

1 **Investigating the evolutionary origins of the first three SARS-CoV-2 variants of**
2 **concern**

3

4 Mahan Ghafari¹, Qihan Liu², Arushi Dhillon², Aris Katzourakis^{1*}, Daniel B Weissman^{2*}

5

6 ¹Department of Zoology, University of Oxford, Oxford, UK

7 ²Department of Physics, Emory University, Atlanta, Georgia, US

8

9 *Correspondence: daniel.weissman@emory.edu (DBW); aris.katzourakis@zoo.ox.ac.uk (AK)

10

11 **Abstract**

12 The emergence of Variants of Concern (VOCs) of SARS-CoV-2 with increased transmissibility,
13 immune evasion properties, and virulence poses a great challenge to public health. Despite
14 unprecedented efforts to increase genomic surveillance, fundamental facts about the
15 evolutionary origins of VOCs remain largely unknown. One major uncertainty is whether the
16 VOCs evolved during transmission chains of many acute infections or during long-term infections
17 within single individuals. We test the consistency of these two possible paths with the observed
18 dynamics, focusing on the clustered emergence of the first three VOCs, Alpha, Beta, and Gamma,
19 in late 2020, following a period of relative evolutionary stasis. We consider a range of possible
20 fitness landscapes, in which the VOC phenotypes could be the result of single mutations, multiple
21 mutations that each contribute additively to increasing viral fitness, or epistatic interactions
22 among multiple mutations that do not individually increase viral fitness—a “fitness plateau”. Our
23 results suggest that the timing and dynamics of the VOC emergence, together with the observed
24 number of mutations in VOC lineages, are in best agreement with the VOC phenotype requiring
25 multiple mutations and VOCs having evolved within single individuals with long-term infections.

26 Introduction

27 For the first 8 months of the SARS-CoV-2 pandemic, the virus exhibited a very slow pace of
28 adaptation, with D614G being the only persistent adaptive substitution that appears to have
29 resulted in an increased transmissibility of the virus [1-3]. However, during the second half of
30 2020, three designated variants of concern (VOCs) of SARS-CoV-2, Alpha, Beta, and Gamma,
31 emerged independently and in quick succession [4-6]. No other VOC emerged until Delta and
32 Omicron in 2021 which appear to be very different, both genetically and phenotypically, from
33 the three original VOCs [7, 8]. The VOCs are characterised by a large number of mutations relative
34 to the genetic background from which they first emerged, and exhibit altered phenotypes
35 resulting in varying combinations of increased transmissibility, virulence, and immune evasion
36 [6, 9-11].

37 Phylogenetic analyses show that a large number of mutations, mostly located in the spike
38 protein, have independently evolved in multiple lineages of SARS-CoV-2 including the Alpha, Beta
39 and Gamma variants and are likely playing a key role in the adaptive evolution of the SARS-CoV-
40 2 [7, 12]. Experimental measurements and molecular dynamics simulations also show that some
41 of these mutations have synergistic interactions for important functional traits [13, 14], indicating
42 that they may have greater combined fitness benefit to the virus. Some of the distinctive
43 mutations in the VOCs, including the E484K and N501Y mutations found in the first three VOCs,
44 have also been observed in chronic infections such as those in certain immunocompromised
45 individuals [15-17], suggesting that the VOCs may have arisen from such infections. Some of the
46 other possible explanations for the emergence of VOCs include prolonged circulation of the virus
47 in areas of the world with poor genomic surveillance or reverse-zoonosis from other animals such
48 as rodents followed by sustained transmission and adaptive evolution within the animal
49 population and a spill over back to the humans (see [18] for a recent review on the possible
50 origins of variants of SARS-CoV-2).

51 While finding the evolutionary process(es) that may have led to the emergence of VOCs has
52 profound consequences for understanding the fate of the SARS-CoV-2 pandemic, there have
53 currently been no systematic investigations to assess the likelihood of any particular evolutionary
54 pathway that would lead to the emergence of VOCs. In this work we investigate whether the
55 emergence of VOCs was the result of evolution via sustained transmission chains between
56 acutely infected individuals or prolonged infections, and evaluate plausible fitness landscapes.

57 We also discuss the potential implications of our results for the future of the pandemic and
58 potential measures that might lower the rate at which new VOCs emerge.

59 **Results**

60 ***Emergence of VOCs: an evolutionary puzzle***

61 The Alpha, Beta, and Gamma VOCs arose independently and in quick succession, with several
62 shared mutations, in three different countries and began to spread globally (**Figure 1**). This long
63 waiting time followed by clustered emergence of a handful of lineages was not predicted by any
64 simple evolutionary theories. Typically, one would assume that either the beneficial mutation
65 supply is small, in which case one expects a long waiting time for the first VOC but also long gaps
66 before subsequent VOCs, or the mutation supply is large, in which case one expects many VOCs
67 with only a short waiting time [19]. Moreover, each VOC had >6-10 mutations distinguishing it
68 from then-dominant genotypes, which was also unexpected. One of the key evolutionary
69 questions is whether VOCs evolved over the course of many acute infections or within single
70 chronic infected hosts. Both possibilities have serious issues. The many-acute-infections
71 hypothesis needs to explain how the virus acquired so many changes, as the mutant lineages
72 would have had to remain at frequencies below the detection threshold in different countries for
73 several months. The chronic-infection hypothesis needs to explain both why adaptation to the
74 within-host environment led to a transmission advantage between hosts, and why there was no
75 'leakage' of some intermediate mutations at the between-host level before the emergence of
76 the VOCs, i.e., why genotypes with some of the VOC mutations did not escape from the
77 chronically infected patients earlier.

78 ***Between-host model of VOC emergence***

79 We assume the effective virus population size is $N_e = N/\sigma^2$ where N is the number of infectious
80 individuals worldwide and σ^2 is the variance in offspring number (secondary cases). We treat
81 each acute infection as one generation, assuming a tight transmission bottleneck of a single virion
82 [20-22]. Viruses mutate at rate μ per base per generation (see **Methods** section). For a mutant
83 virus population with selective advantage s relative to the background, the average number of
84 secondary cases increases by a factor $1+s$. We also assume that the number of secondary cases
85 approximately follows a negative binomial distribution with mean R_t and dispersion parameter
86 k , so that $\sigma^2 \approx R_t(1+R_t/k)$. There is substantial uncertainty in the amount of overdispersion in the
87 pandemic, and consequently similar uncertainty in the effective population size. Therefore, we
88 consider a range of values for k to see if any would be consistent with the observed dynamics of

89 the VOC emergence. We also note that while the importance of spatial structure is clearly visible
90 in the spatially restricted initial spread of the VOCs from real-world data, we expect that we can
91 neglect it when analysing their *emergence*. This is because spatial structure should not have a
92 large impact on viral dynamics until a lineage becomes locally common, and the specific
93 mutations differentiating the VOCs were all locally rare prior to their emergence.

94 ***Within-host model of VOC emergence***

95 Unlike tracking the between-host evolution of SARS-CoV-2 where an unprecedented effort has
96 led to huge numbers of consensus genome sequences [23], our current knowledge of the within-
97 host evolutionary dynamics of SARS-CoV-2 is still very limited, particularly in those with chronic
98 infections. Because there is very limited data with which to constrain the within-host
99 evolutionary dynamics of chronic infections with SARS-CoV-2, we simply treat it as a 'black box'
100 and assume with some probability, P_f , that a new infection is chronic and may lead to the
101 production of a VOC (**Table 1; Methods** section). We also assume that within-host substitutions
102 required for the production of the VOC occur at a constant rate μ_c per generation (see **Table 1**).
103 (Here a generation is still defined as the typical length of an *acute* infection.) Given that we know
104 only three VOC lineages emerged by late 2020, we expect $T_{obs} N P_f \sim 3$ where $T_{obs} \sim 180-317$ days is
105 the expected time to the emergence of the first VOC since the beginning of the pandemic based
106 on phylogenetic estimates (see **Table 1**). Therefore, given the typical variation in the population
107 size throughout the pandemic for biologically relevant parameter combinations $N \sim 1 \times 10^6 - 1 \times 10^7$,
108 we expect that values of $P_f \sim 5 \times 10^{-9} - 1 \times 10^{-7}$ will maximize the likelihood of the within-host model
109 and focus on these.

110 ***Fitness landscapes***

111 One possible explanation for the temporal clustering of VOCs with large numbers of mutations is
112 that the underlying fitness landscape may have some structure that causes the dynamics to
113 deviate from our usual expectations. Unfortunately, the full space of possible fitness landscapes
114 is enormous and impossible to explore exhaustively. To investigate the possible effects of the
115 landscape on the dynamics, we therefore focus on three limiting local fitness landscapes that
116 span a range of biologically plausible scenarios (**Figure 1A**). Importantly, these landscapes
117 describe only between-host fitness, which could be very different from within-host fitness. As
118 mentioned above, we treat within-host dynamics implicitly using an effective substitution rate
119 and so do not need an explicit fitness landscape for it. In all three landscapes, the peak is a VOC
120 phenotype with fitness advantage s over the ancestor. We assume that Alpha, Beta, and Gamma

121 are similar enough that they can be approximately described by the same landscape and the
122 same value of s , which we infer from the early rate of increase of the VOCs (see **Methods**).
123 Landscape 1 is the simplest possibility: a single mutation on the ancestral background is sufficient
124 to confer the full advantage. In Landscape 2, we test whether simply increasing the number of
125 mutations involved can explain the temporal clustering. In this landscape, the VOC phenotype is
126 produced by a combination of $K > 1$ mutations, each providing an independent fitness benefit
127 s/K . In Landscape 3, we test whether epistasis may have an effect: the VOC phenotype again
128 requires K mutations, but we now assume that they provide no fitness benefit until the full
129 combination is acquired, i.e., the population must cross a fitness plateau. As mentioned above,
130 there is experimental evidence for this form of epistasis among the VOC mutations [13, 14]. We
131 expect that shallow fitness valleys will produce similar dynamics to Landscape 3, as will shallow
132 upward slopes with a large jump in fitness at the end [24]. Note that mutations in all the three
133 landscapes can be acquired via the between- or within-host evolutionary pathways (**Figure 1B**).
134 For each evolutionary scenario, we test whether there are parameter values consistent with the
135 data on the timing of the emergence of Alpha, Beta, and Gamma variants of SARS-CoV-2 (see
136 **Methods; Table 1**). For these parameter values, we further investigate whether they correspond
137 to biologically reasonable scenarios in terms of the frequencies of the intermediate mutations
138 prior to the emergence of VOCs, total number of mutations required to produce VOCs, total
139 number of successful VOC lineages produced over time, and the timing between the emergence
140 of different VOC lineages.

141 *Landscape 1: single mutations*

142 We start with the simplest possible fitness landscape, in which a single mutation conferring a
143 fitness advantage s relative to the genetic background of circulating lineages is required for the
144 emergence of VOCs. We first consider the between-host evolutionary pathway. As long as the
145 effective population size of the pandemic was not much smaller than the census size (i.e.,
146 overdispersion was not too large), the mutation supply $N_e \mu$ became large early in 2020. At this
147 point, numerous lineages would have emerged over a short period of time (see the $k=0.2$ scenario
148 in **Figure 3A**), inconsistent with the observed dynamics. We can therefore rule out this scenario.

149 If overdispersion were very large, it could have kept $N_e \mu$ low through the establishment of the
150 VOCs (see the $k=0.005$ and 0.001 scenarios in **Figure 3A**). **Figure 3** shows that under extremely
151 high levels of overdispersion ($k=0.005$ and 0.001) this model can match the long waiting time for
152 the emergence of the first VOC. However, such high levels of overdispersion are not supported
153 by any existing epidemiological studies on SARS-CoV-2 transmission [25]. Moreover, **Figure 3B**

154 shows that this model rarely produces an evolutionary dynamics that would fit the joint waiting
155 time distribution for all three VOCs (also see **Supplementary Figure 1**). Under these mutation-
156 limited conditions, there is an approximately exponential waiting time for the arrival of each VOC
157 lineage (once we reach the point where COVID-19 becomes a pandemic in March 2020). Thus, it
158 predicts similarly long waiting times for the emergence of Alpha, Beta, and Gamma, inconsistent
159 with the observed temporal clustering. Therefore, there is no biologically reasonable
160 combination of parameters that result in the clustered emergence of VOCs in late 2020 via the
161 Landscape 1 between-host evolutionary pathway.

162 On the other hand, if VOCs arose from chronic infections, then their emergence was a two-step
163 process: first, chronic infections had to occur, and then the VOC mutation had to arise in them.
164 The waiting time for the first step is determined by NP_f ; note that the number of chronic
165 infections depends on the census size N rather than N_e , i.e., it is insensitive to the amount of
166 overdispersion. The second step follows an exponential distribution within each chronic host,
167 with rate μ_c . The third step, the spread of the VOC from the original chronic host to the rest of
168 the population, then takes much less time than the first two. **Figure 4** shows that to match
169 observed VOC dynamics we must assume that the level of overdispersion is very high (i.e., very
170 low mutation supply, $N_e\mu$), effectively blocking the between-host evolutionary pathway, while
171 simultaneously assuming that chronic infections are very frequently produced in the population
172 (i.e., $NP_f \sim 1$) and that there is a relatively long waiting time before the production of each VOC
173 mutation ($\mu_c \sim 0.01$). However, like the between-host pathway, this scenario requires very high
174 levels of overdispersion which makes the Landscape 1 within-host evolutionary pathway also an
175 unlikely explanation for the emergence of VOCs (see **Supplementary Figure 2**).

176 *Landscape 2: additive mutations*

177 Landscape 2 corresponds to an evolutionary pathway in which there were $K > 1$ major mutations
178 involved in the emergence of VOCs, each making an additive contribution of $\approx s/K$ to fitness. If
179 evolution occurred at the whole-population level, **Figure 5** and **Supplementary Figure 3** show
180 that, for a range of parameter combinations, the additive fitness landscape requiring up to four
181 mutations can create evolutionary dynamics with appropriately long waiting times before the
182 arrival of the first successful VOC lineage, while for combinations of more than 4 mutations, VOC
183 lineages do not emerge by late 2020 under any biologically reasonable parameter combinations
184 for effective population size, mutation rate, and selective coefficient. However, while $K \leq 4$ can
185 match the observed waiting for the first VOC lineage, for the $K=2$ and 3, this first VOC is usually
186 followed by the establishment of nearly a dozen VOC lineages that emerge in quick succession

187 (see $K=2$ and 3 scenarios in **Figure 5A**; also see **Supplementary Figure 3**), inconsistent with the
188 observation of only 3 VOC lineages emerging in late 2020. However, while for $K=4$ fewer VOC
189 lineages are produced, a closer examination of a typical evolutionary trajectory that matches the
190 long waiting time before the establishment of the first VOC further reveals that the intermediate
191 single-, double-, or triple-mutants reach high frequencies before the emergence of the first
192 successful (quadruple-mutant) VOC lineage (**Figure 5C**). The sequential fixation of adaptive
193 mutations at the population level would imply that the intermediate mutations were detectable
194 many months prior to the emergence of VOCs, again inconsistent with the genomic surveillance
195 data from around the world. The inconsistency is also visible phylogenetically. The sequential
196 fixation dynamics predicted by the model create a ladder-like phylogenetic relationship between
197 the background and mutant populations whereby every new VOC mutation becomes dominant
198 in the population before giving rise to lineages with additional mutations. Even though such
199 phylogenetic relationships may emerge in SARS-CoV-2 over longer evolutionary timescales (as
200 have been observed in human coronaviruses [26]), they do not resemble the observed topology
201 of the phylogeny of the VOCs of SARS-CoV-2, which is more star-like.

202 For the chronic-infection pathway, on the other hand, the intermediate mutants could have fixed
203 within the host while remaining at undetectable frequencies at the between-host level until the
204 production of the VOCs. **Figure 6** shows that for a combination of parameters requiring $K=3$ and
205 6 mutations where the mutation supply is low and the strength of selection is relatively weak
206 such that the intermediate mutants cannot reach fixation before the emergence of the VOC
207 population, the Landscape 2 within-host pathway can lead to the clustered emergence of a few
208 VOC lineages by late 2020. However, if the selective coefficient s/K on single mutants is too high,
209 they will reach observable frequencies before the VOCs emerge, as we discussed above with the
210 between-host pathway. Effectively, this means that there is a minimal K of at least 3 needed so
211 that the strength of selection on each mutant allele is not too strong. Alternatively, lower K is
212 possible but requires extremely large overdispersion, as in the $K = 1$ case.

213 *Landscape 3: fitness plateau crossing*

214 As in Landscape 2, Landscape 3 describes an evolutionary pathway where there are $K>1$ major
215 mutations involved in the generation of VOCs, but in this case, only the full K -mutant VOC
216 genotype has a substantial selective advantage relative to the background population, while the
217 selective advantages of the intermediate genotypes are negligible. This does not necessarily
218 imply that the selective coefficients of the intermediate genotypes are small in the standard
219 weak-selection sense (small relative to $1/N_e$), but only that they are too small to substantially

220 affect the dynamics of the production of the first successful K -mutant VOC lineage, a weaker
221 condition that depends on the mutation rate [24].

222 For the between-host model of VOC emergence, our analysis suggests that only a plateau-
223 crossing of size $K=2$ may be consistent with the timing of the emergence of SARS-CoV-2 VOCs
224 (**Figure 7; Supplementary Figure 5**). Extended plateaus requiring $K>2$ mutations take much
225 longer to cross and for most parameter combinations either zero or one VOC lineage is produced
226 before the end of 2020 (**Figure 7A**). For a typical $K=2$ plateau-crossing trajectory, single-mutant
227 genotypes grow linearly over time and reach a frequency of $\approx 0.1\%$ before producing ~ 1 - 5
228 successful VOC lineages that emerge in quick succession (**Figure 7C**). Therefore, unlike the
229 between-host evolutionary pathway in Landscape 2, a fitness plateau could have led to the
230 clustered emergence of several VOCs after a long waiting time during which none of the
231 intermediate mutations reached high frequency. However, the fact that for biologically plausible
232 parameter values only a narrow plateau of $K=2$ mutations can be crossed seems inconsistent with
233 the high number of mutations found in the VOCs and particularly with the high number of similar
234 mutations shared across unrelated VOC lineages. This inconsistency may be partly reconciled
235 with the possibility of compounded evolutionary effects following the plateau-crossing event
236 such as the emergence of hyper-mutability traits across certain sites or strong within-host
237 selection following the acquisition of the K mutations.

238 If the VOCs arose from chronic infections, the intermediate VOC mutations (which are neutral at
239 the between-host level of selection, but may be selected within-host) can rapidly fix within a
240 host, allowing much wider plateaus to be crossed compared to the between-host evolutionary
241 pathway. Unlike Landscape 2 within-host pathway, the early leakage of intermediate mutations
242 to the population is much less likely as they have no strong selective advantage over the
243 background population. **Figure 8** shows that the within-host evolutionary pathway of Landscape
244 3 creates evolutionary trajectories that are consistent with the clustered emergence of ~ 3 VOCs
245 in late 2020. There is also less seeding of new chronic infections with intermediate mutations,
246 leading to fewer VOC lineages compared to Landscape 2 (also see **Supplementary Figure 6**).

247 **Discussion**

248 The global spread of the Omicron variant of SARS-CoV-2 has given a renewed attention to the
249 underlying evolutionary mechanisms that lead to the emergence of VOCs. Practically, we would
250 like to know whether to expect future VOCs to arise, and if so when and whether there will be
251 early warning signs. Answering this question is not only important for understanding the fate of

252 the pandemic but also may have major public health implications for how to best develop
253 strategies for controlling the spread of the disease. In this study, we provided a quantitative
254 framework for investigating the likelihood of different evolutionary pathways that can give rise
255 to VOCs of SARS-CoV-2. We found that VOCs are unlikely to be driven by a single adaptive change
256 at the population level as this would require significantly high levels of overdispersion which is
257 not supported by any existing epidemiological study on SARS-CoV-2 transmission [25]. We also
258 showed that if multiple VOC mutations combine additively for advantage, they can only emerge
259 on the background of a chronic infection, otherwise individual VOC mutations would reach high
260 frequencies from the early stages of the pandemic and, therefore, would have been picked up
261 from genomic surveillance data. If individual VOC mutations were acquired during chronic
262 infections and had a strong advantage relative to the then-dominant genotypes, they may have
263 still been leaked to the population at large before the emergence of VOCs. Therefore, we showed
264 that only additive mutations with relatively small fractional contribution to VOC fitness may yield
265 evolutionary dynamics that resembles the clustered emergence of SARS-CoV-2 VOCs in late 2020.
266 On the other hand, we showed that cryptic circulation of a mutant lineage for sustained periods
267 of time before producing VOCs is possible via a fitness plateau-crossing landscape. While at the
268 between-host level such a landscape may not yield more than 2 mutations in excess of the
269 background population over a period of 7-12 months under biologically relevant parameter
270 combinations, many more mutations can be accumulated during a chronic infection, for example
271 such as those found in certain immunocompromised individuals, without ever being leaked to
272 the rest of the population. We found that the pattern of the timing of VOC emergence via the
273 fitness plateau-crossing landscape under both the within- and between-host pathways are
274 aligned with the timing of the clustered emergence of Alpha, Beta, and Gamma variants in late
275 2020.

276 Finally, it is important to note that in all of the within-host evolutionary pathways (i.e.,
277 Landscapes 1, 2, and 3), we found parameter combinations that can re-create the clustered
278 emergence of the first three SARS-CoV-2 VOCs. In particular, we showed that either because of
279 having very few mutations that are selectively beneficial at the population level (i.e., Landscape
280 1) or the low prevalence of intermediate mutations before the emergence of the VOCs (i.e.,
281 Landscapes 2 and 3), we would expect the phylogenetic relationship between the VOC lineages
282 and background populations to manifest itself with a long evolutionary distance branch
283 connected to deeper internal nodes of the tree with each VOC clade independently emerging
284 from a unique genetic background and be subsequently replaced by another VOC clade from an

285 entirely different background (**Figure 9**). This creates a phylogenetic relationship between VOC
286 clades that is similar to what we observe for Alpha, Beta, and Gamma variants [4-6].

287 ***Cryptic transmission of VOCs in humans***

288 Another possibility for why the VOCs were not detected until mid to late 2020 is that they may
289 have been circulating cryptically in areas of the world with poor genomic surveillance before
290 becoming globally dominant. While variants of SARS-CoV-2 with multiple spike mutations have
291 been detected through travel surveillance from passengers travelling from areas with little to no
292 genomic surveillance [27, 28], if they were highly transmissible variants and had a potential to
293 become a VOC, given the interconnectivity of the human interactions, it should not take very long
294 before they become globally dominant. Therefore, as we showed in our analysis of Landscape 1,
295 this scenario of VOC emergence seems to be only possible under significant levels of
296 overdispersion such that it prohibits the selectively beneficial mutations from immediately taking
297 off globally relative to other variants.

298 ***Possibility of reverse zoonosis***

299 A somewhat similar idea to the cryptic transmission of variants in human populations is the
300 possibility of a lineage (or multiple lineages) of SARS-CoV-2 jumping from humans to other
301 mammals such as white-tailed deer, mink, hamster, and mouse where they circulate and evolve
302 without being detected for a relatively long period before they jump back to the human
303 population [29-32]. In particular, some recent studies have reported the detection of multiple
304 spillovers of SARS-CoV-2 from humans and onward transmission in deer population with highly
305 divergent genomes being detected in deer population with potential deer-to-human
306 transmission [33, 34]. However, the genomic composition of these divergent genomes in deer
307 population are different from the VOCs with a much lower ratio of non-synonymous to
308 synonymous changes which suggests they may be following a completely different evolutionary
309 path. Nevertheless, these studies indicate that it is possible for a highly divergent set of genomes
310 to evolve in another species with a potential for deer-to-human transmission without ever being
311 detected. Mink and hamster sequences offer some of the more compelling examples of
312 transmission from humans to a non-human species and back, supported by phylogenetic
313 evidence [30, 31]. None of the currently identified sequences from animals appear as sister taxa
314 to any of the circulating VOCs. While we cannot rule out evolution in an animal reservoir, one
315 might expect the contribution of animals to human transmission chains to be dwarfed by the
316 amount of human-to-human transmission currently happening.

317 ***Role of recombination***

318 Recombination can bring together mutations from different backgrounds, potentially expediting
319 the rate of adaptation by creating viable and more pathogenic hybrid new variants of a pathogen.
320 Coronaviruses are also known to recombine with one another during mixed infections [35, 36].
321 While during the early stages of the pandemic SARS-CoV-2 sequences typically differed by only a
322 handful of mutations from each other thereby making the effects of recombination
323 indistinguishable from those of recurrent mutation [37], as more viral genetic diversity built-up
324 in the population, the generation and transmission of interlineage recombinants of SARS-CoV-2
325 in humans were reported in multiple studies [38, 39]. Even though there is currently no definitive
326 evidence for recombination being involved in the emergence of VOCs, including the Omicron
327 variant of SARS-CoV-2 [40], we would expect the role of recombination to be more pronounced
328 as the virus continues accumulate more genetic diversity by sustained circulation around the
329 world.

330 ***Shifting landscape***

331 We have assumed a static fitness landscape prior to the emergence of the first three VOCs; here
332 we consider the plausibility of that assumption. During the first year of the pandemic, a novel
333 virus was spreading in an immunologically naïve population [12]. As more individuals became
334 infected and developed natural immunity, it is possible that the fitness landscape for the virus
335 shifted as selection for immune escape increased [41]. However, by the time the first three VOCs
336 emerged in late 2020, the majority of the world's population were still susceptible to the disease
337 and may not have even been exposed to it. Therefore, it is unlikely that the build-up of natural
338 immunity alone was the reason behind their increased selective advantage. In contrast, the
339 global dominance of Omicron in late 2021 was largely due to its immune escape properties
340 relative to previous variants of SARS-CoV-2 [42] and, therefore, its emergence was likely the
341 result of a changing viral fitness landscape.

342 ***Future VOCs***

343 Emergence of new VOCs with increased transmissibility, immune evasion properties, and
344 virulence poses a great challenge to managing SARS-CoV-2. Based on our findings, one of the
345 major implications of chronic infections being the main source of generating VOCs is that finding
346 and treating chronic infections should be a top priority, not just for the benefit of chronically ill
347 patients but also from a public health standpoint. One of the main challenges with assessing the
348 likelihood of VOC emergence during chronic infections would be to quantify the prevalence of

349 immunosuppressed individuals within a population and determine which forms of
350 immunosuppression are associated with chronic infections.

351 Several studies have now shown evidence of recurrent SARS-CoV-2 mutations in
352 immunocompromised patients [15-17], with some suggesting the detection of a variant-like
353 lineage which arose from a chronic infection that spilled over into a local population [43]. Another
354 major implication of our work is that we can now quantitatively explain the possibility of such
355 events and find the expected time that it takes for a new VOC to emerge from a within-host
356 evolutionary pathway that involves any number of mutations. We showed that a typical within-
357 host plateau-crossing or additive mutation pathway involving 3-6 mutations requires a within-
358 host fixation rate of $\mu_c \sim 0.1-0.3$ per generation which corresponds to a period of 50-300 days
359 since the start of a chronic infection. The timing of such an event aligns with the time frame over
360 which some of the major mutations involved in VOCs have been observed in patients with chronic
361 infections [15-17]. This also implies that if a VOC emerges from the within-host evolutionary
362 pathway, it is more likely to reflect the genetic diversity of the virus population from several
363 months ago. It can explain why, for instance, the Omicron variants were not descendents of
364 Delta, which was the most prevalent variant at the time of emergence of Omicron. It also
365 suggests that while the next VOC could emerge from the prevalent Omicron background, it could
366 also come from, e.g., a chronic infection with Delta that started prior to the Omicron wave. Some
367 of the key remaining questions involve how much more of the fitness landscape the virus will be
368 able to explore as more chronic cases accumulate and existing chronic cases last longer. For
369 instance, if it has already crossed a 6-mutant fitness-plateau, how much longer would it take to
370 explore 7-mutant fitness plateaus?

371 **Methods**

372 *Between-host model of VOC emergence*

373 *Effective population size*

374 We approximate the between-host evolution of SARS-CoV-2 as a haploid population of size $N(t)$
375 which is equal to the number of daily infectious individuals with SARS-CoV-2 worldwide. Since
376 the number of confirmed cases is often a significant underestimation of the true number of
377 infections (e.g., see: [44, 45]), we use the number of daily confirmed deaths [46] to back-calculate
378 the number of infectious individuals, $N(t)$, from the global median infection fatality rate (IFR) of
379 COVID-19 [47]. We note this approach is still subject to several potential sources of bias including
380 variation in IFR over time (e.g. due to various pharmaceutical interventions) and across different
381 demographics [48]. Using confirmed COVID-19-related deaths may still underestimate the true
382 number of deaths associated with the disease due to under-reporting of deaths particularly in
383 areas of the world where there is limited testing from suspect cases [49]. Nevertheless, by
384 allowing for a wide variation in global IFR (from 0.2% to 1.5%), we can capture most of the
385 uncertainty in the number of infectious individuals worldwide. We also note that for the
386 timespan of interest in our work (i.e., start of the pandemic until the emergence of the first three
387 VOCs), the impact of pharmaceutical interventions such as vaccination on lowering the global IFR
388 is likely to have been negligible given that vaccination campaigns mostly started in 2021. The
389 confirmed global deaths started being reported from 2020-01-23. Assuming a 20-day delay from
390 the onset of symptoms to death [50], we set 2020-01-03 as the first timepoint in the simulation.

391 *Advantage of mutants*

392 The selective coefficient of a mutant individual depends on its number of mutations and the
393 fitness landscape (see **Figure 2A**). For instance, in the case of an additive fitness landscape of size
394 $K=3$, the fitness advantage of the single-, double-, and triple-mutants are $s/3$, $2s/3$, and s relative
395 to the wild-type population, respectively. During one generation, the frequency f_i of individuals
396 with genotype i and selective advantage s_i relative to the wild-type increases by a factor $(1 +$
397 $s_i)$, along with further adjustments to their frequency due to mutations from/to other
398 genotypes. Upon normalisation ($\sum_i f_i = 1$), these frequencies are used for the Dirichlet-
399 multinomial sampling step. After the sampling step, the numbers of cases are converted to
400 frequencies for sampling in the next generation.

401 *Epidemic spread*

402 Due to a high degree of individual-level variation in the transmission of SARS-CoV-2 (i.e.,
403 overdispersion) [25, 51], we use a Dirichlet-multinomial (instead of a multinomial) distribution to
404 assign offspring in generation $t+1$ to parents from generation t . The Dirichlet-multinomial is
405 parametrized by $N(t+1)$ (the number of offspring to draw for the next generation) and $A\vec{f}$, the
406 weights of the different genotypes, where \vec{f} is the normalized vector giving their frequencies in
407 the current generation. The scalar A controls the amount of dispersion, with smaller A
408 corresponding to increased demographic noise. To match it to observations, we note that under
409 the Dirichlet-multinomial model, the number of secondary cases produced by an infection
410 approximately follows a negative binomial distribution with mean $R_t = N(t+1) / N(t)$. In terms of
411 the Dirichlet multinomial parameters, the variance of this negative binomial is

412 $\sigma^2 = \frac{N(t+1)}{N(t)} \left(1 - \frac{1}{N(t)}\right) \frac{N(t+1)+A}{1+A} \approx R_t \frac{N(t+1)+A}{1+A}$. This should match the variance in the number of

413 secondary cases written in terms of the dispersion parameter k , $\sigma^2 = R_t(1 + R_t/k)$. Equating these

414 two expressions gives $A = k N(t) \left(1 - \frac{1}{N(t+1)}\right) - 1 \approx k N(t)$.

415 *Mutation rate*

416 Assuming a constant generation time 5.2 days for all variants of SARS-CoV-2 over time [52], we
417 use the phylogenetically estimated substitution rate per site per year [53] to calculate the
418 mutation rate per site per generation time, μ . We also note that generation time may vary over
419 time depending on the behavioural changes in the population and emergence of variants, that is
420 why we allow for some variation in the mutation rate parameter in our model $(0.87 - 2.0) \times 10^{-5}$
421 based on phylogenetic estimates. We assume each site has two states: wild-type and mutant.
422 Therefore, for a group of K sites, there are 2^K genotypes.

423 *Inferring the selective advantage of VOCs*

424 Finally, the selective advantage s of the VOCs is determined by fitting an exponential function,
425 $f(t)$, of the form, $f(t)=ae^{st}$, to the proportion of Alpha, Beta, and Gamma variants sampled in the
426 country where they were first detected (i.e., UK, South Africa, and Brazil) using the
427 NonlinearModelFit function in Mathematica 11.0 (cite Mathematica). We find that the selective
428 advantage s for Alpha, Beta, and Gamma are 0.37 (95% confidence interval: 0.33 - 0.41), 0.74
429 (95% confidence interval: 0.65 - 0.83), and 0.84 (95% confidence interval: 0.58 - 1.08),
430 respectively (see **Supplementary Figure 7**). The confidence interval is obtained by multiplying the
431 standard error by the value of Student's t for the given confidence level and degrees of
432 freedom. Given the uncertainty in our estimates due to noise in the observations, potential

433 sampling bias, and spatio-temporal heterogeneities, we make the assumption that the value of s
434 is roughly the same for the different VOCs and use the same estimate for all three trajectories
435 (**Table 1**).

436 ***Within-host model of VOC emergence***

437 Each VOC mutation is fixed within the host at rate μ_c such that the fixation time is an
438 exponentially distributed number with mean $1/\mu_c$. Each mutation may then spread to the rest of
439 the population with a probability that is proportional to its fitness as determined by the Dirichlet-
440 multinomial sampling. At any time during the pandemic, a chronic infection can be seeded by
441 other infectious individuals within the population, $N(t)$, with a probability P_f . Therefore, at every
442 generation, the number of chronic infections is given by a binomial distribution with success
443 probability P_f . Once a chronic infection emerges, it remains in the population for the remainder
444 of the simulation period.

445 ***Simulation setup***

446 For both within-host and between-host models of VOC emergence, we run each evolutionary
447 scenario for a given combination of model parameters 1,000 times. We then measure total
448 number of established VOC lineages, M , the time that it takes for the establishment of the first
449 VOC, T_0 , and the time between the establishment of the i^{th} and $(i+1)^{\text{th}}$ VOC, $T_{i:(i+1)}$, for the first 6
450 established VOC lineages in each scenario. An *established* VOC lineage is defined as a lucky
451 lineage with selective advantage s that survives drift upon reaching a size $1/s$. Similarly, the
452 establishment time of a VOC lineage is defined as the time that it takes for that lineage to reach
453 size $1/s$. Each run stops once the frequency of the VOC population reaches 75%.

454 **Code and data availability**

455 All software code and analysis scripts to reproduce the figures are available at
456 https://github.com/weissmanlab/Valley_Crossing.

457 References

- 458 1. L. Z. Zhang, et al., *SARS-CoV-2 spike-protein D614G mutation increases virion spike density and infectivity*.
459 Nature Communications, 2020. **11**: p. 6013.
- 460 2. Hodcroft, E.B., *CoVariants: SARS-CoV-2 Mutations and Variants of Interest*. 2021, URL:
461 <https://covariants.org/>.
- 462 3. T. P. Peacock, et al., *SARS-CoV-2 one year on: evidence for ongoing viral adaptation*. Journal of General
463 Virology, 2021. **102**(4): p. 001584.
- 464 4. A. Rambaut, et al., *Preliminary genomic characterisation of an emergent SARS-CoV-2 lineage in the UK*
465 *defined by a novel set of spike mutations*. Virological, 2020. [https://virological.org/t/preliminary-genomic-](https://virological.org/t/preliminary-genomic-characterisation-of-an-emergent-sars-cov-2-lineage-in-the-uk-defined-by-a-novel-set-of-spike-mutations/563)
466 [characterisation-of-an-emergent-sars-cov-2-lineage-in-the-uk-defined-by-a-novel-set-of-spike-](https://virological.org/t/preliminary-genomic-characterisation-of-an-emergent-sars-cov-2-lineage-in-the-uk-defined-by-a-novel-set-of-spike-mutations/563)
467 [mutations/563](https://virological.org/t/preliminary-genomic-characterisation-of-an-emergent-sars-cov-2-lineage-in-the-uk-defined-by-a-novel-set-of-spike-mutations/563).
- 468 5. H. Tegally, et al., *Detection of a SARS-CoV-2 variant of concern in South Africa*. Nature, 2021. **592**(7854):
469 p. 438-443.
- 470 6. N. R. Faria, et al., *Genomics and epidemiology of the P.1 SARS-CoV-2 lineage in Manaus, Brazil*. Science,
471 2021. **372**(6544): p. 815-821.
- 472 7. K. M. Tao, et al., *The biological and clinical significance of emerging SARS-CoV-2 variants*. Nature Reviews
473 Genetics, 2021. **22**(12): p. 757-773.
- 474 8. D. Planas, et al., *Considerable escape of SARS-CoV-2 Omicron to antibody neutralization*. Nature, 2022.
475 **602**: p. 671-675
- 476 9. T. N. Starr, et al., *Deep Mutational Scanning of SARS-CoV-2 Receptor Binding Domain Reveals Constraints*
477 *on Folding and ACE2 Binding*. Cell, 2020. **182**(5): p. 1295-1310.
- 478 10. N. G. Davies, et al., *Estimated transmissibility and impact of SARS-CoV-2 lineage B.1.1.7 in England*.
479 Science, 2021. **372**: p. 149-154.
- 480 11. S. Cele, et al., *Escape of SARS-CoV-2 501Y.V2 from neutralization by convalescent plasma*. Nature, 2021.
481 **593**: p. 142-146.
- 482 12. D. P. Martin, et al., *The emergence and ongoing convergent evolution of the SARS-CoV-2 N501Y lineages*.
483 Cell, 2021. **184**(20): p. 5189-5200.e7.
- 484 13. J. Zahradnik, et al., *SARS-CoV-2 variant prediction and antiviral drug design are enabled by RBD in vitro*
485 *evolution*. Nature Microbiology, 2021. **6**(9): p. 1188-1198.
- 486 14. G. Nelson, et al., *Molecular dynamic simulation reveals E484K mutation enhances spike RBD-ACE2 affinity*
487 *and the combination of E484K, K417N and N501Y mutations (501Y.V2 variant) induces conformational*
488 *change greater than N501Y mutant alone, potentially resulting in an escape mutant*. bioRxiv, 2021.
489 <https://doi.org/10.1101/2021.01.13.426558>.
- 490 15. B. Choi, M. Cernadas, and J.Z. Li, *Persistence and Evolution of SARS-CoV-2 in an Immunocompromised*
491 *Host*. New England Journal of Medicine, 2020. **383**(23): p. 2291-2293.
- 492 16. S. A. Kemp, et al., *SARS-CoV-2 evolution during treatment of chronic infection*. Nature, 2021. **592**: p. 277-
493 282.
- 494 17. F. Karim, et al., *Persistent SARS-CoV-2 infection and intra-host evolution in association with advanced HIV*
495 *infection*. medRxiv, 2021. <https://doi.org/10.1101/2021.06.03.21258228>.
- 496 18. S. P. Otto, et al., *The origins and potential future of SARS-CoV-2 variants of concern in the evolving COVID-*
497 *19 pandemic*. Current Biology, 2021. **31**(14): p. R918-R929.
- 498 19. T. Karasov, P.W. Messer, and D.A. Petrov, *Evidence that Adaptation in Drosophila Is Not Limited by*
499 *Mutation at Single Sites*. Plos Genetics, 2010. **6**(6): p. e1000924.
- 500 20. M. A. Martin and K. Koelle, *Comment on "Genomic epidemiology of superspreading events in Austria*
501 *reveals mutational dynamics and transmission properties of SARS-CoV-2"*. Science Translational Medicine,
502 2021. **13**(617): p. eabe2555.
- 503 21. K. A. Lythgoe, et al., *SARS-CoV-2 within-host diversity and transmission*. Science, 2021. **372**(6539): p.
504 eabg0821.
- 505 22. K. M. Braun, et al., *Acute SARS-CoV-2 infections harbor limited within-host diversity and transmit via tight*
506 *transmission bottlenecks*. Plos Pathogens, 2021. **17**(8): p. e1009849.
- 507 23. J. Hadfield, et al., *Nextstrain: real-time tracking of pathogen evolution*. Bioinformatics, 2018. **34**(23): p.
508 4121-4123.
- 509 24. D. B. Weissman, et al., *The rate at which asexual populations cross fitness valleys*. Theoretical Population
510 Biology, 2009. **75**(4): p. 286-300.

- 511 25.A. Endo, et al., *Estimating the overdispersion in COVID-19 transmission using outbreak sizes outside China*.
512 Wellcome Open Research, 2020. **5**: p. 67.
- 513 26.R. T. Eguia, et al., *A human coronavirus evolves antigenically to escape antibody immunity*. Plos Pathogens,
514 2021. **17**(4): p. e1009453.
- 515 27.T. de Oliveira, et al., *A novel variant of interest of SARS-CoV-2 with multiple spike mutations detected*
516 *through travel surveillance in Africa*. medRxiv, 2021. <https://doi.org/10.1101/2021.03.30.21254323>.
- 517 28.P. Colson, et al., *Emergence in Southern France of a new SARS-CoV-2 variant of probably Cameroonian*
518 *origin harbouring both substitutions N501Y and E484K in the spike protein*. medRxiv, 2021.
519 <https://doi.org/10.1101/2021.12.24.21268174>.
- 520 29.J. C. Chandler, et al., *SARS-CoV-2 exposure in wild white-tailed deer (*Odocoileus virginianus*)*. Proceedings
521 of the National Academy of Sciences of the United States of America, 2021. **118**(47): p. e2114828118.
- 522 30.B. B. O. Munnink, et al., *Transmission of SARS-CoV-2 on mink farms between humans and mink and back*
523 *to humans*. Science, 2021. **371**(6525): p. 172-177.
- 524 31.H. L. Yen, et al., *Transmission of SARS-CoV-2 delta variant (AY.127) from pet hamsters to humans, leading*
525 *to onward human-to-human transmission: a case study*. Lancet, 2022. **399**(10329): p. 1070-1078.
- 526 32.C. S. Wei, et al., *Evidence for a mouse origin of the SARS-CoV-2 Omicron variant*. Journal of Genetics and
527 Genomics, 2021. **48**(12): p. 1111-1121.
- 528 33.S. V. Kuchipudi, et al., *Multiple spillovers from humans and onward transmission of SARS-CoV-2 in white-*
529 *tailed deer*. Proceedings of the National Academy of Sciences of the United States of America, 2022.
530 **119**(6): p. e2121644119.
- 531 34.B. Pickering, et al., *Highly divergent white-tailed deer SARS-CoV-2 with potential deer-to-human*
532 *transmission*. bioRxiv, 2022. <https://doi.org/10.1101/2022.02.22.481551>.
- 533 35.R. L. Graham and R.S. Baric, *Recombination, Reservoirs, and the Modular Spike: Mechanisms of*
534 *Coronavirus Cross-Species Transmission*. Journal of Virology, 2010. **84**(7): p. 3134-3146.
- 535 36.S. Lytras, et al., *Exploring the Natural Origins of SARS-CoV-2 in the Light of Recombination*. Genome
536 Biology and Evolution, 2022. **14**(2): p. evac018.
- 537 37.G. McVean, P. Awadalla, and P. Fearnhead, *A coalescent-based method for detecting and estimating*
538 *recombination from gene sequences*. Genetics, 2002. **160**(3): p. 1231-1241.
- 539 38.B. Jackson, et al., *Generation and transmission of interlineage recombinants in the SARS-CoV-2 pandemic*.
540 Cell, 2021. **184**(20): p. 5179-5188.
- 541 39.B. Gutierrez, et al., *Emergence and widespread circulation of a recombinant SARS-CoV-2 lineage in North*
542 *America*. medRxiv, 2021. <https://doi.org/10.1101/2021.11.19.21266601>.
- 543 40.R. Viana, et al., *Rapid epidemic expansion of the SARS-CoV-2 Omicron variant in southern Africa*. Nature,
544 2022. **603**(7902): p. 679-682.
- 545 41.R. D. H. Barrett and D. Schluter, *Adaptation from standing genetic variation*. Trends in Ecology &
546 Evolution, 2008. **23**(1): p. 38-44.
- 547 42.D. Planas, et al., *Considerable escape of SARS-CoV-2 Omicron to antibody neutralization*. Nature, 2022.
548 **602**(7898): p. 671-675.
- 549 43.S.A. Wilkinson, et al., *Recurrent SARS-CoV-2 Mutations in Immunodeficient Patients*. medRxiv, 2021.
550 <https://doi.org/10.1101/2022.03.02.22271697>.
- 551 44.A. M. Caliendo, et al., *Better Tests, Better Care: Improved Diagnostics for Infectious Diseases*. Clinical
552 Infectious Diseases, 2013. **57**: p. S139-S170.
- 553 45.M. Ghafari, et al., *Lessons for preparedness and reasons for concern from the early COVID-19 epidemic in*
554 *Iran*. Epidemics, 2021. **36**: p. 100472.
- 555 46. *World Health Organization COVID-19 Dashboard*. . 2021, URL: <https://covid19.who.int>.
- 556 47.A.T. Levin, et al., *Assessing the Burden of COVID-19 in Developing Countries: Systematic Review, Meta-*
557 *Analysis, and Public Policy Implications*. medRxiv, 2021. <https://doi.org/10.1101/2021.09.29.21264325>.
- 558 48.M. O'Driscoll, et al., *Age-specific mortality and immunity patterns of SARS-CoV-2*. Nature, 2020. **590**: p.
559 140-145.
- 560 49.A. Karlinsky and D. Kobak, *Tracking excess mortality across countries during the COVID-19 pandemic with*
561 *the World Mortality Dataset*. Elife, 2021. **10**: p. e69336.
- 562 50.J. T. Wu, et al., *Estimating clinical severity of COVID-19 from the transmission dynamics in Wuhan, China*.
563 Nature Medicine, 2020. **26**(4): p. 506-510
- 564 51.J. O. Lloyd-Smith, et al., *Superspreading and the effect of individual variation on disease emergence*.
565 Nature, 2005. **438**(7066): p. 355-359.

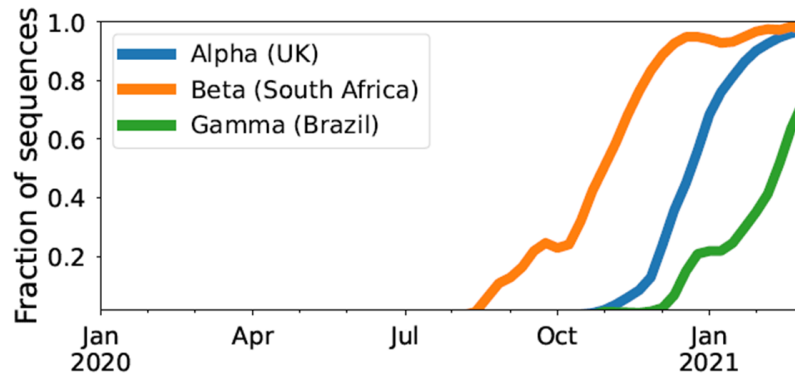
- 566 52.L. Ferretti, et al., *Quantifying SARS-CoV-2 transmission suggests epidemic control with digital contact*
567 *tracing*. *Science*, 2020. **368**(6491): p. 619-625.
- 568 53.M. Ghafari, et al., *Purifying Selection Determines the Short-Term Time Dependency of Evolutionary Rates*
569 *in SARS-CoV-2 and pH1N1 Influenza*. *Molecular Biology and Evolution*, 2022. **39**(2): p. msac009.
- 570 54.K.A. Lythgoe, et al., *Lineage replacement and evolution captured by the United Kingdom Covid Infection*
571 *Survey*. medRxiv, 2022. <https://doi.org/10.1101/2022.01.05.21268323>.

Tables

Table 1: Model parameters.

Symbol	Description	Value (range)	Source
t	Time in units of generations	assuming 5.2 days per generation	[52]
IFR	Global median infection fatality rate of COVID-19	2% - 1.5%)	[47]
N	Number of daily infectious individuals worldwide	daily confirmed deaths / median global IFR	--
μ	Mutation rate per nucleotide per generation	$1.0 (0.87 - 2.0) \times 10^{-5}$	[53]
s	Selective advantage of the VOCs	(0.3 - 1.1)	--
k	Dispersion in distribution of number of secondary infections	0.1 (0.05 - 0.2)	[25]
T_{obs}	Time to the emergence of the first VOC (number of days since 2020-01-03)	(180 - 317) days	[4-6, 54]
ΔT_{obs}	Time between the emergence of the first and second VOC	(0 - 137) days	[4-6, 54]
P_f	Probability of a chronic SARS-CoV-2 infection in an ICI producing a VOC	--	--
μ_c	Within-host fixation rate of VOC mutations per generation	--	--

Figures



Variant\Site	Spike: 18	Spike: 417	Spike: 484	Spike: 501	Spike: 614	NSP-6: 106	N: 203-204
Alpha			E*	Y	G	deletion	K-R
Beta	F	N	K	Y	G	deletion	
Gamma	F	T	K	Y	G	deletion	K-R

*E484K has been reported in some of the Alpha variant genomes sampled within the UK and elsewhere.

Figure 1: **The three initial Variants of Concern arose in quick succession after a long period of limited adaptation.** For each VOC, the curve shows its frequency among the SARS-CoV-2 sequences collected each week from its country of origin. The table shows the amino acid changes across the SARS-CoV-2 genome that are shared between at least two of the three VOCs [7].

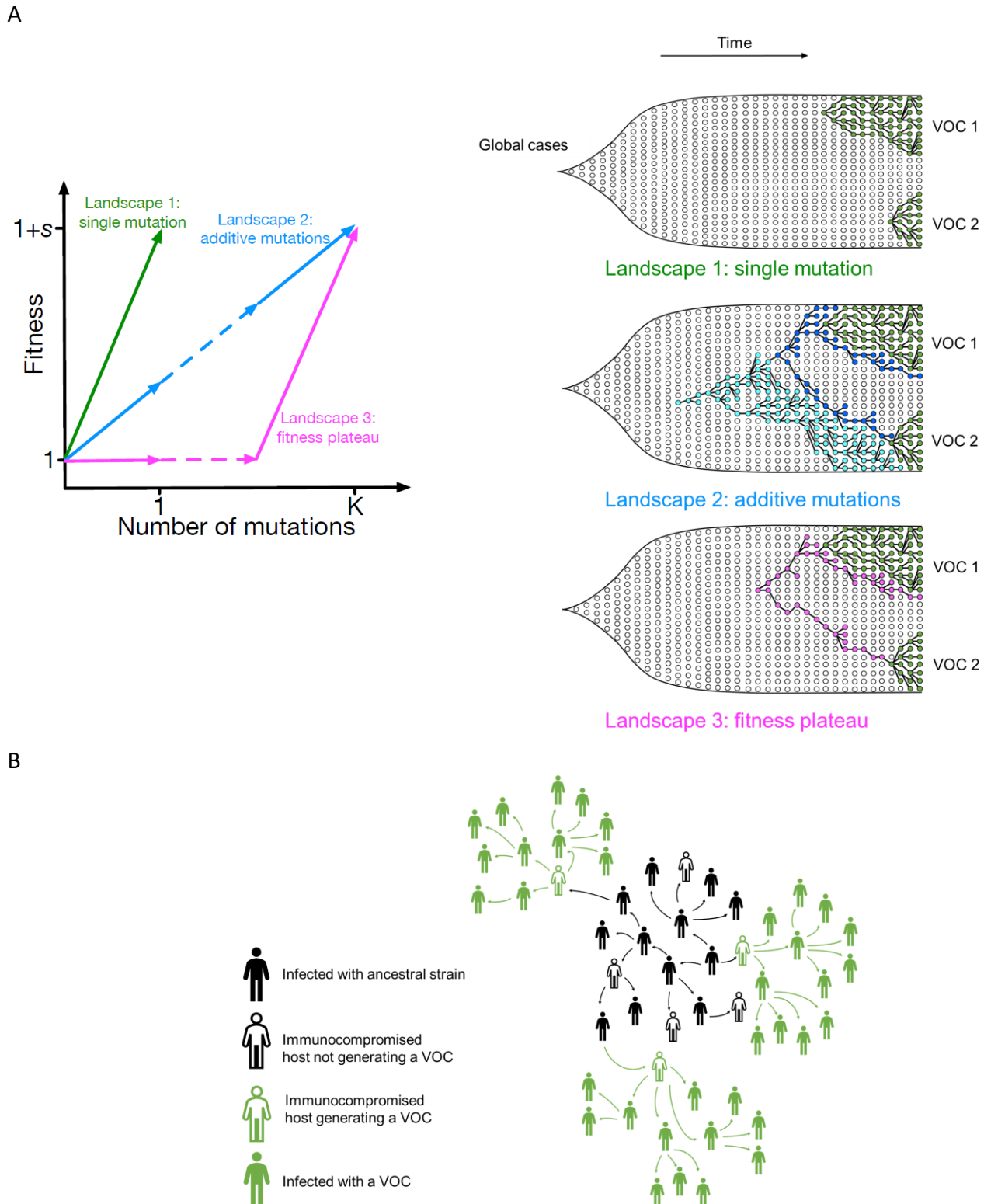


Figure 2: Possible evolutionary pathways to the emergence of SARS-CoV-2 VOCs. (A, left) The three limiting fitness landscapes for the emergence of VOCs as a function of the relevant number of mutations required, K . **(A, right)** VOCs can emerge from either a single advantageous mutation (green) or multiple mutations that each contribute independently to increasing fitness (blue) or only in combination (magenta). **(B)** Emergence of VOCs via the within-host evolutionary path such that an infectious individual passes on a wild-type variant of the virus to an immunocompromised individual where the virus may acquire the relevant mutations during the chronic phase of the infection and later be passed on to the rest of the population.

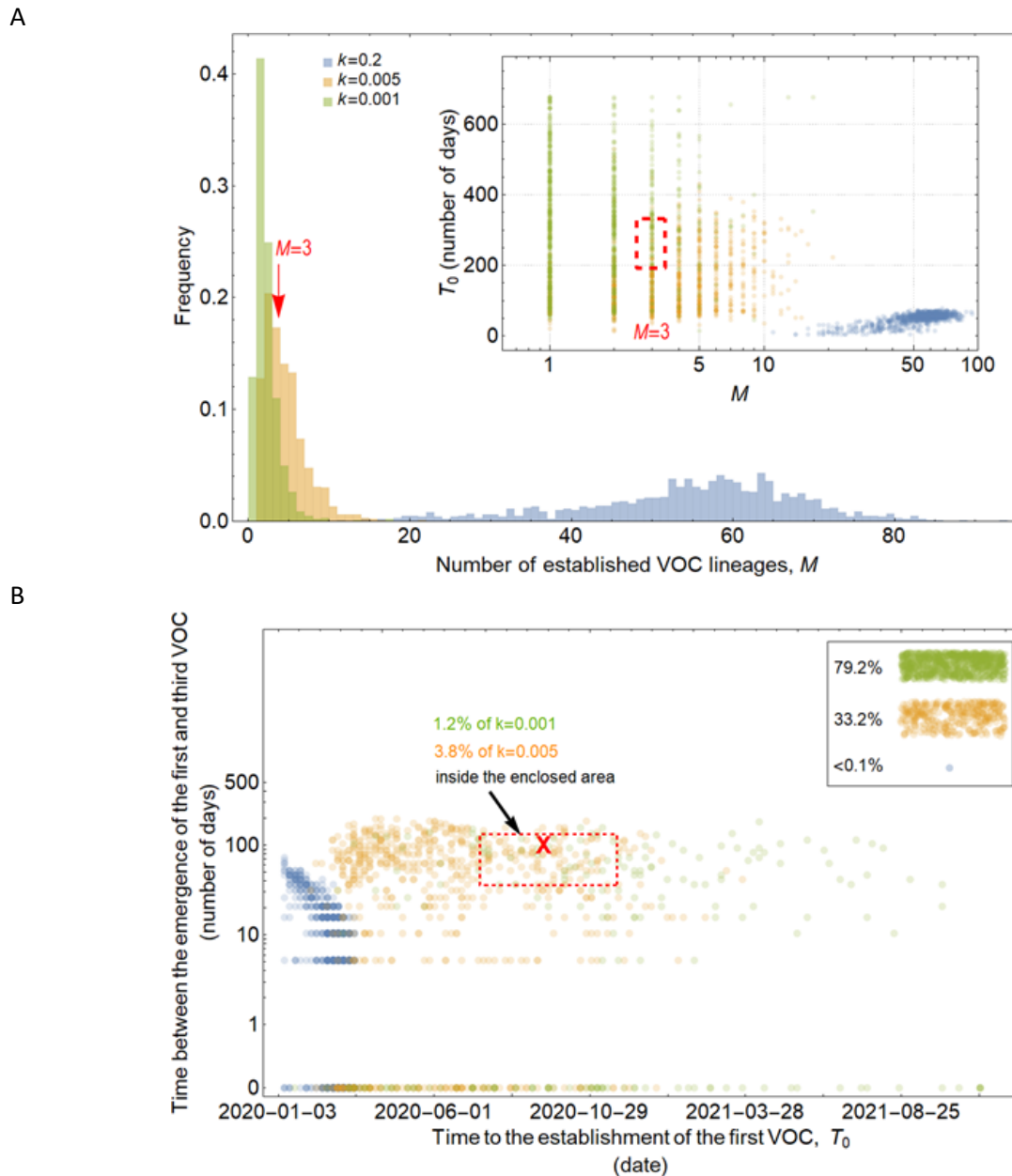
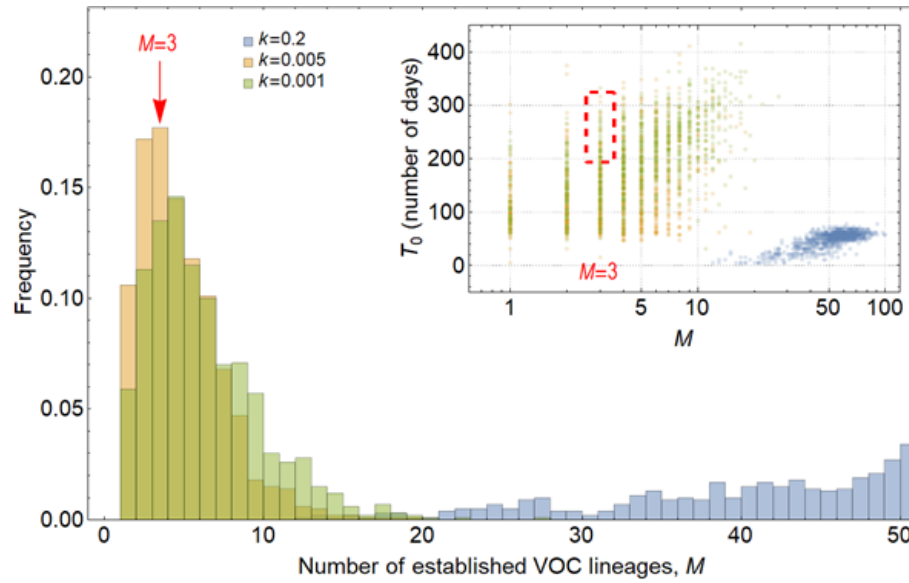
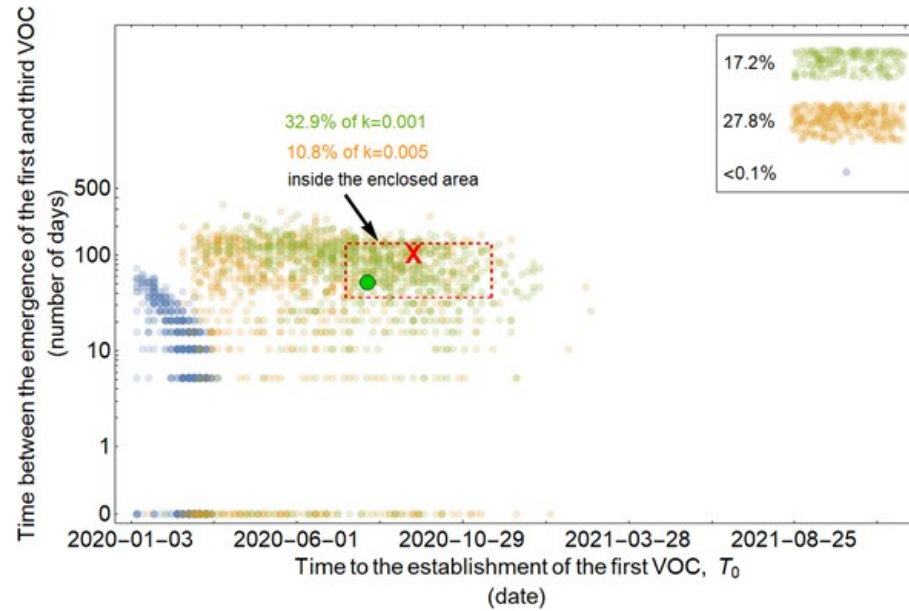


Figure 3: Evolution between hosts on a single-mutation landscape ($K = 1$) rarely reproduces the observed VOC dynamics, even with extreme overdispersion. (A) Total number of established VOC lineages (M) measured under varied levels of overdispersion, k , such that IFR=1.5%, $\mu=0.87 \times 10^{-5}$, $K=1$, and $s=0.4$. The inset shows M with respect to the waiting time for the establishment of the first VOC lineage since the start of the pandemic, T_0 . The region corresponding to the waiting time for the emergence of the first three SARS-CoV-2 VOC is highlighted in red. Under low levels of overdispersion (blue, $k=0.2$), too many VOC lineages are produced very early on in the pandemic. On the other hand, as we increase overdispersion (orange and green), fewer VOCs can establish in the population. It also takes them much longer to establish and reach high frequencies in the population. **(B)** Evaluating the temporal clustering of the first three VOC lineages. For each simulation run, represented by a point on the graph, we measure T_0 and the time difference between the establishment of the first and third successful VOC lineages. The red dashed rectangle shows the region corresponding to the emergence of the first three SARS-CoV-2 VOCs with the cross sign ("X") representing the mean value. We see that as the level of overdispersion increases, the emergence time of VOCs are more scattered and rarely exhibit temporal clustering in late 2020 -- Only 1.2% and 3.8% of the evolutionary dynamics corresponding to overdispersion $k=0.005$ and $k=0.001$ fall inside the enclosed area, respectively. The inset shows that 33.2% and 79.2% of the runs for $k=0.005$ and $k=0.001$ scenarios produce fewer than three successful VOC lineages by the end of the simulation period. Each run stops once the frequency of the VOC population reaches 75%. See also supplementary figure 1.

A



B



C

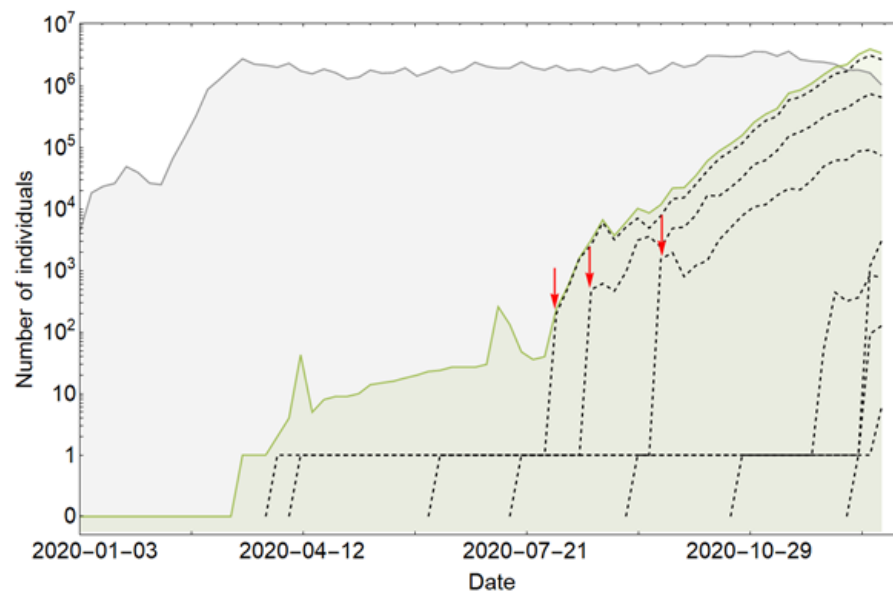
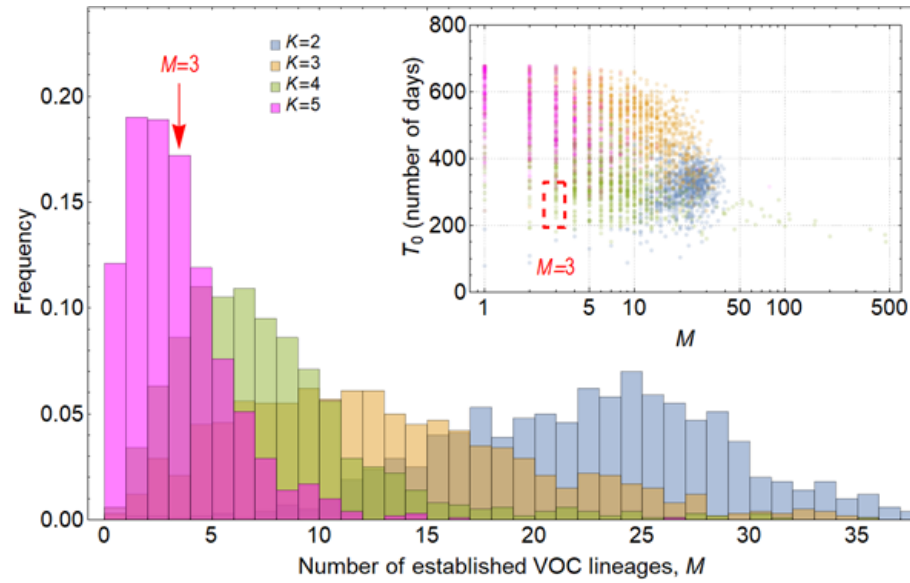


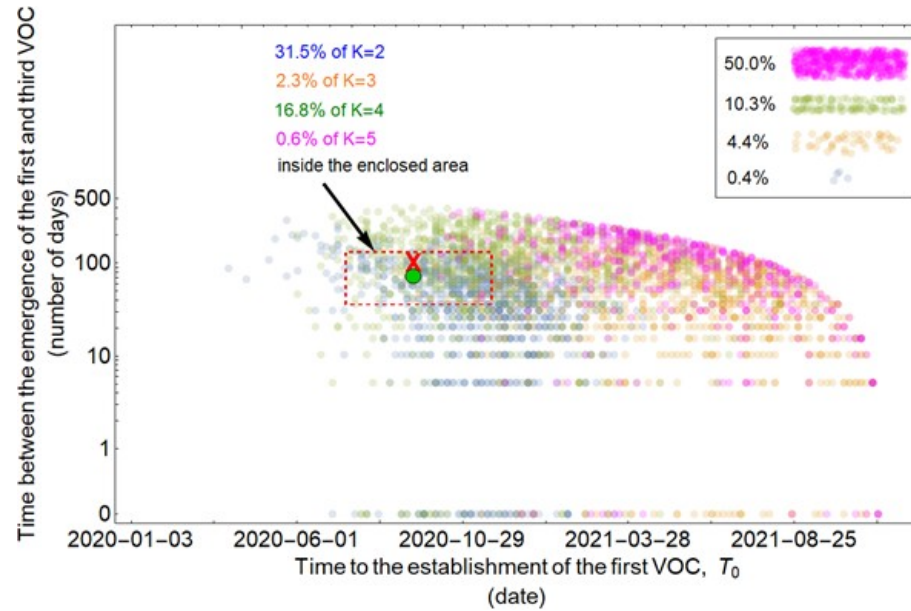
Figure 4: Evolution within hosts on a single-mutation ($K = 1$) landscape can match the observed VOC dynamics, but only with extreme overdispersion to prevent between-host evolution. (A) Total number

of established VOC lineages (M) measured under varied levels of overdispersion, k , where the within-host parameters for the $k=0.2$ scenario (blue) are $P_f=5 \times 10^{-10}$ and $\mu_c=0.001$. For the $k=0.005$ scenario (orange), $P_f=6 \times 10^{-8}$ and $\mu_c=0.1$. Finally, for the $k=0.001$ scenario (green), $P_f=6 \times 10^{-6}$ and $\mu_c=0.01$. For all the three scenarios, the between-host parameters $\mu=0.87 \times 10^{-5}$, IFR=1.5%, $K=1$, and $s=0.4$ are the same. The inset shows M with respect to the waiting time for the establishment of the first VOC lineage since the start of the pandemic, T_0 . The region corresponding to the waiting time for the emergence of the first three SARS-CoV-2 VOC is highlighted in red. We see that under low levels of overdispersion, $k=0.2$ (blue), too many VOC lineages are produced very early on in the pandemic. On the other hand, as we increase overdispersion (orange and green), fewer VOC lineages can establish in the population, and it generally takes longer for them to do so. **(B)** Evaluating the temporal clustering of the first three VOC lineages. For each simulation run, represented by a point on the graph, we measure the time that it takes for a single adaptive mutation to establish in the population and the time difference between the establishment of the first and third successful VOC lineage. The red dashed rectangle shows the region of the parameter space corresponding to the emergence of the first three SARS-CoV-2 VOCs with the cross sign ("X") representing the mean value. The graph shows that by increasing the level of overdispersion and lowering the evolutionary contribution from the between-host pathway, multiple VOCs can emerge in quick succession via the within-host pathway such that a larger fraction of the simulation runs yield the correct timing for the emergence of the first three VOCs in late 2020 (i.e., they fall inside the enclosed area). The inset shows that 27.8% and 17.2% of the runs for $k=0.005$ and $k=0.001$ scenarios produce fewer than three successful VOC lineages by the end of the simulation period. Each run stops once the frequency of the VOC population reaches 75%. **(C)** A typical evolutionary trajectory corresponding to the $k=0.001$ scenario (green) highlighted with a bold green circle in panel (B). The graph shows the VOC population (green) along with the individual VOC lineages (black dashed lines) emerging from the background population (gray). Red vertical arrows show the establishment time of the first three VOCs. We see that the VOC mutation is first produced in a single individual within the population (chronically infected case) for a relatively long time before successfully spreading to the rest of the population. See also supplementary figure 2.

A



B



C

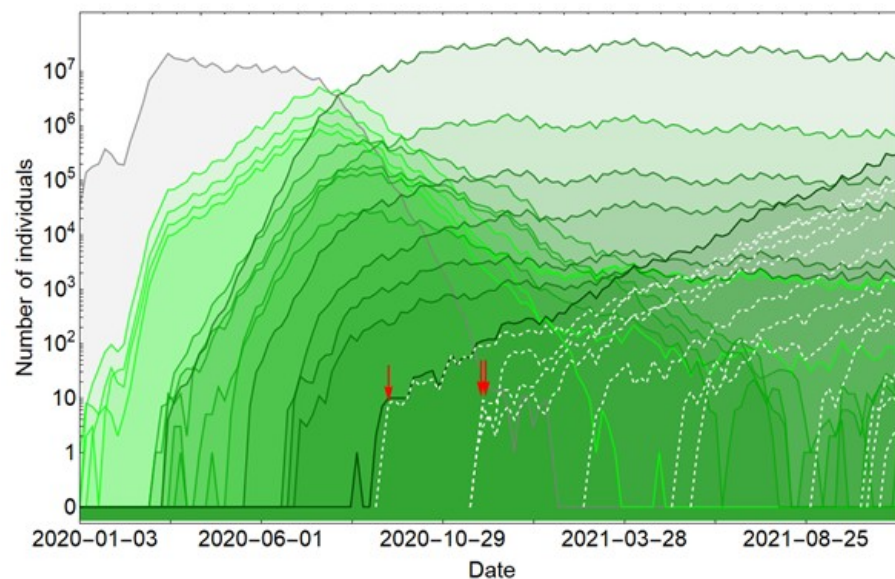


Figure 5: **Between-host evolution on an additive fitness landscape can match the observed VOC dynamics, but only by having intermediate mutants reach unrealistically high frequencies. (A)** Total

number of established VOC lineages (M) for different number of mutations, K , involved in the production of a VOC. For the $K=2$ scenario (blue), IFR=1.5%, $\mu=0.87 \times 10^{-5}$, $k=0.05$, and $s=0.3$. For the $K=3$ scenario (orange), IFR=0.5%, $\mu=0.87 \times 10^{-5}$, $k=0.05$, and $s=0.5$. For both the $K=4$ (green) and $K=4$ (magenta) scenarios, IFR=0.2%, $\mu=2 \times 10^{-5}$, $k=0.2$, and $s=1.0$. The inset shows M with respect to the waiting time for the establishment of the first VOC lineage since the start of the pandemic, T_0 . Under $K=2$ and 3, a very large number of successful VOC lineages are produced by late 2020, with the $K=2$ scenario producing, on average, more than 20 VOC lineages that establish in the population. On the other hand, for the $K=5$ scenario, on average, fewer than 3 lineages are produced. It also takes much longer for them to establish in the population. **(B)** Evaluating the temporal clustering of the first three VOC lineages. For each simulation run, represented by a point on the graph, we measure the time that it takes for a single adaptive mutation to establish in the population and the time difference between the establishment of the first and third successful VOC lineage. The red dashed rectangle shows the region of the parameter space corresponding to the emergence of the first three SARS-CoV-2 VOCs with the cross sign ("X") representing the mean value. We see a noticeable overlap between the $K=2$ and 4 scenarios and the red rectangle suggesting that a larger fraction of the simulation runs exhibit temporal clustering dynamics for VOC emergence. The inset shows that 10.3% of the runs for the $K=4$ scenario produce fewer than three successful VOC lineages by the end of the simulation period. Each run stops once the frequency of the VOC population reaches 75%. **(C)** A typical evolutionary trajectory corresponding to the $K=4$ scenario highlighted with a bold green circle in panel (B). The graph shows the background population in gray and the i -mutant populations ($1 < i \leq K$) in different shades of green from light (fewer mutations) to dark (more mutations). Note that for the $K=4$ scenario, there are four single-mutant, six double-mutant, four triple-mutant, and one quadruple-mutant genotypes. The dashed lines show the dynamics of all the established VOC lineages over time. Red vertical arrows show the establishment time of the first three VOCs. We can see that some of the intermediate mutant genotypes reach close to fixation before giving rise to the VOC population. See also supplementary figure 3.

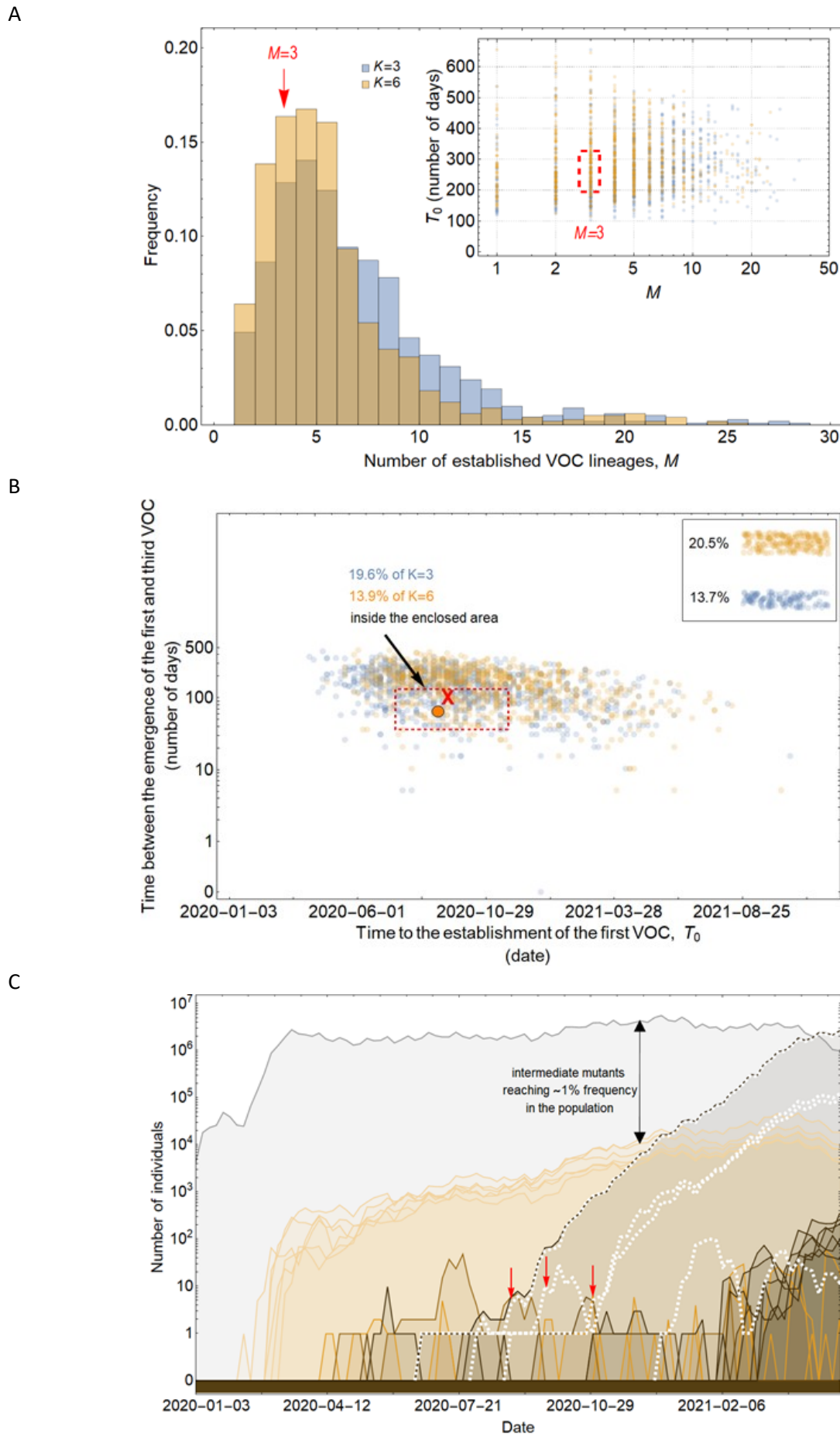
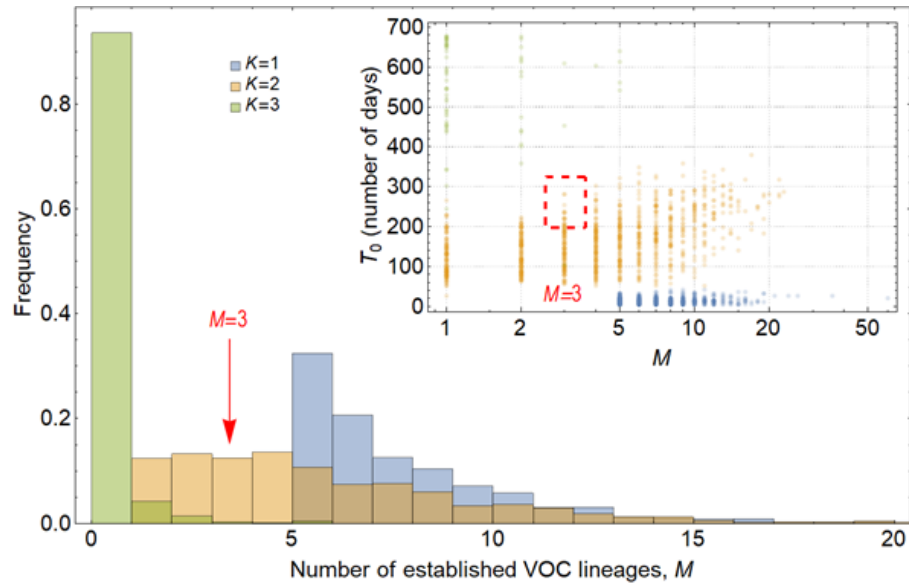


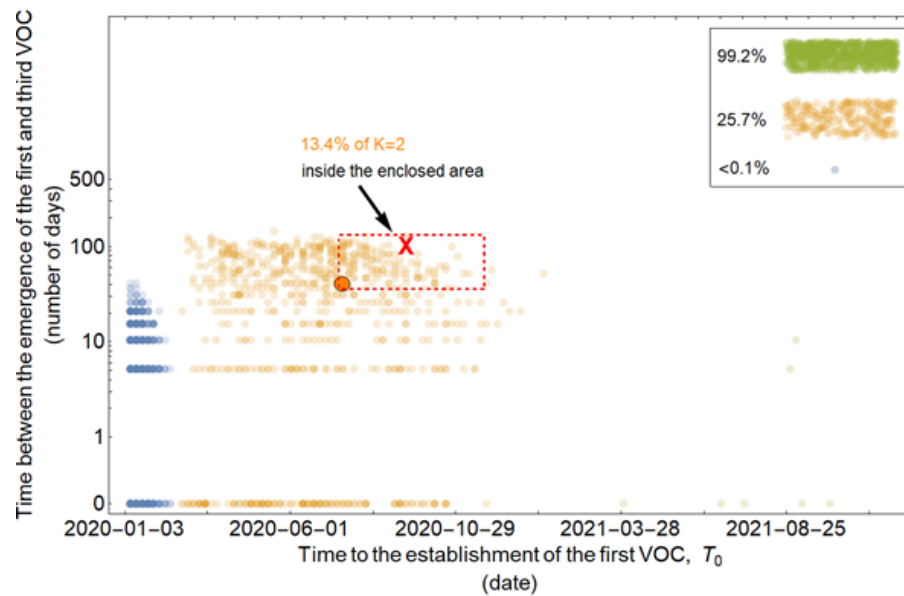
Figure 6: Evolution within hosts on an additive fitness landscape can match the observed VOC dynamics as long as K is large enough that between-host evolution is ineffective. (A) Total number of established

VOC lineages (M for different number of mutations, K , involved in the production of a VOC. For $K=3$ (blue), the within-host parameters are $P_i=3.5 \times 10^{-8}$, and $\mu_c=0.15$. For $K=6$ (orange), $P_i=3 \times 10^{-8}$, and $\mu_c=0.3$. In both scenarios, the between-host parameters $\mu=0.87 \times 10^{-5}$, IFR=1.5%, $k=0.05$, and $s=0.3$ are the same. The inset shows M with respect to the waiting time for the establishment of the first VOC lineage since the start of the pandemic, T_0 . The region corresponding to the waiting time for the emergence of the first three SARS-CoV-2 VOC is highlighted in red. Both scenarios produce roughly the same number of VOC lineages. However, on average, T_0 is slightly longer for the $K=6$ scenario. **(B)** Evaluating the temporal clustering of the first three VOC lineages. For each simulation run, represented by a point on the graph, we measure the time that it takes for a single adaptive mutation to establish in the population and the time difference between the establishment of the first and third successful VOC lineage. The red dashed rectangle shows the region of the parameter space corresponding to the emergence of the first three SARS-CoV-2 VOCs with the cross sign ("X") representing the mean value. We can see that by having a combination of relatively high level of overdispersion, high IFR, and low between-host mutation rate, there is a lower chance of intermediate mutations reaching fixation via the between-host path. Instead, multiple VOCs can emerge in quick succession during chronic infections such that a relatively large fraction of the simulation runs yield a temporal clustering that matches the emergence of the first three VOCs in late 2020 (i.e., they fall inside the enclosed area). The inset shows that 20.5% and 13.7% of the runs for $K=3$ and 6 scenarios produce fewer than three successful VOC lineages by the end of the simulation period, respectively. Each run stops once the frequency of the VOC population reaches 75%. **(C)** A typical evolutionary trajectory corresponding to the $K=6$ scenario highlighted with a bold orange circle in panel (B). The graph shows the background population in gray and the i -mutant populations ($1 < i \leq K$) in different shades of green from light (fewer mutations) to dark (more mutations). The dashed lines show the dynamics of all the established VOC lineages over time. Red vertical arrows show the establishment time of the first three VOCs. We can see that the single-mutant genotypes (lines in light orange) are produced via the between-host pathway but never reach above 1% prevalence before the emergence of the VOCs (white dashed lines). See also supplementary figure 4.

A



B



C

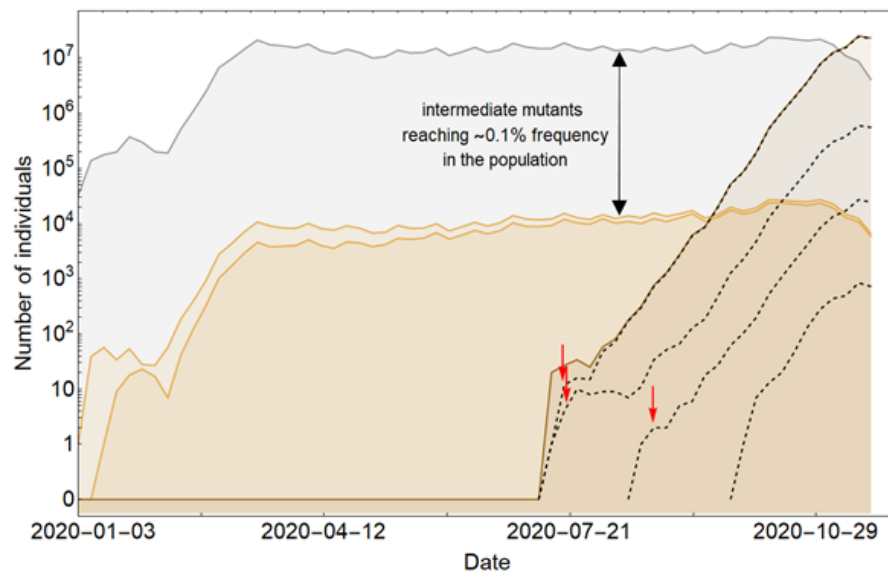
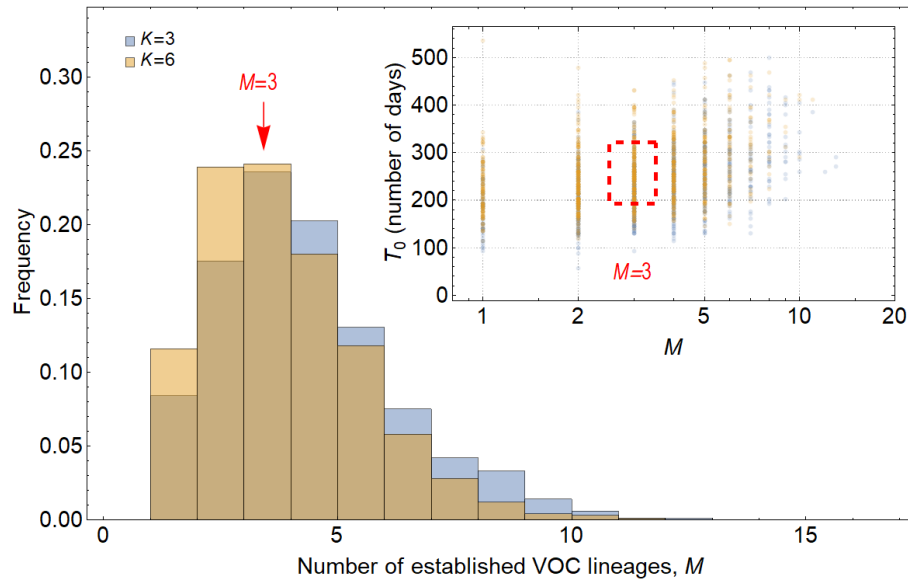


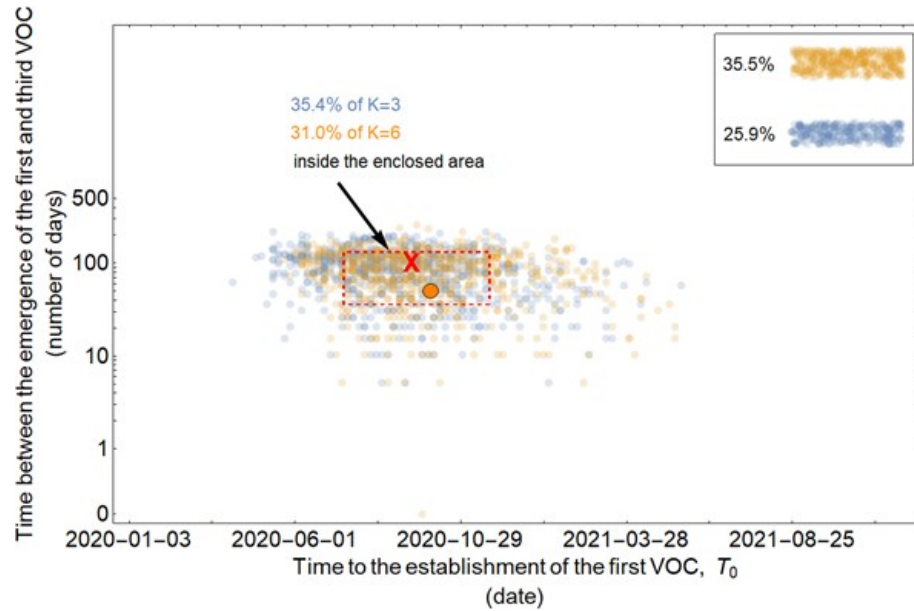
Figure 7: **Between-host evolution on a fitness plateau can match the observed VOC dynamics, but only for $K = 2$.** (A) Total number of established VOC lineages (M) for different number of mutations, K , involved

in the production of a VOC, such that $IFR=0.2\%$, $\mu=2 \times 10^{-5}$, $k=0.2$, and $s=1.0$. The inset shows M with respect to the waiting time for the establishment of the first VOC lineage since the start of the pandemic, T_0 . For $K=1$ and 3 scenarios, there are too many and too few VOC lineages are produced by late 2020. Only for the $K=2$ scenario we can see an intermediate number of VOC lineages being produced in the right time span. **(B)** Evaluating the temporal clustering of the first three VOC lineages. For each simulation run, represented by a point on the graph, we measure the time that it takes for a single adaptive mutation to establish in the population and the time difference between the establishment of the first and third successful VOC lineage. The red dashed rectangle shows the region of the parameter space corresponding to the emergence of the first three SARS-CoV-2 VOCs with the cross sign ("X") representing the mean value. We see a noticeable overlap between the $K=2$ scenario and the red rectangle suggesting that a fraction of the simulation runs exhibit temporal clustering dynamics for VOC emergence. The inset shows that 99.2% and 25.7% of the runs for the $K=3$ and 2 scenarios produce fewer than three successful VOC lineages by the end of the simulation period. Each run stops once the frequency of the VOC population reaches 75%. **(C)** A typical evolutionary trajectory corresponding to the $K=6$ scenario highlighted with a bold orange circle in panel (B). The graph shows the background population in gray, single-mutants in light orange, and double-mutants in dark orange. Note that for the $K=2$ scenario, there are two single-mutant and one double-mutant genotypes. The dashed lines show the dynamics of all the established VOC lineages over time. Red vertical arrows show the establishment time of the first three VOCs. We can see that the single-mutant genotypes reach close to 0.1% before giving rise to the VOC population. See also supplementary figure 5.

A



B



C

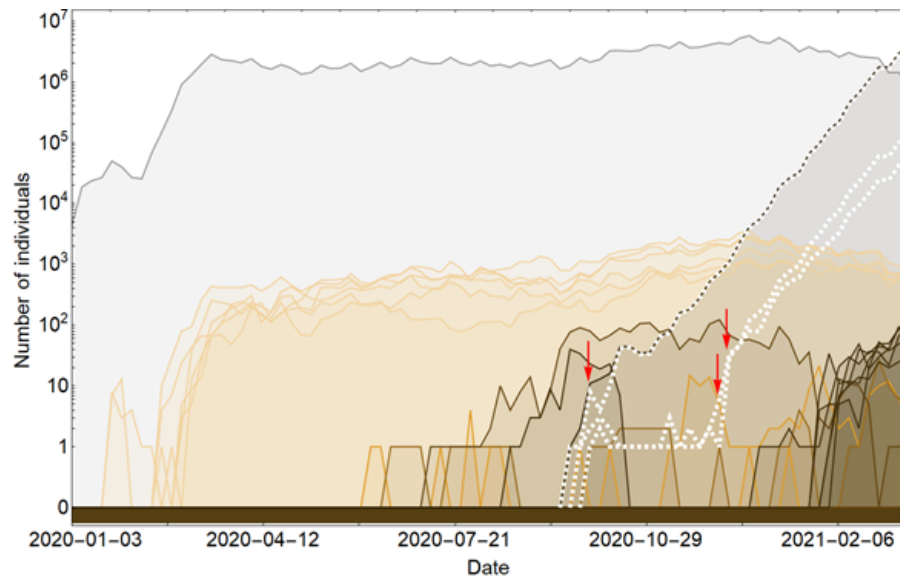


Figure 8: Within-host evolution on a fitness plateau can match the observed VOC dynamics for a large range of plateau widths. (A) Total number of established VOC lineages (M for different number of

mutations, K , involved in the production of a VOC. For $K=3$ (blue), the within-host parameters are $P_f=2 \times 10^{-8}$, and $\mu_c=0.1$. For $K=6$ (orange), $P_f=4.5 \times 10^{-8}$, and $\mu_c=0.25$. In both scenarios, the between-host parameters $\mu=1 \times 10^{-5}$, IFR=0.5%, $k=0.1$, and $s=0.7$ are the same. The inset shows M with respect to the waiting time for the establishment of the first VOC lineage since the start of the pandemic, T_0 . The region corresponding to the waiting time for the emergence of the first three SARS-CoV-2 VOC is highlighted in red. Both scenarios produce roughly the same number of VOC lineages. However, on average, T_0 is slightly longer for the $K=6$ scenario. **(B)** Evaluating the temporal clustering of the first three VOC lineages. For each simulation run, represented by a point on the graph, we measure the time that it takes for a single adaptive mutation to establish in the population and the time difference between the establishment of the first and third successful VOC lineage. The red dashed rectangle shows the region of the parameter space corresponding to the emergence of the first three SARS-CoV-2 VOCs with the cross sign ("X") representing the mean value. We can see that a noticeable fraction of simulation runs for both scenarios yield a temporal clustering that matches the emergence of the first three VOCs in late 2020 (i.e., they fall inside the enclosed area). The inset shows that 35.5% and 25.9% of the runs for $K=3$ and 6 scenarios produce fewer than three successful VOC lineages by the end of the simulation period, respectively. Each run stops once the frequency of the VOC population reaches 75%. **(C)** A typical evolutionary trajectory corresponding to the $K=6$ scenario highlighted with a bold orange circle in panel (B). The graph shows the background population in gray and the i -mutant populations ($1 < i \leq K$) in different shades of orange from light (fewer mutations) to dark (more mutations). The dashed lines show the dynamics of all the established VOC lineages over time. Red vertical arrows show the establishment time of the first three VOCs. We can see that the single-mutant genotypes (lines in light orange) are produced via the between-host pathway from very early on in the pandemic but are at very low prevalence before the emergence of the VOCs (white dashed lines). See also supplementary figure 6.

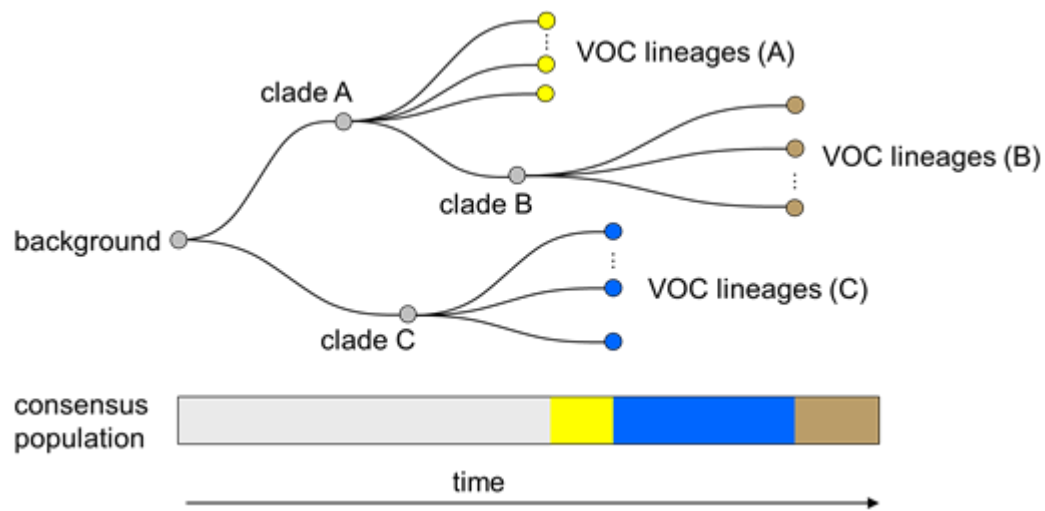
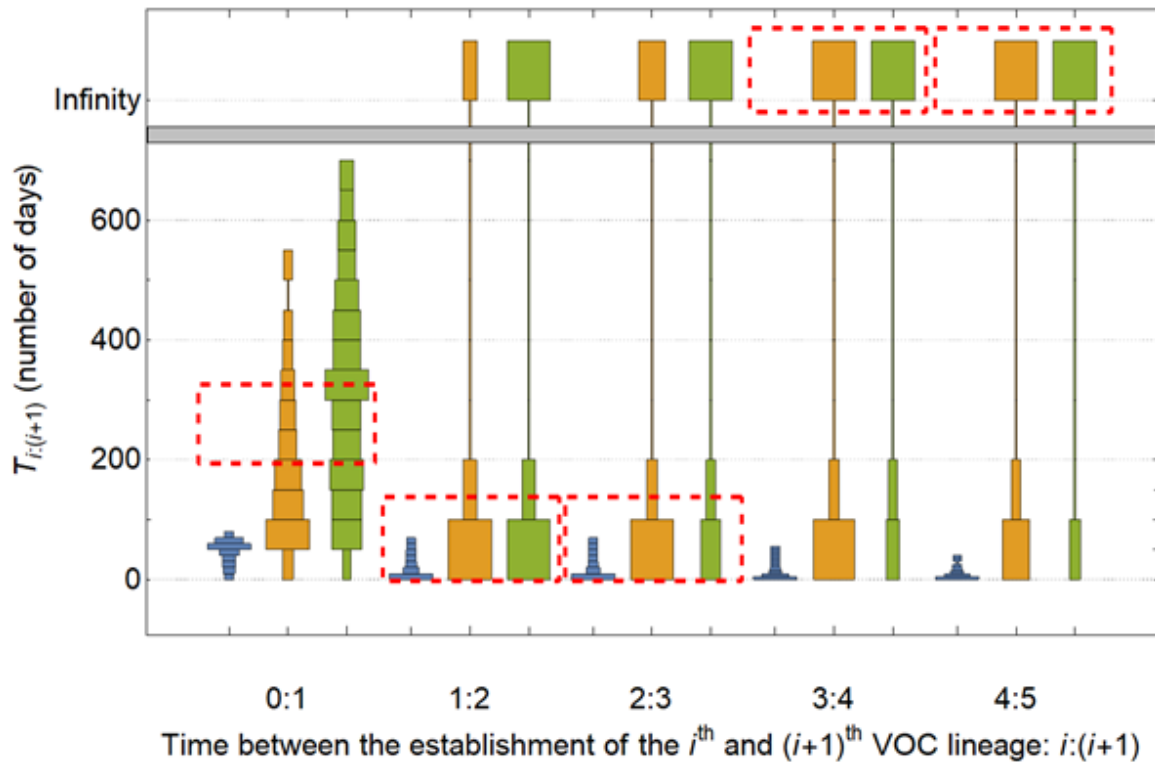
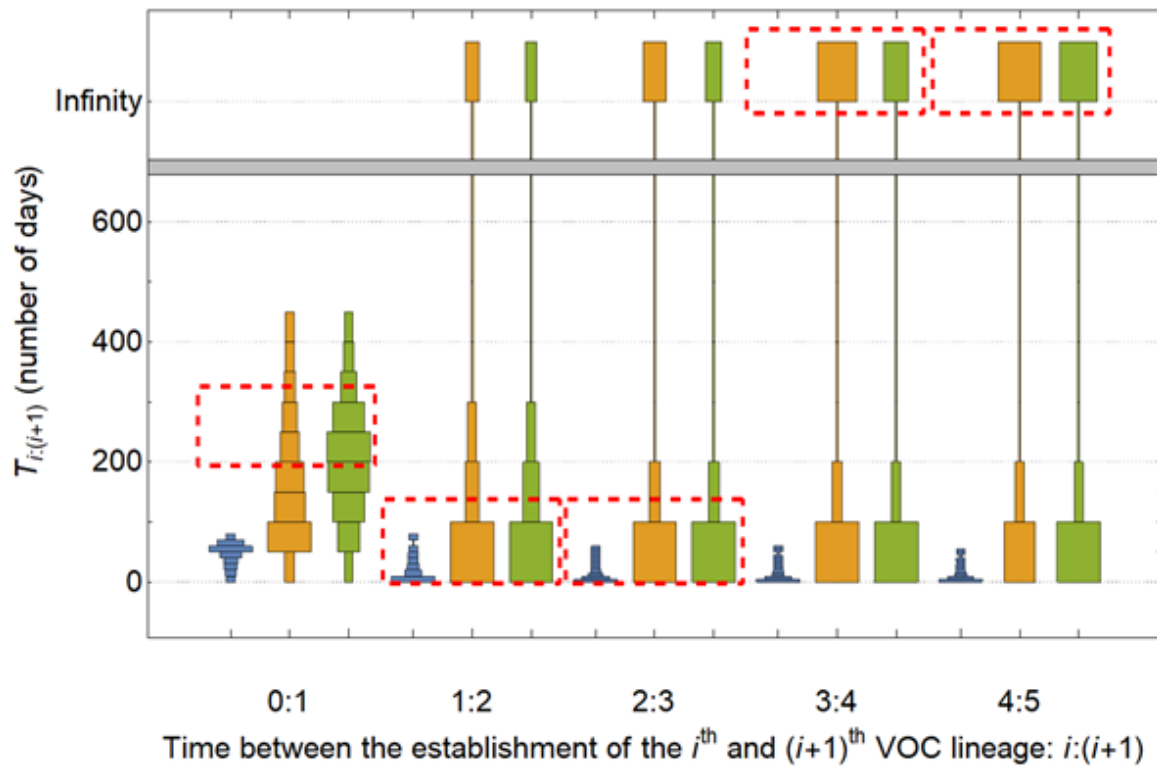


Figure 9: **Within-host evolution reproduces the star-like genealogy of the VOCs.** For a wide range of parameter combinations in Landscape 1, 2, and 3, we showed that the within-host evolutionary dynamics can become virtually uncoupled from the the evolution at the between-host level enabling each VOC lineage to arise independently on the background of a different clade which leads to a star-like tree topology.

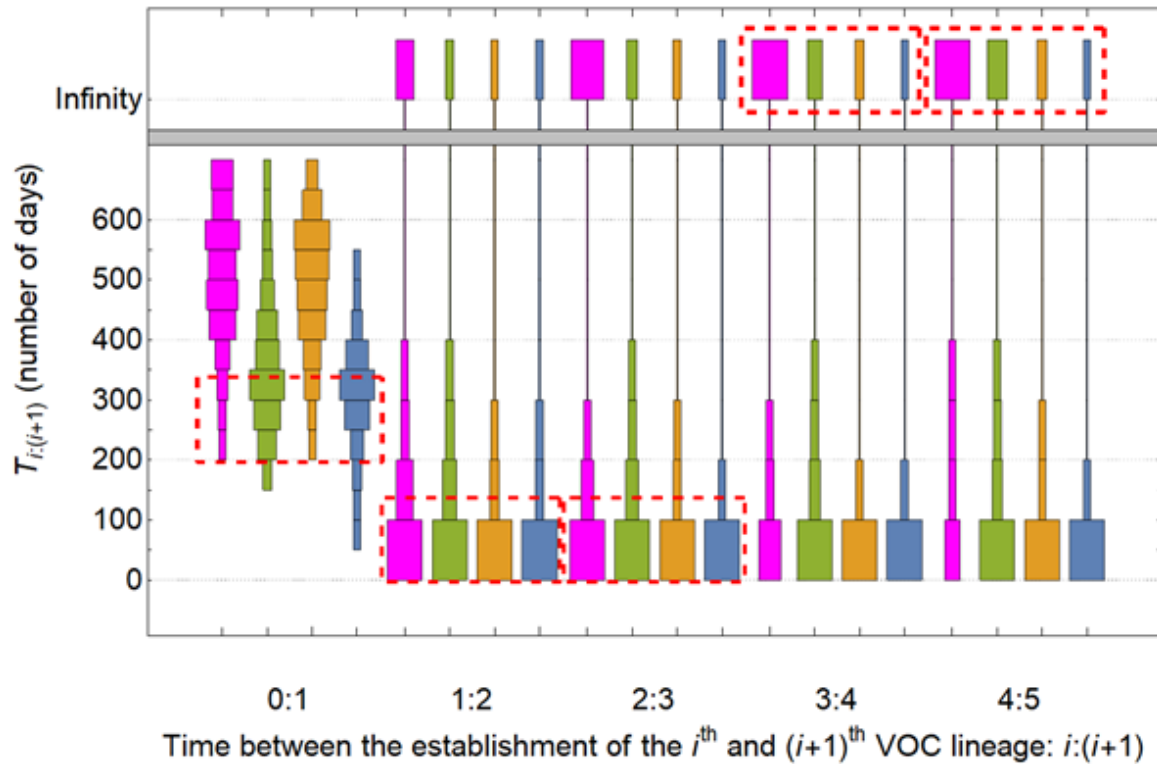
Supplementary figures



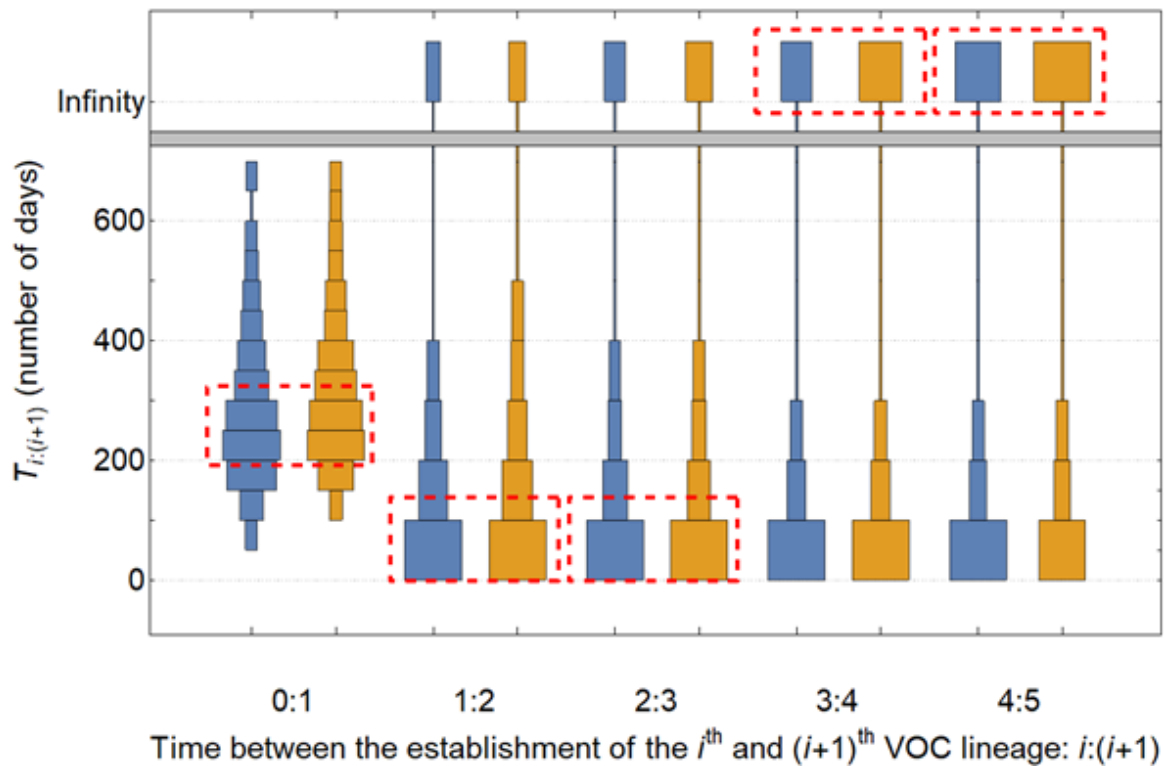
Supplementary figure 1: **Distribution of waiting times for the establishment of consecutive pairs of VOC lineages via the between-host pathway assuming a fitness landscape with a single adaptive mutation.** The distribution of times that it takes between the production of the i^{th} and $(i+1)^{\text{th}}$ lineage, $T_{i:(i+1)}$, for the first 5 established VOC lineages described in Figure 3. $T_{0:1}$ is the waiting time for the production for the establishment of the first VOC lineage (equivalent to T_0).



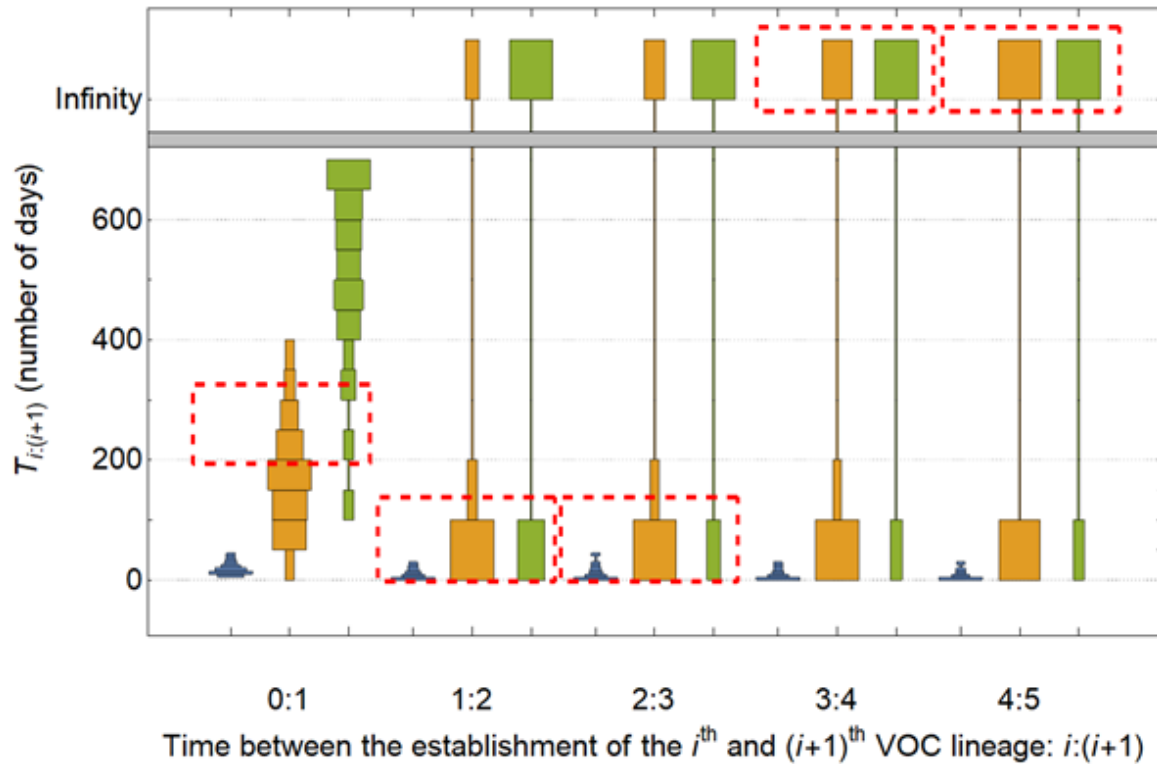
Supplementary figure 2: **Distribution of waiting times for the establishment of consecutive pairs of VOC lineages via the within-host pathway assuming a fitness landscape with a single adaptive mutation.** The distribution of times that it takes between the production of the i^{th} and $(i+1)^{\text{th}}$ lineage, $T_{i:(i+1)}$, for the first 5 established VOC lineages described in Figure 4. $T_{0:1}$ is the waiting time for the production for the establishment of the first VOC lineage (equivalent to T_0).



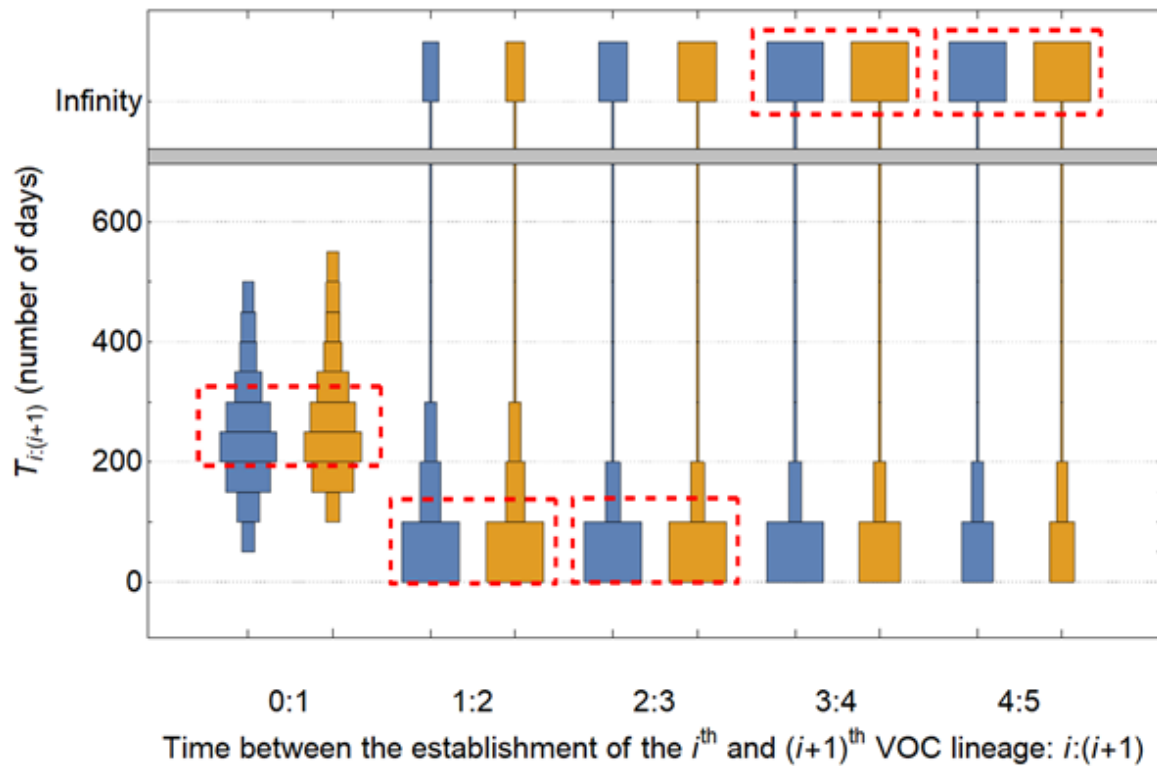
Supplementary figure 3: **Distribution of waiting times for the establishment of consecutive pairs of VOC lineages via the between-host pathway assuming an additive fitness landscape.** The distribution of times that it takes between the production of the i^{th} and $(i+1)^{\text{th}}$ lineage, $T_{i:(i+1)}$, for the first 5 established VOC lineages described in Figure 5. $T_{0:1}$ is the waiting time for the production for the establishment of the first VOC lineage (equivalent to T_0).



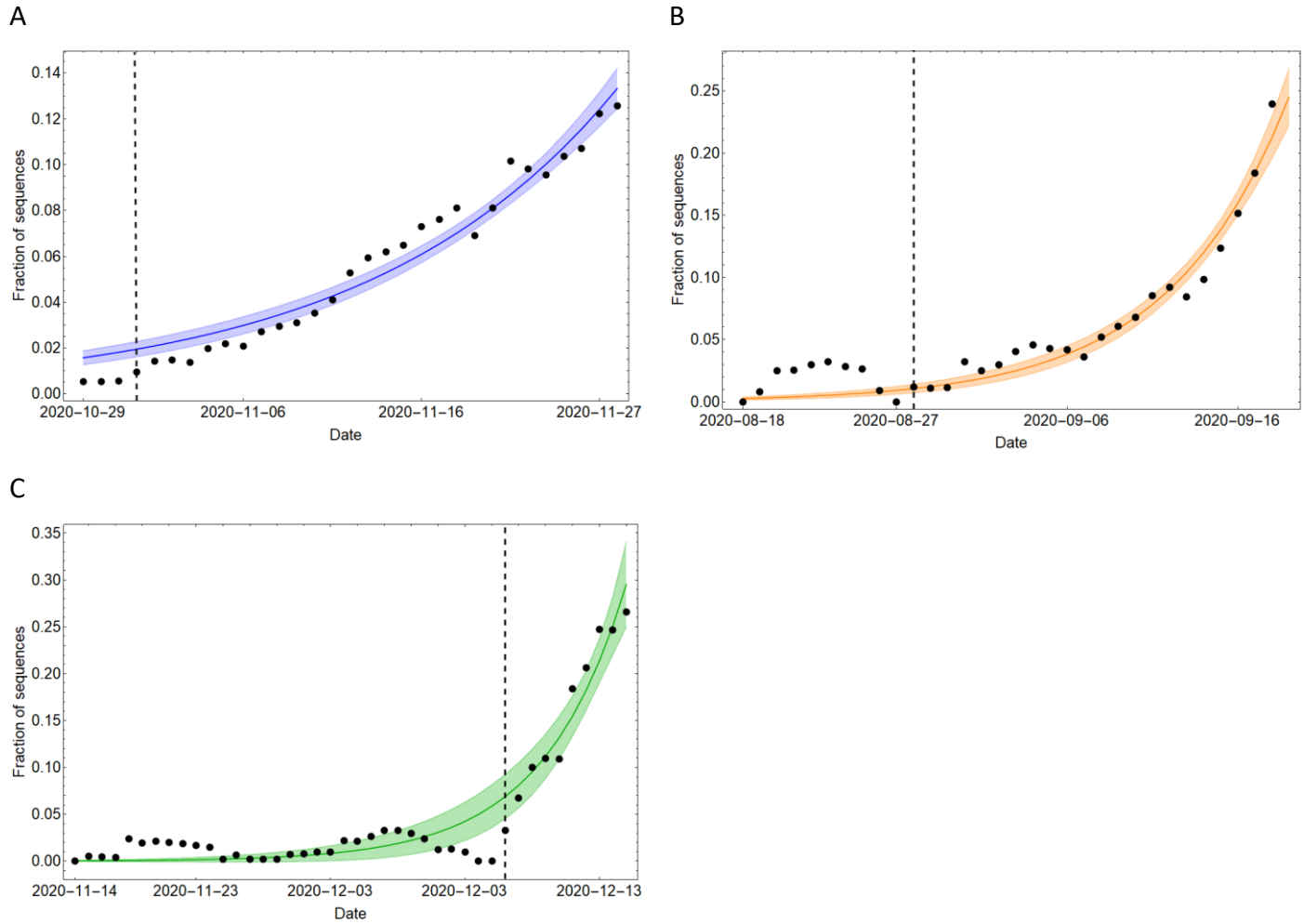
Supplementary figure 4: **Distribution of waiting times for the establishment of consecutive pairs of VOC lineages via the within-host pathway assuming an additive fitness landscape.** The distribution of times that it takes between the production of the i^{th} and $(i+1)^{\text{th}}$ lineage, $T_{i:(i+1)}$, for the first 5 established VOC lineages described in Figure 6. $T_{0:1}$ is the waiting time for the production for the establishment of the first VOC lineage (equivalent to T_0).



Supplementary figure 5: **Distribution of waiting times for the establishment of consecutive pairs of VOC lineages via the between-host pathway assuming a fitness plateau landscape.** The distribution of times that it takes between the production of the i^{th} and $(i+1)^{\text{th}}$ lineage, $T_{i:(i+1)}$, for the first 5 established VOC lineages described in Figure 7. $T_{0:1}$ is the waiting time for the production for the establishment of the first VOC lineage (equivalent to T_0).



Supplementary figure 6: **Distribution of waiting times for the establishment of consecutive pairs of VOC lineages via the within-host pathway assuming a fitness plateau landscape.** The distribution of times that it takes between the production of the i^{th} and $(i+1)^{\text{th}}$ lineage, $T_{i:(i+1)}$, for the first 5 established VOC lineages described in Figure 8. $T_{0:1}$ is the waiting time for the production for the establishment of the first VOC lineage (equivalent to T_0).



Supplementary figure 7: Exponential model fits to the frequency of individual SARS-CoV-2 VOC sequences sampled in its country of origin. (A-C) Fitting an exponential function of the form, $f(t)=ae^{bt}$, to the frequency of Alpha, Beta, and Gamma sequences sampled in the UK, South Africa, and Brazil, respectively. Vertical dashed line shows the starting timepoint used for the fitting. The shaded area shows the mean prediction bands.