19 vaccination with dynamical social distancing. *Nature*
643 **507**, 57–61 (2022).
- 644 24. Luksza, M. & Lässig, M. A predictive fitness model for influenza. *Nature Human Behaviour* **6**, 57–61 (2014).
- 645 25. Wen, F. *et al.* The potential beneficial effects of vaccination on antigenically evolving pathogens. *American*
646 *Naturalist* **199** 2 193–206 (2022).
- 647 26. Wen, F. T., Bell, S. M., Bedford, T. & Cobey, S. Estimating vaccine-driven selection in seasonal influenza.
648 *Viruses* **10** (2018).
- 649 27. Mullen, J. L. *et al.* Outbreak.info (2020). URL <https://outbreak.info>.
- 650 28. van der Straten, K. *et al.* Mapping the antigenic diversification of SARS-CoV-2. *bioRxiv* (2022). URL
651 <https://doi.org/10.1101/2022.01.03.21268582>.
- 652 29. Wilks, S. H. *et al.* Mapping SARS-CoV-2 antigenic relationships and serological responses. *bioRxiv* (2022).
653 URL <https://doi.org/10.1101/2022.01.28.477987>.
- 654 30. Polack, F. P. *et al.* Safety and efficacy of the BNT162b2 mRNA Covid-19 vaccine. *New England Journal of*
655 *Medicine* **383**, 2603–2615 (2020).
- 656 31. Iyer, A. S. *et al.* Persistence and decay of human antibody responses to the receptor binding domain of
657 SARS-CoV-2 spike protein in COVID-19 patients. *Science Immunology* **5** (2020).
- 658 32. Israel, A. *et al.* Large-scale study of antibody titer decay following BNT162b2 mRNA vaccine or SARS-CoV-2
659 infection. *Vaccines* **10** (2022).
- 660 33. Coudeville, L. *et al.* Relationship between haemagglutination-inhibiting antibody titres and clinical pro-
661 tection against influenza: development and application of a bayesian random-effects model. *BMC Medical*
662 *Research Meth.* **10** (2010).
- 663 34. Dunning, A. J. *et al.* Correlates of protection against influenza in the elderly: Results from an influenza
664 vaccine efficacy trial. *Clinical and Vaccine Immunology* **23**, 228–235 (2016).
- 665 35. Rotem, A. *et al.* Evolution on the biophysical fitness landscape of an rna virus. *Molecular Biology and*
666 *Evolution* **35**, 2390–2400 (2018).
- 667 36. Meijers, M., Vanshylla, K., Gruell, H., Klein, F. & Laessig, M. Predicting in vivo escape dynamics of HIV-1
668 from a broadly neutralizing antibody. *PNAS* **118** (2021).
- 669 37. Morris, D. H. *et al.* Predictive modeling of influenza shows the promise of applied evolutionary biology
670 (2018).
- 671 38. Huddleston, J. *et al.* Integrating genotypes and phenotypes improves long-term forecasts of seasonal influenza
672 A/H3N2 evolution. *eLife* **9**, 1–48 (2020).
- 673 39. Gog, J. R. & Grenfell, B. T. Dynamics and selection of many-strain pathogens. *PNAS* **99**, 17209–17214
674 (2002).
- 675 40. Lipsitch, M. *et al.* Mapping person-to-person variation in viral mutations that escape polyclonal serum
676 targeting influenza hemagglutinin (2019).
- 677 41. Cobey, S. & Hensley, S. E. Immune history and influenza virus susceptibility. *Current Opinion in Virology*
678 **22**, 105–111 (2017).
- 679 42. Bates, T. A. *et al.* Vaccination before or after SARS-CoV-2 infection leads to robust humoral response and
680 antibodies that effectively neutralize variants *Science Immunology* **7** 68 (2022).
- 681 43. Yuan, S. *et al.* Pathogenicity, transmissibility, and fitness of SARS-CoV-2 Omicron in Syrian hamsters.
682 *Science* (2022). URL <https://www.science.org/doi/10.1126/science.abn8939>.
- 683 44. Gruell, H. *et al.* mRNA booster immunization elicits potent neutralizing serum activity against the SARS-
684 CoV-2 Omicron variant. *Nature Medicine* **28**, 477–480 (2022).
- 685 45. Hachmann, N. P. *et al.* Neutralization escape by SARS-CoV-2 Omicron subvariants BA.2.12.1, BA.4, and
686 BA.5. *New England Journal of Medicine* (2022).
- 687 46. Earn, D. J. D., Dushoff, J. & Levin, S. A. Ecology and evolution of the flu. *Trends in Ecology and Evolution*
688 **17**, 334–340 (2002).
- 689 47. Boni, M. F., Gog, J. R., Andreasen, V. & Feldman, M. W. Epidemic dynamics and antigenic evolution in a
690 single season of influenza a. *Proceedings of the Royal Society B: Biological Sciences* **273**, 1307–1316 (2006).
- 691 48. Khan, K. *et al.* Omicron sub-lineages ba.4/ba.5 escape ba.1 infection elicited neutralizing immunity *bioRxiv*
692 (2022). URL <https://doi.org/10.1101/2022.04.29.22274477>.
- 693 49. Fonville, J. M. *et al.* Antibody landscapes after influenza virus infection or vaccination. *Science* **346**,
694 996–1000 (2014).
- 695 50. Strelkowa, N. & Lässig, M. Clonal interference in the evolution of influenza. *Genetics* **192**, 671–682 (2012).
- 696 51. Gong, L. I., Suchard, M. A. & Bloom, J. D. Stability-mediated epistasis constrains the evolution of an
697 influenza protein. *eLife* **2013** (2013).
- 698 52. Lässig, M., Mustonen, V. & Walczak, A. M. Predicting evolution. *Nature Ecology and Evolution* **1** (2017).

- 699 53. Moulana, A. *et al.* Compensatory epistasis maintains ace2 affinity in sars-cov-2 omicron ba.1 *bioRxiv* (2022).
700 URL <https://doi.org/10.1101/2022.06.17.496635>.
- 701 54. Mustonen, V. & Lässig, M. Eco-evolutionary control of pathogens. *PNAS* **117**, 19694–19704 (2020).
- 702 55. Benson, D. A. *et al.* GenBank. *Nucleic Acids Research* **43**, D30–D35 (2015).
- 703 56. Katoh, K. & Standley, D. M. MAFFT multiple sequence alignment software version 7: Improvements in
704 performance and usability. *Molecular Biology and Evolution* **30**, 772–780 (2013).
- 705 57. Minh, B. Q. *et al.* IQ-TREE 2: New models and efficient methods for phylogenetic inference in the genomic
706 era. *Molecular Biology and Evolution* **37**, 1530–1534 (2020).
- 707 58. Hoang, D. T. *et al.* UFBoot2: Improving the ultrafast bootstrap approximation. *Mol. Biol. Evol* **35**, 518–522
708 (2017).
- 709 59. Sagulenko, P., Puller, V. & Neher, R. A. TreeTime: Maximum-likelihood phylodynamic analysis. *Virus
710 Evolution* **4** (2018).
- 711 60. Kingman, J. F. C. The coalescent. *Stochastic Processes and their Applications* **13**, 235–248 (1982).
- 712 61. Smith, D. J., Forrest, S., Hightower, R. R. & Perelson, A. S. Deriving shape space parameters from im-
713 munological data. *J. theor. Biol* **189**, 141–150 (1997).
- 714 62. Smith, D. J., Lapedes, A. S. & Jong, J. C. D. Mapping the antigenic and genetic. *Science* **305**, 371–377
715 (2004).
- 716 63. Bates, T. A. *et al.* Neutralization of SARS-CoV-2 variants by convalescent and BNT162b2 vaccinated serum.
717 *Nature Communications* **12** (2021).
- 718 64. Zhou, D. *et al.* Evidence of escape of SARS-CoV-2 variant B.1.351 from natural and vaccine-induced sera.
719 *Cell* **184**, 2348–2361.e6 (2021).
- 720 65. Garcia-Beltran, W. F. *et al.* Multiple SARS-CoV-2 variants escape neutralization by vaccine-induced humoral
721 immunity. *Cell* **184**, 2372–2383.e9 (2021).
- 722 66. Liu, C. *et al.* Reduced neutralization of SARS-CoV-2 B.1.617 by vaccine and convalescent serum. *Cell* **184**,
723 4220–4236.e13 (2021).
- 724 67. Wang, P. *et al.* Antibody resistance of SARS-CoV-2 variants B.1.351 and B.1.1.7. *Nature* **593**, 130–135
725 (2021).
- 726 68. Liu, Y. *et al.* BNT162b2-elicited neutralization against new SARS-CoV-2 spike variants. *New England
727 Journal of Medicine* **385**, 472–474 (2021).
- 728 69. Planas, D. *et al.* Sensitivity of infectious SARS-CoV-2 B.1.1.7 and B.1.351 variants to neutralizing antibodies.
729 *Nature Medicine* **27**, 917–924 (2021).
- 730 70. Uriu, K. *et al.* Neutralization of the sars-cov-2 mu variant by convalescent and vaccine serum. *New England
731 Journal of Medicine* **385**, 2395–2397 (2021).
- 732 71. Rössler, A., Riepler, L., Bante, D., von Laer, D. & Kimpel, J. SARS-CoV-2 Omicron variant neutralization
733 in serum from vaccinated and convalescent persons. *New England Journal of Medicine* **386**, 698–700 (2022).
- 734 72. Cameroni, E. *et al.* Broadly neutralizing antibodies overcome sars-cov-2 omicron antigenic shift. *Nature*
735 **602**, 664–670 (2022).
- 736 73. Wang, Q. *et al.* Antibody evasion by sars-cov-2 omicron subvariants ba.2.12.1, ba.4; ba.5. *Nature* (2022).
- 737 74. Iketani, S. *et al.* Antibody evasion properties of sars-cov-2 omicron sublineages. *Nature* **604**, 553–556 (2022).
- 738 75. Muik, A. *et al.* Neutralization of SARS-CoV-2 lineage B.1.1.7 pseudovirus by BNT162b2 vaccine-elicited
739 human sera *bioRxiv* (2021). URL <https://doi.org/10.1101/2020.12.30.20249034>.
- 740 76. Bowen, J. E. *et al.* Omicron ba.1 and ba.2 neutralizing activity elicited by a comprehensive panel of human
741 vaccines. *bioRxiv* (2022). URL <http://www.ncbi.nlm.nih.gov/pubmed/35313570>.
- 742 77. Cao, Y. *et al.* Ba.2.12.1, ba.4 and ba.5 escape antibodies elicited by omicron infection division of hiv/aids
743 and sex-transmitted virus vaccines, institute for biological product *bioRxiv* (2022). URL <https://doi.org/10.1101/2022.04.30.489997>.
- 744 78. Mykytyn, A. Z. *et al.* Omicron ba.1 and ba.2 are antigenically distinct sars-cov-2 variants *bioRxiv* (2022).
745 URL <https://doi.org/10.1101/2022.02.23.481644>.
- 746 79. Neher, R. A., Bedford, T., Daniels, R. S., Russell, C. A. & Shraiman, B. I. Prediction, dynamics, and
747 visualization of antigenic phenotypes of seasonal influenza viruses. *Proceedings of the National Academy of
748 Sciences* **113**, E1701–E1709 (2016).
- 749 80. Hart, W. S. *et al.* Generation time of the alpha and delta SARS-CoV-2 variants: an epidemiological analysis.
750 *The Lancet Infectious Diseases* **22**, 603–610 (2022).
- 751 81. Hart, W. S. *et al.* Inference of the SARS-CoV-2 generation time using UK household data. *eLife* **11** (2022).
- 752 82. Schwarz, G. Estimating the dimension of a model. *The Annals of Statistics* **6**, 461–464 (1978).

754 Supplementary Tables and Figures

Table S1: Antigenic data.

| | α | δ | $o (o1)$ | $o2$ | $o45$ | vac | bst |
|----------|----------|----------|------------|------------|------------|-------|------------|
| α | 0 | 1.8 | <i>5.0</i> | <i>5.0</i> | <i>5.0</i> | 0.8 | <i>0.8</i> |
| δ | 1.5 | 0 | <i>4.8</i> | <i>4.8</i> | <i>4.8</i> | 1.7 | 1.5 |
| $o (o1)$ | 5.0 | 4.8 | 0 | 2.1 | <i>2.2</i> | 5.6 | 2.7 |
| $o2$ | < 5.0 | < 4.8 | 1.4 | 0 | 1.2 | < 5.6 | 2.5 |
| $o45$ | < 5.0 | < 4.8 | 2.2 | 1.2 | 0 | < 5.6 | 3.9 |

We list log titer drops, or neutralisation fold changes, $\Delta T_i^k = T_*^k - T_i^k$, of strains from variant i assayed against human antisera induced by primary immunisation (infection or vaccination) with strains of variant k (columns). Numbers are average values of primary data from ref. [3, 8, 9, 28, 29, 44, 45, 48, 63–78]. All vaccination titers refer to mRNA vaccines. Where no primary data is available, titer drops are inferred by symmetry or (as lower bounds) by genetic similarity (numbers in italics, Methods). Absolute titers T_i^k are shifted by the reference titers $T_*^k = 6.5$, ($k = \alpha, \delta, o(o1), o2, o45$), $T_*^{\text{vac}} = 7.8$, $T_*^{\text{bst}} = 9.8$ obtained from ref. [9, 10, 30, 44]; see Methods and Fig. 2a.

Table S2: Ranking of fitness models.

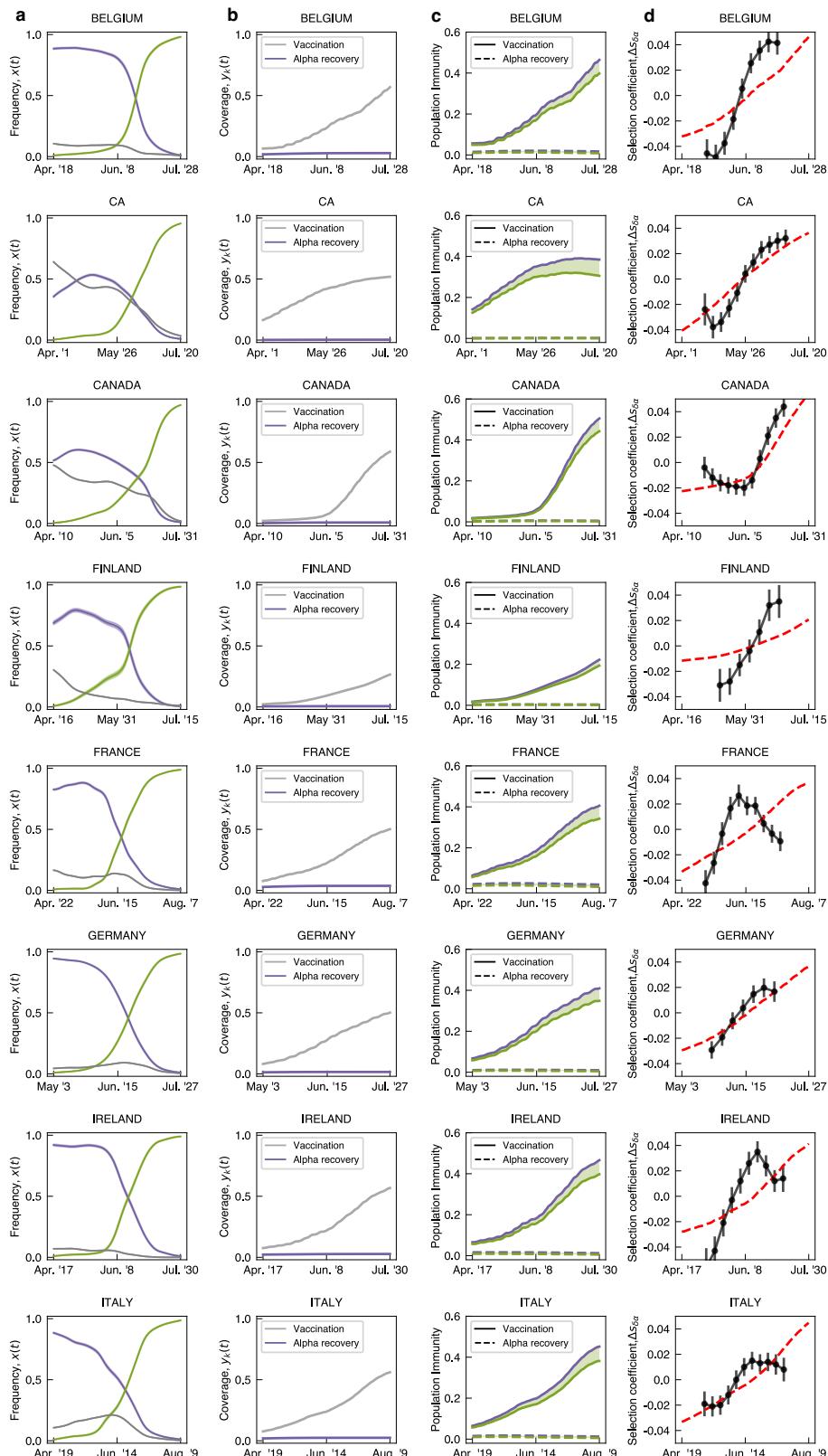
| model | antigenic parameters | | | posterior scores | |
|-------|-----------------------|-----------------|---------------|------------------|------------|
| | γ_{vac} | a | b | ΔL | ΔH |
| VI | 1.22 ± 0.03 | 0.24 ± 0.03 | 2.0 ± 0.5 | 947 | -1883 |
| V | 1.22 | 0.34 | - | 882 | -1758 |
| I | 4.9 | 0.21 | - | 378 | -750 |
| 0 | - | - | - | 0 | 0 |

We compare the full fitness model used in the main text (VI: vaccination + infection + intrinsic selection) with partial models (V: vaccination + intrinsic selection, I: infection + intrinsic selection) and a null model (0: intrinsic selection only). Columns from left to right: model parameters, γ_{vac} , a , b , ML values and 95% confidence intervals (definitions are given in Methods); log likelihood difference to the null model, ΔL ; BIC score difference to the null model, ΔH .

Table S3: Intrinsic and antigenic selection components.

| clade shift | selection coefficients | | | | | | | |
|-------------------|------------------------|----------------|-----------------|----------------|----------------|-----------------|----------------|---------------|
| | α | δ | $o(o1)$ | $o2$ | vac | bst | 0 | s |
| $1 - \alpha$ | < | - | - | - | < | - | $.08 \pm .001$ | $.08 \pm .01$ |
| $\alpha - \delta$ | $.01 \pm .001$ | < | - | - | $.04 \pm .002$ | - | $.05 \pm .002$ | $.09 \pm .02$ |
| $\delta - o$ | < | $.02 \pm .006$ | $-.01 \pm .002$ | - | $.06 \pm .01$ | $-.01 \pm .002$ | $.06 \pm .01$ | $.14 \pm .03$ |
| $o1 - o2$ | < | < | $.01 \pm .004$ | < | < | < | $.08 \pm .002$ | $.08 \pm .01$ |
| $o2 - o45$ | < | < | $.01 \pm .002$ | $.01 \pm .002$ | < | $.04 \pm .01$ | $.06 \pm .008$ | $.12 \pm .02$ |

Selection coefficients between the invading and the ancestral variant, $s = f_{\text{inv}} - f_{\text{anc}}$, and their decomposition into antigenic and intrinsic components are inferred for the full fitness model; all values are time averages for each clade shift. Rows from top to bottom: major clade shifts, $1 - \alpha$, $\alpha - \delta$, $\delta - o$; recent clade shifts, $o1 - o2$, $o2 - o45$ (shift incomplete, entries refer to initial period). Columns from left to right: average antigenic selection in immune channels $k = \alpha, \delta, o(o1), o2$, vac, bst; intrinsic selection (0); total selection (s). Selection coefficients are given in units [1/day]; the symbol “<” marks values $s < 0.01$. We list ML values with 95% confidence intervals (for selection components) or with rms cross-region variation of selection (for s ; cf. Fig. 1b).



(continued on next page)

Fig. S1: Empirical and model-based trajectories of the $\alpha - \delta$ shift. Evolutionary, epidemiological, and cross-immune trajectories are shown for all regions of this study. **(a)** Observed frequency trajectories of relevant clades, $x_i(t)$; rms sampling error is indicated by shading. **(b)** Cumulative coverage of primary vaccination, $y_{\text{vac}}(t)$ (light gray), and of booster vaccination, $y_{\text{bst}}(t)$ (dark gray); cumulative population fraction of α infections, $y_\alpha(t)$ (purple), and of δ infections, $y_\delta(t)$ (green). **(c)** Population immunity functions, $C_i^k(t)$ (as in Fig. 1c). **(d)** Empirical selection change, $\Delta\hat{s}(t)$ (dots, with rms statistical errors indicated by bars), together with ML model prediction, $\Delta s(t)$ (dashed line). Criteria for inclusion of regions are given in Methods.

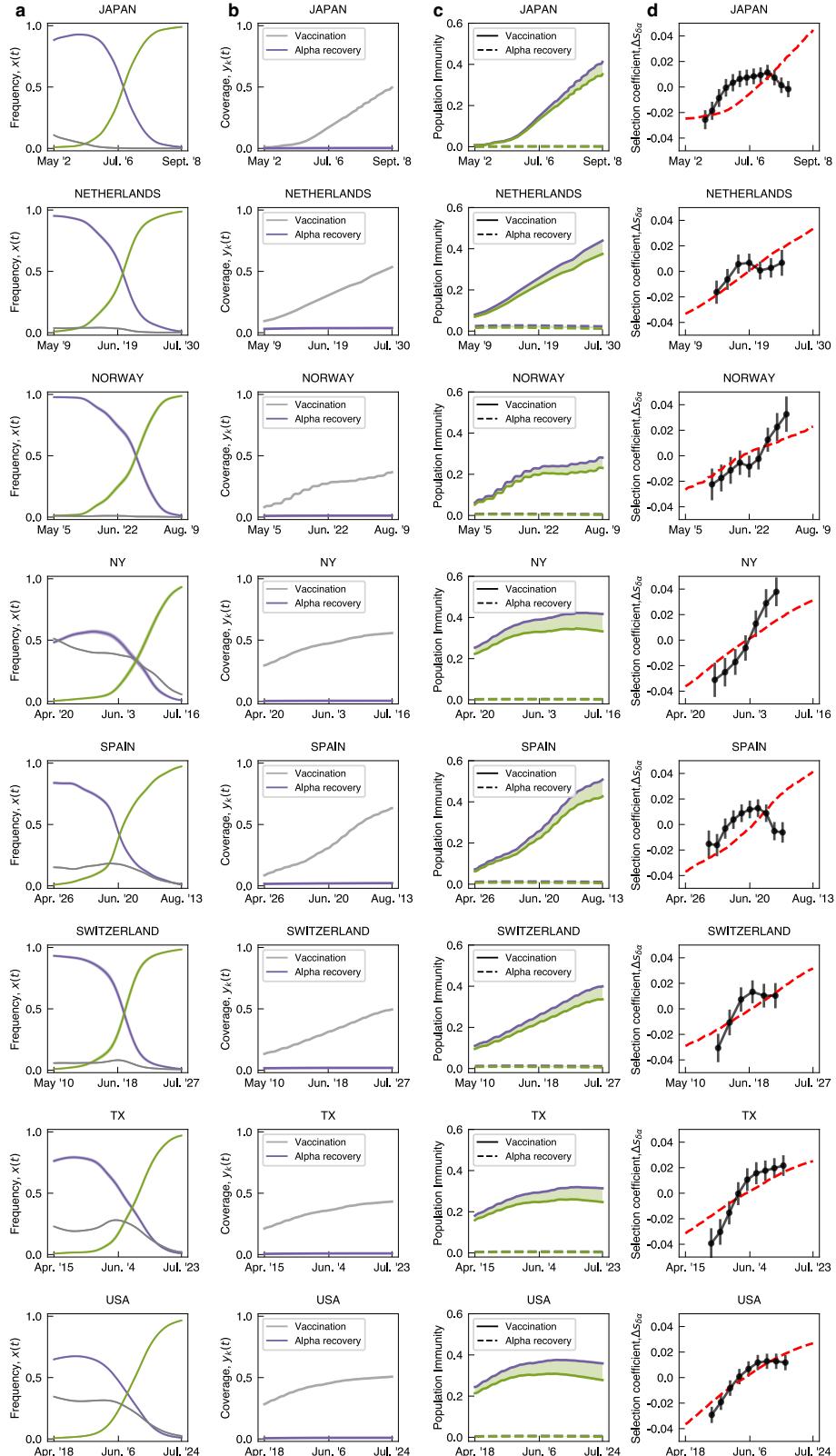
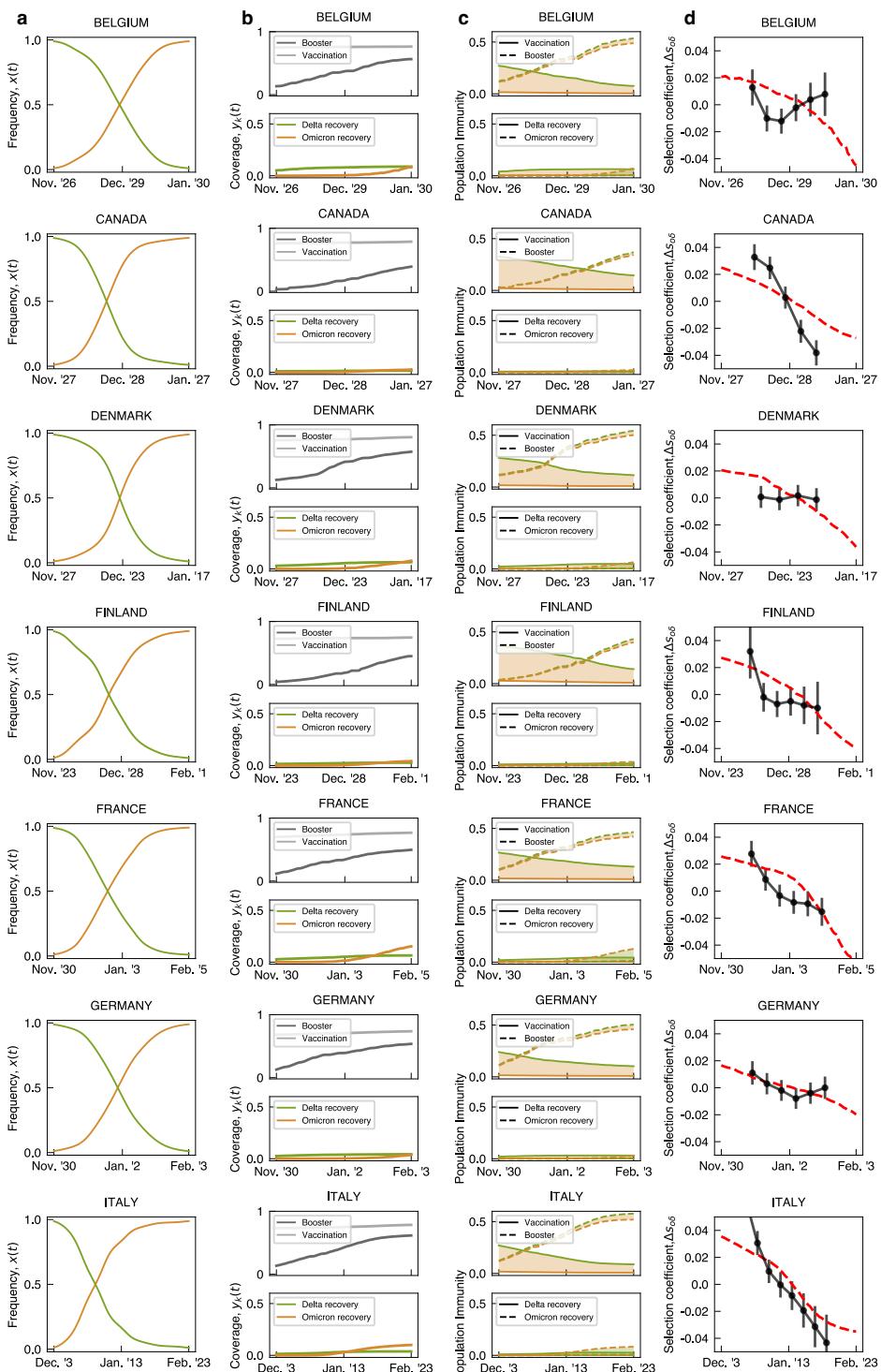


Fig. S1: (continued)



(continued on next page)

Fig. S2: Empirical and model-based trajectories of the $\delta - o$ shift. Evolutionary, epidemiological, and cross-immune trajectories are shown for all regions of this study. (a) Observed frequency trajectories of relevant clades, $x_i(t)$; rms sampling error is indicated by shading. (b) Cumulative coverage of primary vaccination, $y_{\text{vac}}(t)$ (light gray), and of booster vaccination, $y_{\text{bst}}(t)$ (dark gray); cumulative population fraction of δ infections, $y_\delta(t)$ (green), and of o infections, $y_o(t)$ (orange). (c) Population immunity functions, $C_i^k(t)$ (as in Fig. 1c). (d) Empirical selection change, $\Delta \hat{s}(t)$ (dots, with rms statistical errors indicated by bars), together with ML model prediction, $\Delta s(t)$ (dashed line). Criteria for inclusion of regions are given in Methods.

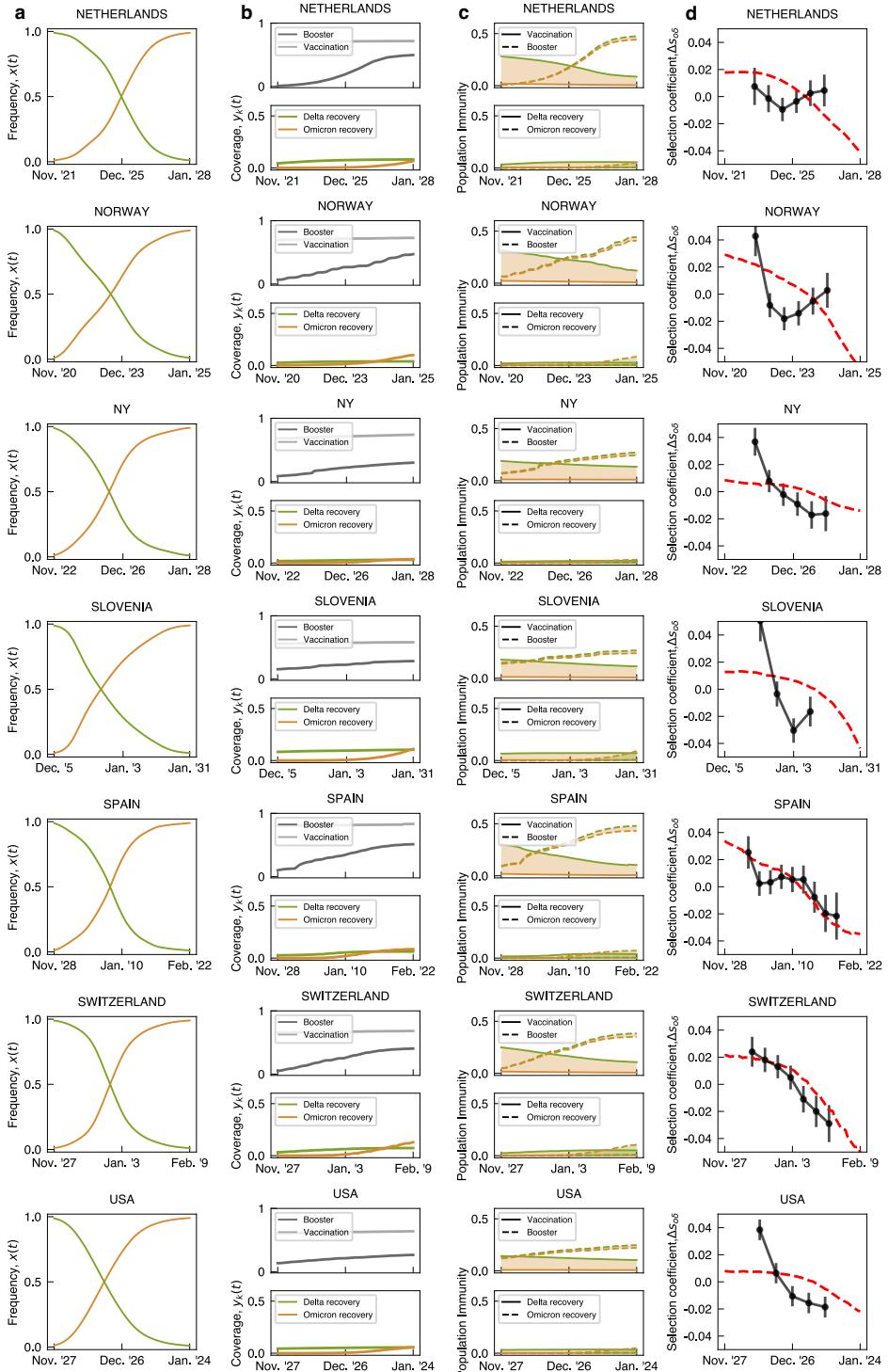


Fig. S2: (continued)

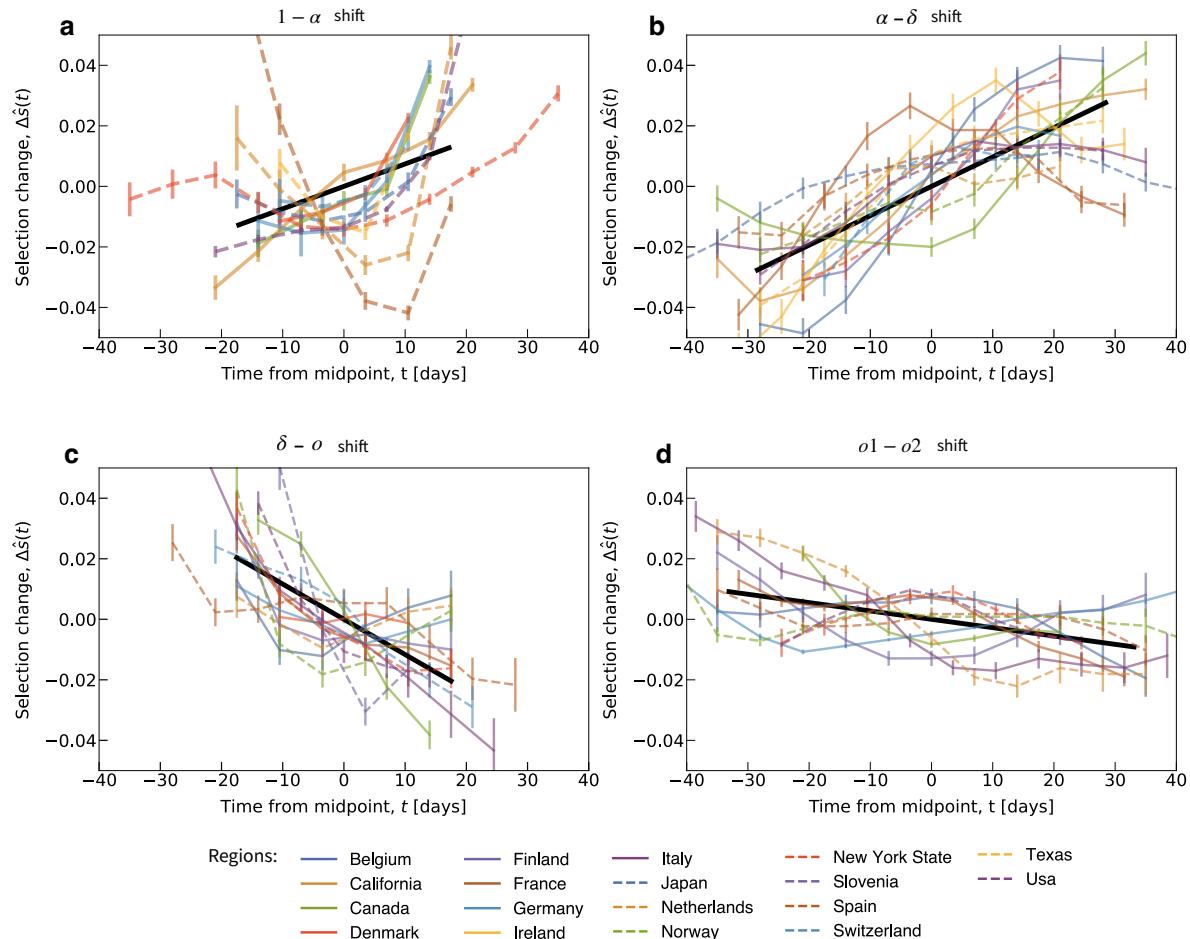


Fig. S3: Selection tracking in multiple regions and clade shifts. Empirical selection change between invading and ancestral clade, $\Delta\hat{s}(t) = \hat{s}(t) - \langle s \rangle$, for all complete clade shifts and all regions of this study (brackets denote time averages for each trajectory). Selection trajectories are derived from regional frequency trajectories and plotted against time counted from the midpoint (colored lines); rms statistical error is indicated by shading. Summary statistics: cross-region linear regression, $s_{lin}(t)$ (black solid line, length gives r.m.s. time span of trajectories). (a) $1 - \alpha$ shift: small, statistically insignificant time dependence, $\text{Var}(s_{lin}) = 3.6 \times 10^{-4}$, $P > 0.01$; (b) $\alpha - \delta$ shift: substantial, statistically significant time dependence, $\text{Var}(s_{lin}) = 2. \times 10^{-3}$, $P < 10^{-16}$; (c) $\delta - o$ shift: substantial, statistically significant time dependence, $\text{Var}(s_{lin}) = 1.4 \times 10^{-3}$, $P < 10^{-5}$; (d) $o1 - o2$ shift: small, but statistically significant time dependence, $\text{Var}(s_{lin}) = 2.9 \times 10^{-4}$, $P < 10^{-4}$. All P values are computed using a two-sided Wald test. The statistical grading of shifts is described and criteria for inclusion of regions are given in Methods.

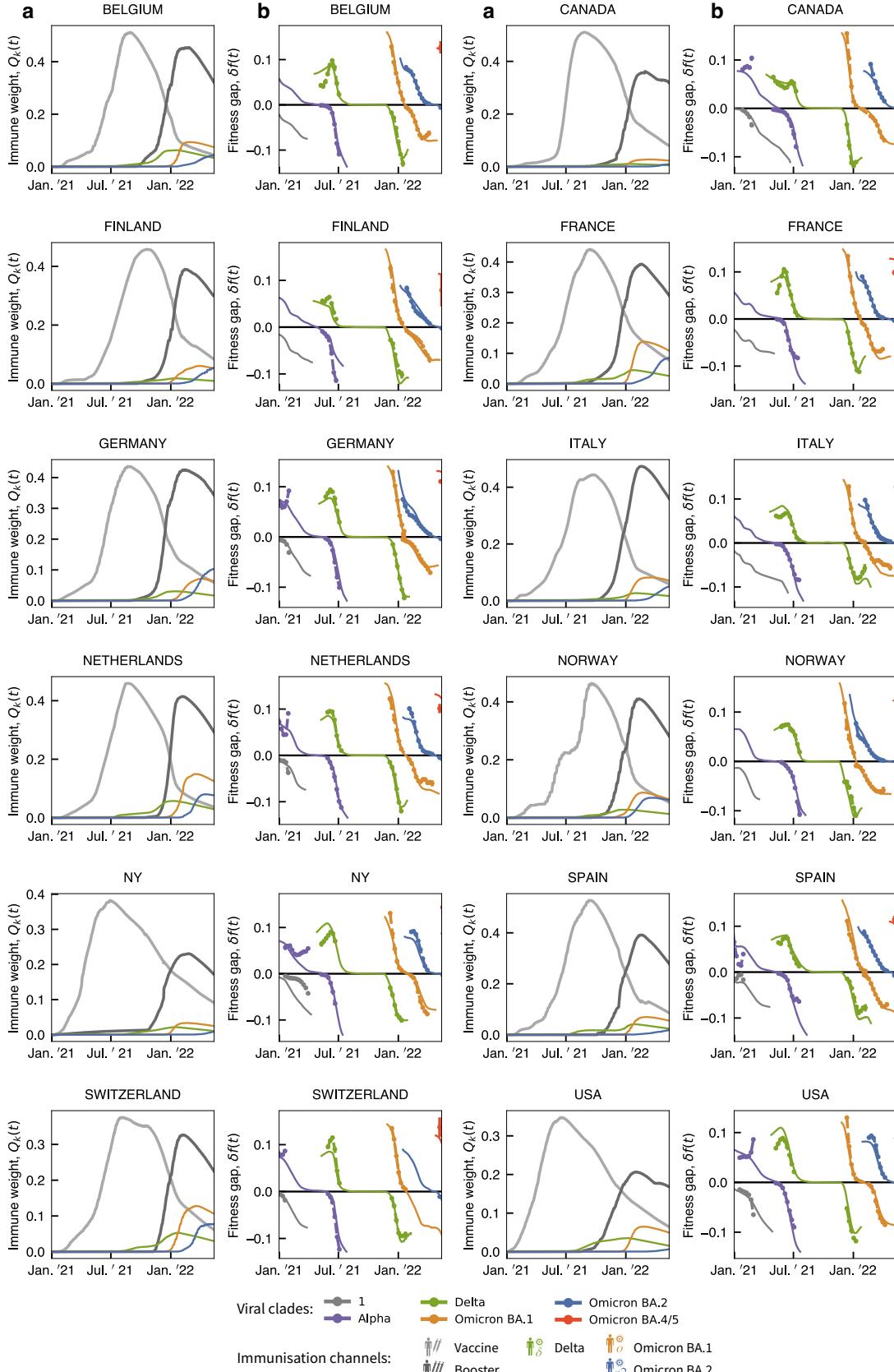


Fig. S4: Regional long-term trajectories of immune weight and fitness. (a) Time-dependent weight factors of different immune classes, $Q_k(t)$. (b) Time-dependent fitness gap, $\delta f_i(t)$. Criteria for inclusion of regions are given in Methods; see Fig. 4 for averaged trajectories.