

1 Immune Heterogeneity and Epistasis 2 Explain Punctuated Evolution of 3 SARS-CoV-2

4 Bjarke Frost Nielsen^{1‡,3*}, Yimei Li^{2†}, Kim Sneppen^{3§}, Lone Simonsen^{1‡}, Cécile
5 Viboud^{4¶}, Simon A. Levin^{2†}, Bryan T. Grenfell^{2†*}

***For correspondence:**

bjarkefrost@ruc.dk (BFN);
grenfell@princeton.edu (BTG)

Present address: [†]Dept of Ecology
and Evolutionary Biology, 106A
Guyot Hall, Princeton University,
Princeton, NJ 08544-2016, United
States of America.; [‡]Department of
Science and Environment, Roskilde
University, Universitetsvej 1, 4000
Roskilde, Denmark; [§]Niels Bohr
Institute, University of
Copenhagen, Blegdamsvej 17,
2100 Copenhagen Ø, Denmark;
[¶]Division of International
Epidemiology and Population
Studies, Fogarty International
Center, National Institutes of
Health, Bethesda, MD 20892-2220;
United States of America.

6 ¹PandemiX Center, Roskilde University; ²Department of Ecology & Evolutionary Biology,
7 Princeton University; ³Niels Bohr Institute, University of Copenhagen; ⁴Fogarty
International Center, National Institutes of Health

Abstract Identifying drivers of viral diversity is key to understanding the evolutionary as well as
epidemiological dynamics of the COVID-19 pandemic. Using rich viral genomic data sets, we show
that periods of steadily rising diversity have been punctuated by sudden, enormous increases
followed by similarly abrupt collapses of diversity. We introduce a mechanistic model of
saltational evolution with epistasis and demonstrate that these features parsimoniously account
for the observed temporal dynamics of inter-genomic diversity. Our results provide support for
recent proposals that saltational evolution may be a signature feature of SARS-CoV-2, allowing
the pathogen to more readily evolve highly transmissible variants. These findings lend theoretical
support to a heightened awareness of biological contexts where increased diversification may
occur. They also underline the power of pathogen genomics and other surveillance streams in
clarifying the phylodynamics of emerging and endemic infections. In public health terms, our
results further underline the importance of equitable distribution of up-to-date vaccines.

23 Introduction

24 During the coronavirus disease 2019 (COVID-19) pandemic, the responsible pathogen, severe acute
25 respiratory syndrome coronavirus 2 (SARS-CoV-2), has continuously evolved. However, evolution
26 has by no means happened at an even pace, but rather through a pattern of steady diversification
27 punctuated by sudden large jumps involving dozens of point mutations. Indeed, it has been sug-
28 gested that SARS-CoV-2 exhibits saltational evolution, a process where evolution proceeds by large
29 multimutational jumps, rather than gradually (*Corey et al., 2021*).

30
31 A simple way to quantify the genomic diversity existing at a given time is through the pairwise
32 Hamming distance. Given two genomes, the pairwise Hamming distance simply measures how
33 many nucleotides the two sequences disagree on. This rather crude measure turns out to reveal
34 surprisingly robust patterns of viral diversification.

35 Due to the large amount of full genome sequencing performed on SARS-CoV-2 specimens during
36 the COVID-19 pandemic, Hamming distances can be computed not just at the level of summary
37 statistics, but as temporally varying distributions (Figure 1; Video 1), revealing a pattern of slowly
38 increasing diversity punctuated by abrupt increases and subsequent collapses in diversity.

39 That much can be gleaned from considering the time-development of the mean (or median) Ham-
40 ming distance. However, the dynamics of the often multimodal distribution is not captured by the

41 mean Hamming distance, even if temporally resolved, and much less by the usual static treatment.
42 The full, time-dependent Hamming distribution possesses further structure that reveals that suc-
43 cessive variants are well-separated in sequence space, suggesting that one did not arise from the
44 other by a string of single-point mutations accruing in successive hosts. Rather, an evolutionary
45 jump – a saltation – seems to have taken place at each major transition (see Video 1).

46 **Dynamical explanations**

47 There are several plausible mechanisms that may contribute to saltational evolution in SARS-CoV-2,
48 including increased build-up of mutations in immunocompromised individuals infected with SARS-
49 CoV-2 (*Corey et al., 2021; Nussenblatt et al., 2022; Avanzato et al., 2020; Choi et al., 2020; Truong*
50 *et al., 2021; Kemp et al., 2021; Harari et al., 2022; Kumata and Sasaki, 2022*) and evolution in animal
51 reservoirs followed by animal-to-human transmission (*Kupferschmidt, 2021; Larsen et al., 2021*).

52 In this paper, we present a mathematical model aimed at capturing the particular punctuated
53 evolutionary pattern of SARS-CoV-2. Our goal is to recapitulate the main features of the temporal
54 Hamming distribution observed during the COVID-19 pandemic (see Figure 1A) as parsimoniously
55 as possible in a dynamical model.

56 We show that the overall pattern can be captured by combining epistasis with heterogeneous
57 within-host evolution. The model is sufficiently general that it does not make any assumptions
58 about the biological mechanism behind saltations.

59 The proposed model is conceptually related to the NK Model of *Kauffman and Levin (1987)* and
60 *Kauffman and Weinberger (1989)* in that it operates on the space of possible genotypes, with each
61 genotype corresponding to a preassigned fitness value. This is in contrast to phenotypic fitness
62 landscape models which operate directly on the space of possible values of some quantifiable
63 trait. The most well-known among those is perhaps Fisher’s geometric model (*Fisher, 1930*) which
64 assumes a continuous phenotypic (‘trait’) space with a single optimum (*Blanquart and Bataillon,*
65 *2016*) and that the effects of single mutations are mild (*Martin, 2014*). The NK Model, a genotypic
66 fitness landscape model, instead explicitly allows for a rough (epistatic) fitness landscape. The
67 NK model, however, does not include the concept of *neutral space* – in that model, mutations are
68 generically accompanied by a change in fitness. Our model includes neutral mutations and is, in
69 that respect, closer to the models of *Koelle et al. (2006); Newman and Engelhardt (1998)*.

70 However, a crucial component of our model is the presence of sign epistasis, i.e. that the fitness
71 contribution of a point mutation may change sign (going from deleterious to beneficial or vice-
72 versa) depending on the presence of other mutations. This property turns out to offer an explana-
73 tion for the role of saltation in evolving high-fitness SARS-CoV-2 genotypes.

74 In a recent study, *Starr et al. (2022)* showed by deep mutational scanning that epistasis – includ-
75 ing sign epistasis – is an important important feature of SARS-CoV-2 evolution. As a concrete exam-
76 ple, they show that the N501Y mutation (which is present in the Alpha, Beta and Omicron variants)
77 and Q498R exhibit sign epistasis. In this case, the presence of the N501Y substitution changes the
78 contribution of Q498R from deleterious to advantageous, as measured by angiotensin-converting
79 enzyme 2 (ACE2) receptor binding affinity.

80 In general, the fitness landscape of an organism is combinatorially large, and the number of
81 possible evolutionary paths from one genotype to another fitter one is, a priori, enormous. How-
82 ever, in a seminal paper, *Weinreich et al. (2006, 2013)* showed that only very few such paths are
83 in fact accessible. The interpretation of this finding in terms of fitness landscapes is that epistasis
84 or the *ruggedness* of the landscape is highly important for understanding evolutionary trajectories
85 (*Blanquart and Bataillon, 2016*). However, even if evolutionary paths seem blocked, this conclu-
86 sion may only hold in the weak mutation limit, i.e. when the probability of multiple mutations
87 arising in the same genome within a generation is low (*Katsnelson et al., 2019*). If saltational evo-
88 lution is possible, even seemingly inaccessible regions of the fitness landscape may be explored
89 by the organism. Our model suggests that such saltations may thus increase – or, in some cases,
90 altogether enable – the emergence of new concerning variants.

91 Results

92 SARS-CoV-2 genomic diversity is characterized by punctuated evolution

93 On the basis of UK sequences (a particularly rich data set), we have computed a time-dependent
94 Hamming distribution for SARS-CoV-2, which is presented in Figure 1. Figure 1A shows the full
95 Hamming distance histogram as a function of time, from March 2020 to mid-2022, with the colour
96 and height indicating the frequency of observing genome pairs with a particular Hamming distance.
97 The peaks corresponding to saltational variant transitions are clearly visible as isolated 'islands' at
98 large Hamming distance. The insert in the same panel shows a 2D representation of the data.

99 In panel B, time series of the mean and median Hamming distances are shown, revealing clear
100 spikes associated with each of the major variant transitions to date. Each transition event is marked
101 by a very sudden spike in the typical Hamming distance, as is especially clear when considering the
102 median (dashed line). It should be noted that data quality is highest after the end of 2020, when
103 sequencing capacity was greatly increased, and before February 2022. As a concrete example,
104 4,945 sequences were included for June of 2020, while 72,292 sequences were included for the
105 month of June, 2021.

106 In Figure 1C (left), a snapshot from June 1st 2021 shows three well-defined peaks. Each peak
107 corresponds to comparisons between pairs of genomes, with the members of each pair belonging
108 to either the Alpha or Delta variant. The peak corresponding to the highest Hamming distance
109 is of course that due to comparisons between the 'new' and 'old' variant, since these are furthest
110 from each other in a genomic sense. Similarly, the plot clearly shows that variation within the Delta
111 variant is, at that point in time, much lower than within the Alpha variant, since each of the Delta
112 variant genomes belong to a clade with a recent common ancestor. In the right half of Figure 1C,
113 the situation 46 days later is shown, once the Hamming distribution has collapsed to a single peak,
114 corresponding to the then-dominant Delta variant.

115 During the month of March, 2020, the Hamming distribution appears bimodal, but there are
116 no signs of saltation. This transient bimodality, present in the early pandemic, can most clearly
117 be seen in Video 1. This can be explained by the D614G substitution, which was associated with
118 a clade that dominated from around the end of March/beginning of April 2020 (Volz et al., 2021).
119 This early, saltation-free transition is reminiscent of a result by Kauffman and Levin (1987), who
120 suggested that adaptation on a rugged fitness landscape is associated with two separate time
121 scales. First, the pathogen searches its neighbourhood in the fitness landscape until it finds a local
122 maximum. This does not require saltation and happens rather rapidly. Then, on a slower time
123 scale, the pathogen may transition to new fitness peaks by saltation.

124 Due to the relatively high quality of SARS-CoV-2 genomic surveillance in the United Kingdom,
125 both in terms of the absolute number of publicly available sequences and per capita coverage, we
126 have based the bulk of our observations on UK sequence data. However, patterns similar to those
127 presented here can be observed in US data, the analysis of which is included in Appendix 1.

128 For each day in the included range (2020-03-01 to 2022-05-10) a 7-day window (consisting of the
129 indicated day and the 6 following days) was considered. All high-quality sequences obtained within
130 that window were pooled, and a distribution of Hamming distances was compiled by repeatedly
131 picking out random pairs from the sequence pool and comparing.

132 While the Hamming distance is a somewhat crude measure of the variance between circulating
133 genomes, it turns out to offer a surprisingly powerful window into the evolution of SARS-CoV-2
134 when large amounts of sequence data are available. The transitions ancestral¹ variant→Alpha,
135 Alpha→Delta, Delta→Omicron (BA.1) as well as Omicron BA.1→BA.2 all show the tell-tale signs of
136 saltational evolution. This is perhaps especially clear when considering the median Hamming dis-
137 tance (Figure 1B, dashed line) which increases almost discontinuously at these transitions, while
138 the full distribution (Figure 1A) shows clearly defined, disconnected 'islands'. The Omicron BA.2→BA.5

¹ By ancestral variant, we mean the lineages circulating before the Alpha transition, whether including the D614G substitution or not (Hou et al., 2020; Isabel et al., 2020).

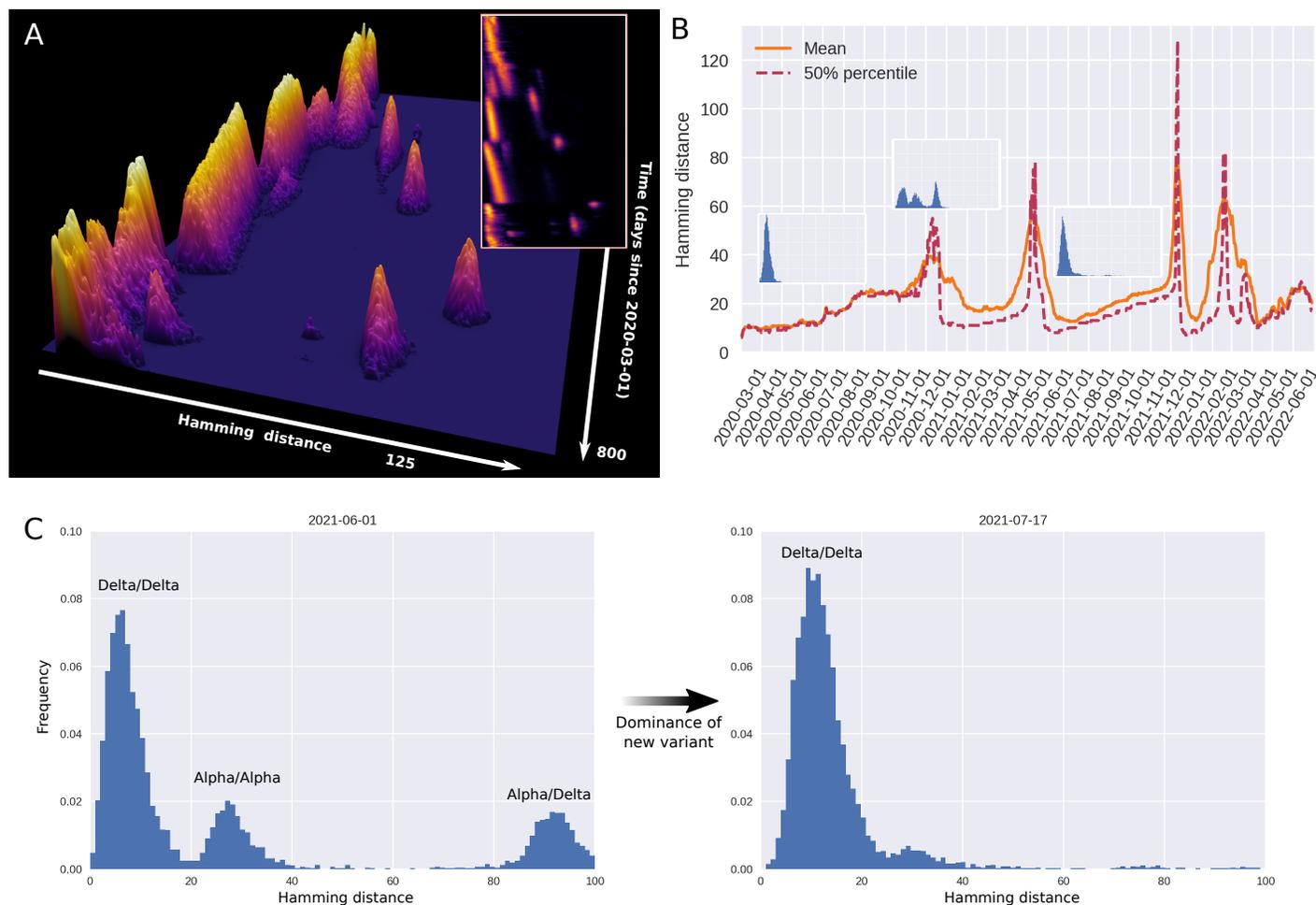
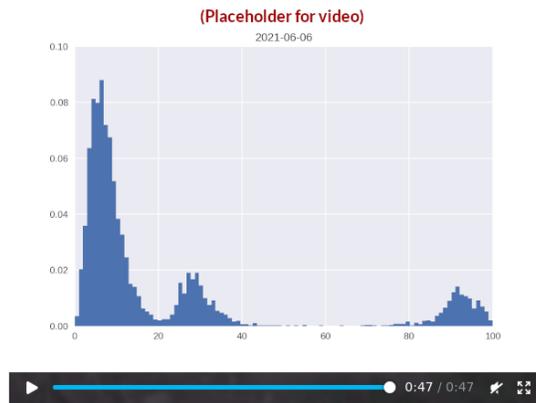


Figure 1. Genomic diversity over time in SARS-CoV-2, UK genomic surveillance data. (A) Full, time-dependent Hamming distance distribution between 2020-03-01 and 2022-05-10 (UK data, GenBank via Nextstrain (*Hadfield et al., 2018*)). Insert: Two-dimensional heatmap representation of the time-dependent Hamming distribution. **(B)** Time evolution of the mean and median Hamming distance. Each time point represents Hamming distances between genomes sampled within a one-week window beginning on that date. The three miniature inserts show Hamming distance histograms at three different time points. **(C)** Left: A snapshot of the Hamming distance distribution for genomes sampled during a one-week window starting on June 1st, 2021. The three distinct peaks correspond to the Hamming distances between pairs of genomes from each of the prevailing variants at the time, Alpha and Delta. Right: 46 days later, a single variant (Delta) dominates. See also **Video 1** for the animated Hamming histogram.

Figure 1-Figure supplement 1. Data analysis workflow.

Figure 1-Figure supplement 2. Hamming distributions for influenza H3N2 and H1N1.



Video 1. Day-by-day time development of the Hamming distribution for UK samples obtained between March 2020 and May 2022. Each snapshot is based on samples obtained within a one-week time window. **Click image to access video.**

139 transition is somewhat less clearly defined, although a moderately sized saltation does appear to
140 be present in the data. It should be noted that UK sequencing has become less dense since Febru-
141 ary 2022, explaining why the picture is somewhat murkier for the BA.2→BA.5 transition².

142 As shown in Appendix 2 Figure 1, all but one of these saltational transitions are also associ-
143 ated with a discontinuous increase in the distance to the *origin* (Wuhan-Hu-1, GenBank reference
144 sequence accession number MN908947.3). The exception is the Alpha→Delta transition, where a
145 moderate decrease is observed. In other words, the Delta variant is closer to the ancestral variant
146 than Alpha is. In Appendix 2, we model one possible explanation for this phenomenon, namely
147 the occurrence of persistent infections.

148 The plots of Figure 1 are based on the entire SARS-CoV-2 genome, meaning that a substitution
149 leading to an amino acid change in the spike protein (a major antigen) counts just as much as a
150 synonymous mutation elsewhere in the genome. In Figure 2, we probe to what extent the observed
151 drift-boom-bust pattern of diversity is driven by changes in the S-gene (coding for the spike protein)
152 or by (non-)synonymous mutations. Overall, the pattern is present whether considering only the
153 S-gene (Figure 2B), non-synonymous mutations (2C) or the entire genome (2A). We interpret this
154 to mean that

- 155 1. The drift seen between saltations is not driven solely by synonymous mutations but affects
156 the amino acid sequence as well.
- 157 2. When saltations occur, mutations are observed within the spike protein as well as outside it.
- 158 3. The observed pattern is quite robust, being observed within the whole genome, in the amino
159 acid sequence as well as within the S-gene itself.

160 It is notable, however, that the S-gene does not undergo quite as much drift as the whole genome,
161 relatively speaking. That is to say, when the whole genome is considered, the Delta→Omicron jump
162 is associated with a peak that is approximately 5.5 times larger than the typical Hamming distances
163 in the weeks that preceded it, while the ratio is closer to 11 for the spike protein. We interpret this
164 to mean that, while the S-gene is subject to large saltations, it undergoes less drift than an average,
165 similarly sized section of the genome.

166 Mechanistic modelling captures the essential dynamics

167 Our goal is to capture the overall temporal pattern of diversity observed in Figure 1 in a mathe-
168 matical model that is as parsimonious as possible. The model we have formulated consists of two
169 parts: a branching process and an evolutionary algorithm incorporating sign epistasis and salta-

²The BA.2→BA.5 transition was also muddled somewhat by the B.1.1.529 subvariant briefly making up as much as 10% of UK sequences (Hodcroft, 2021)

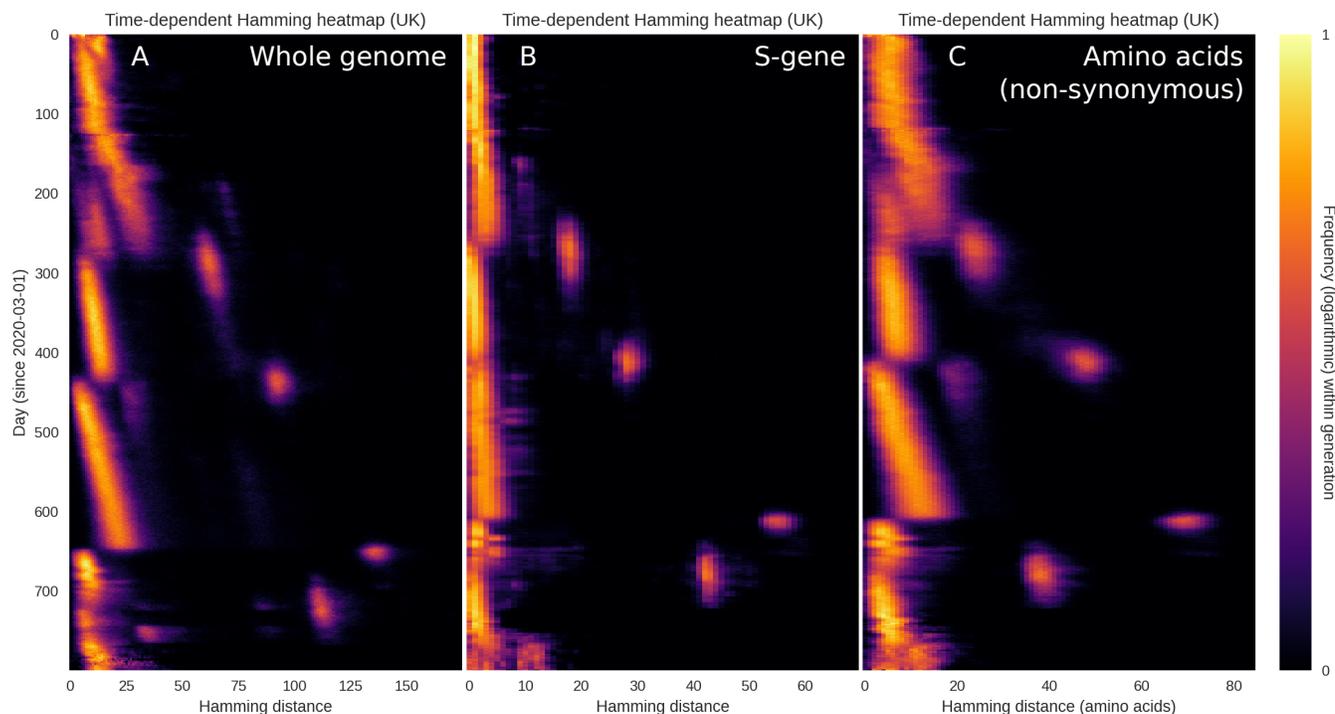


Figure 2. Restricting the Hamming distribution to the S-gene or the amino acid sequence. The overall temporal pattern of diversity seen in Figure 1 is found to persist when the analysis is restricted to the S-gene or non-synonymous mutations. **(A)** Temporal Hamming distance distribution based on the whole genome, included for reference. For each time on the vertical axis, the colour encodes the histogram of Hamming distances between genomes sampled within a one-week window starting on that date. **(B)** Temporal Hamming distance distribution confined to the S-gene which encodes the SARS-CoV-2 Spike protein. **(C)** Time evolution of the Hamming distance distribution as measured by the number of amino acid changes. We use this as a proxy for non-synonymous mutations, since a synonymous mutation would not produce an amino acid change.

170 tional evolution. The details of both of these elements can be found in the Methods and Materials
171 section. See also Figure 3–Figure supplement 1 for a schematic description of the model elements.

172 The model assumes the existence of a number of possible high-fitness genotypes, but that each
173 of them are ‘screened’ by epistasis. From a fitness landscape viewpoint, this can be thought of as a
174 landscape with a number of peaks, each of which is surrounded by a fitness trough or valley. The
175 extent of sign epistasis is then determined by the depth (and width) of these valleys.

176 To get from a moderate-fitness genotype to a local fitness peak, it is thus necessary to either tra-
177 verse a region of low fitness, with its potential for extinction, or to somehow jump across that
178 valley.

179 Evolutionary models typically assume that the ‘weak mutation limit’ holds, meaning that the
180 probability of several mutations arising in one genome in one generation is negligible (*Katsnel-
181 son et al., 2019*), leading to gradual evolution. However, as described in the introduction, there
182 are several mechanisms which can introduce a sudden burst of novelty within a single host, in-
183 cluding by recombination (*Burel et al., 2022; Duerr et al., 2022; Jackson et al., 2021*). The most
184 well-documented is perhaps elevated mutation in immunocompromised individuals. Our model,
185 however, is agnostic with respect to the precise etiology, but includes saltation simply as rare oc-
186 currences of drastically increased evolution within a single host.

187 As shown in Figure 3, the model replicates the main features observed in Figure 1, including
188 the long periods of drift (linearly increasing pairwise Hamming distances) punctuated by rapid
189 rises and subsequent collapses of diversity. Just as in the empirical data, each variant transition is
190 accompanied by three distinct peaks in the Hamming distribution.

191 The pattern shown in Figure 3 is the typical outcome of a model simulation, but occasional co-
192 existence of two variants does occur in the model, see Figure 3–Figure supplement 2. This happens
193 when two distinct variants with the same fitness happen to arise close to each other in time.

194 In the interest of simplicity, we have assumed an epidemic of constant size (constant incidence³),
195 however we explore the consequences of relaxing this assumption in the section *Epidemic dynamics
196 and spatial structure*.

197 If epistasis and saltation are turned off, evolution and variant transitions still happen within
198 the model. The temporal pattern changes, however. In Figure 4, we explore this regime by set-
199 ting ϵ , the frequency of saltations, to zero and letting $\delta R_L = 0$, thus disabling sign epistasis. The
200 resulting behaviour is characterized by periods of increasing diversity – essentially, genetic drift –
201 interrupted by sudden collapses of the typical Hamming distance. No sudden spikes are seen in
202 Figure 4A, rendering the dynamics fundamentally different from that of Figures 1 and 3. The be-
203 haviour observed in this regime is more reminiscent of the dynamics observed for H3N2 influenza
204 in *Koelle et al. (2006)*. However, one could object that the temporal resolution of the empirical
205 time series shown in *Koelle et al. (2006)* is not sufficiently high to allow one to discriminate be-
206 tween the scenarios of our Figures 3 and 4 – after all, the periods of drastically increased pairwise
207 nucleotide Hamming distance seen in Figure 1 are brief and require high temporal resolution to
208 discern. While the amount of genomic data available for SARS-CoV-2 enables this, the picture is
209 murkier for seasonal influenza. In Figure 1–Figure supplement 2, we present the result of applying
210 the analysis of Figure 1A to influenza types H3N2 and H1N1. While there is no apparent evidence
211 of saltation, the available data is relatively coarse-grained.

212 Another influential evolutionary model of influenza is due to *Ferguson et al. (2003)*. In their
213 model, the appearance of new variants is driven by immune system memory and a non-linear
214 relation between Hamming distance and cross-immunity, the latter in the form of short-lived strain-
215 transcending immunity. While a sensible model for seasonal influenza, it gives rise to diversity
216 dynamics which are closer to Figure 4 than to the pattern observed for SARS-CoV-2.

217 In the simulations of Figure 4, saltation and epistasis are completely lacking, but in Figure 4–
218 Figure supplement 1, we consider what happens if some saltational evolution *does* occur, without

³Note that time in the model is measured in generations, meaning that constant incidence and prevalence both hold.

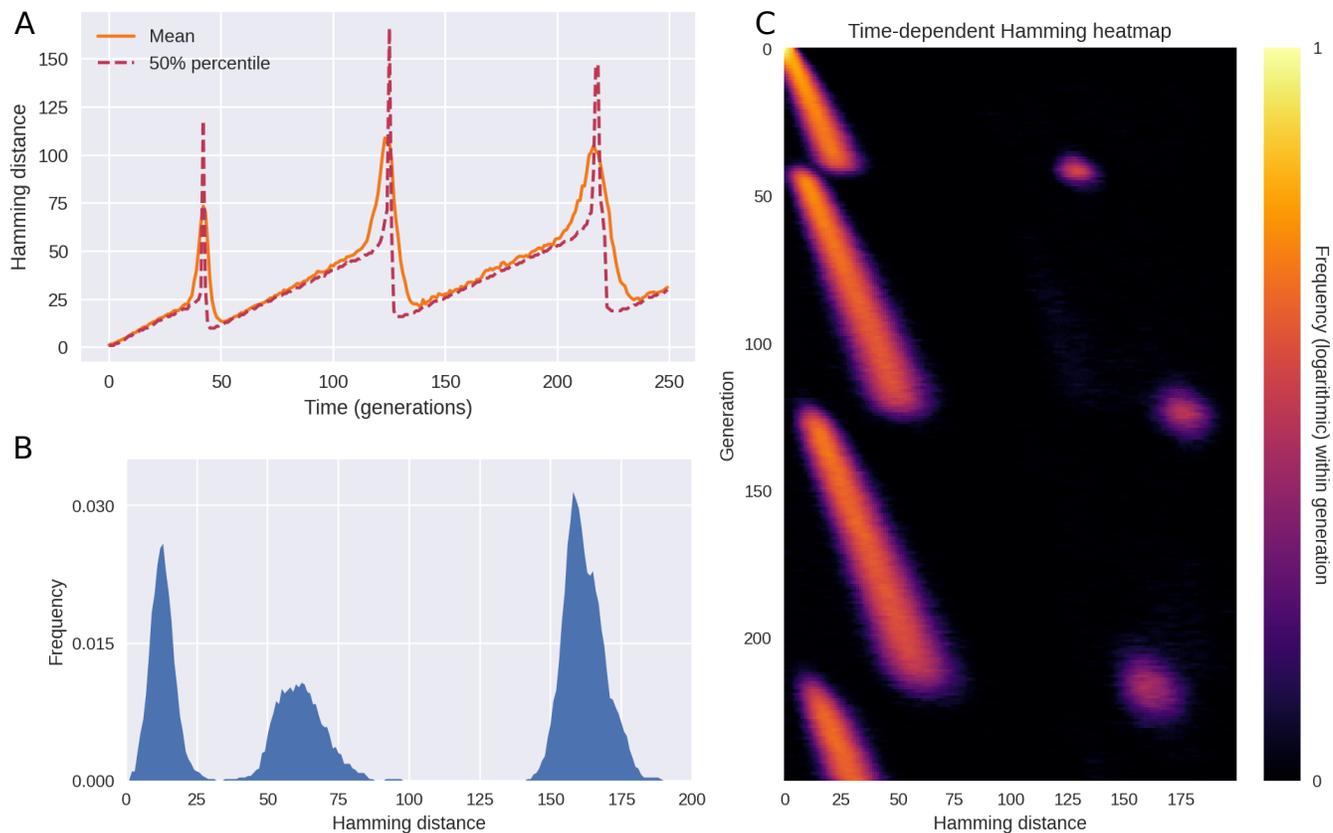


Figure 3. Simulated outbreak with saltation (heterogeneous mutation rates) and epistasis. (A) Time evolution of the mean and median Hamming distance between bitstring genomes present in any given generation of the model simulation. The pattern of genetic drift punctuated by sudden increases and subsequent collapses in diversity is similar to what is observed in SARS-CoV-2 (see Figure 1). **(B)** A snapshot of the Hamming distance distribution in generation $t = 218$ of the simulated outbreak. Just as in Figure 1, the three distinct peaks correspond to the distances between pairs of genomes from each of the two prevailing variants at the time. **(C)** Time evolution of the full Hamming distance distribution. For each generation on the vertical axis, the colour encodes the histogram of Hamming distances between genomes within that generation. The parameters used in these simulations were $\epsilon = 0.001$, $d_0 = 3$, $\delta R_H = 1.0$, $\delta R_L = -\infty$ (i.e. deleterious mutations were fatal to the pathogen).

Figure 3–Figure supplement 1. Model schematics.

Figure 3–Figure supplement 2. Coexistence of equally fit variants.

219 sign epistasis. Qualitatively, the picture most resembles the saltation-free scenario of Figure 4, but
 220 occasional Hamming spikes are observed. Overall, this scenario does not conform to the empirical
 221 observations in the form of Figure 1. In the next section, we systematically probe how different
 222 levels of epistasis and saltation affects the evolution of new, highly transmissible variants.

223 Our focus is mainly on the dynamics of diversity, and for this reason we have emphasized the
 224 distribution of Hamming distances between viral genomes present in the population at the same
 225 time. This goes for the empirical observations (Figure 1) as well as our model simulations (Figure
 226 3). However, in Appendix 2, we explore the distributions of Hamming distance relative to the origin
 227 (meaning Wuhan-Hu-1, GenBank reference sequence accession number MN908947.3).

228 **Saltation facilitates the evolution of highly transmissible variants**

229 Saltational evolution may not only be a way to generate vastly different variants, but may indeed be
 230 necessary for the virus to evolve highly fit variants at all. In the presence of strong epistasis, gradual
 231 evolution towards a high fitness genotype can be blocked (see Figure 3–Figure supplement 1A).
 232 Conceptually, such gradual evolution under strong epistasis would correspond to traversing a deep
 233 valley in the fitness landscape by a series of small steps before reaching a peak (*Katsnelson et al.*,

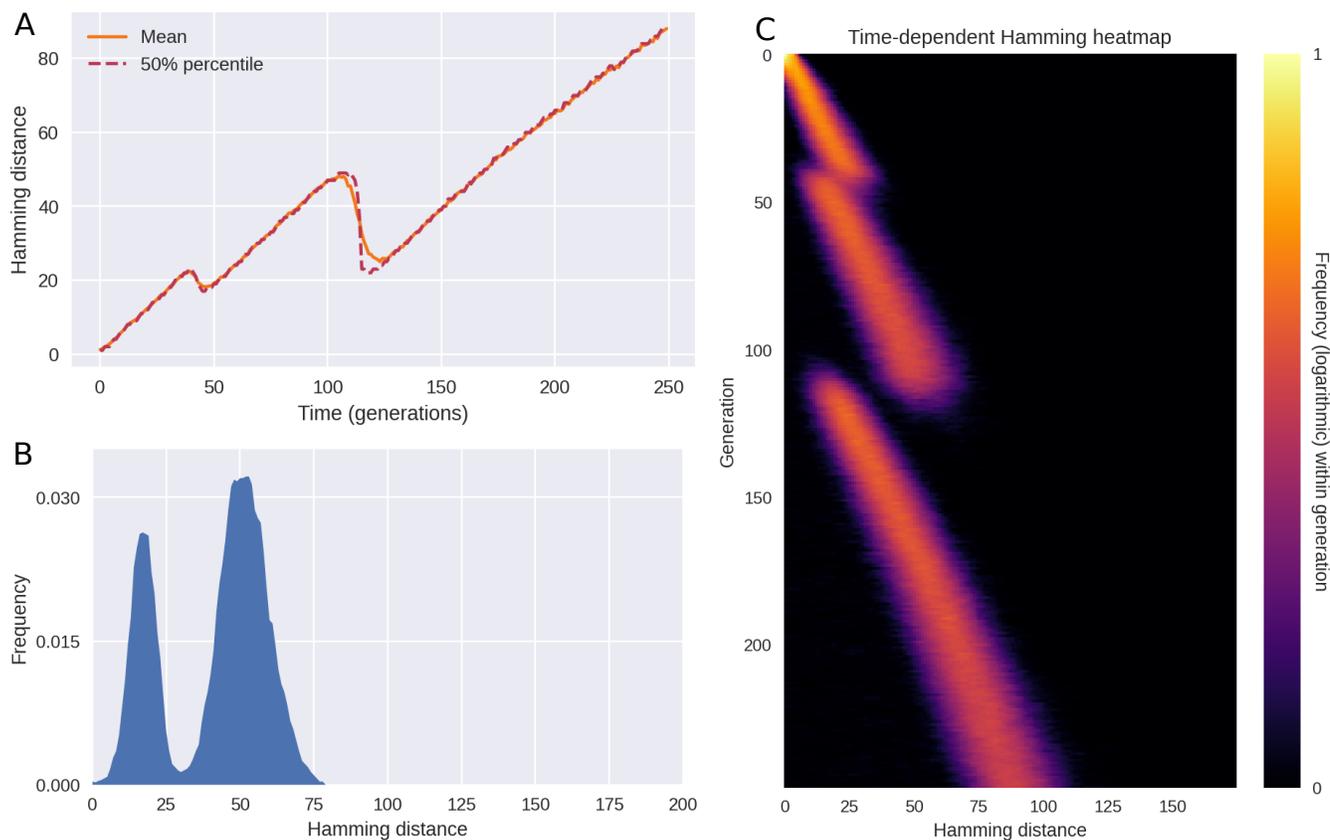


Figure 4. In the absence of epistasis and saltation, model results do not match observations. When saltations do not occur ($\epsilon = 0$) and sign epistasis is absent ($\delta R_L = 0$), the pathogen will only rarely overcome troughs in the fitness landscape. When it does, the transitions are marked by a collapse of diversity (as measured by the typical Hamming distance), giving a drift-bust-drift dynamics as opposed to the drift-boom-bust pattern seen in SARS-CoV-2. **(A)** Time evolution of the mean and median Hamming distance between genomes present in any given generation of the model simulation. **(B)** A snapshot of the Hamming distance distribution for bitstring genomes at generation $t = 112$ of the simulated outbreak. **(C)** Time evolution of the Hamming distance distribution. For each generation indicated on the vertical axis, the colour encodes the histogram of Hamming distances between genomes within that generation.

Figure 4-Figure supplement 1. Saltational evolution in the absence of sign epistasis.

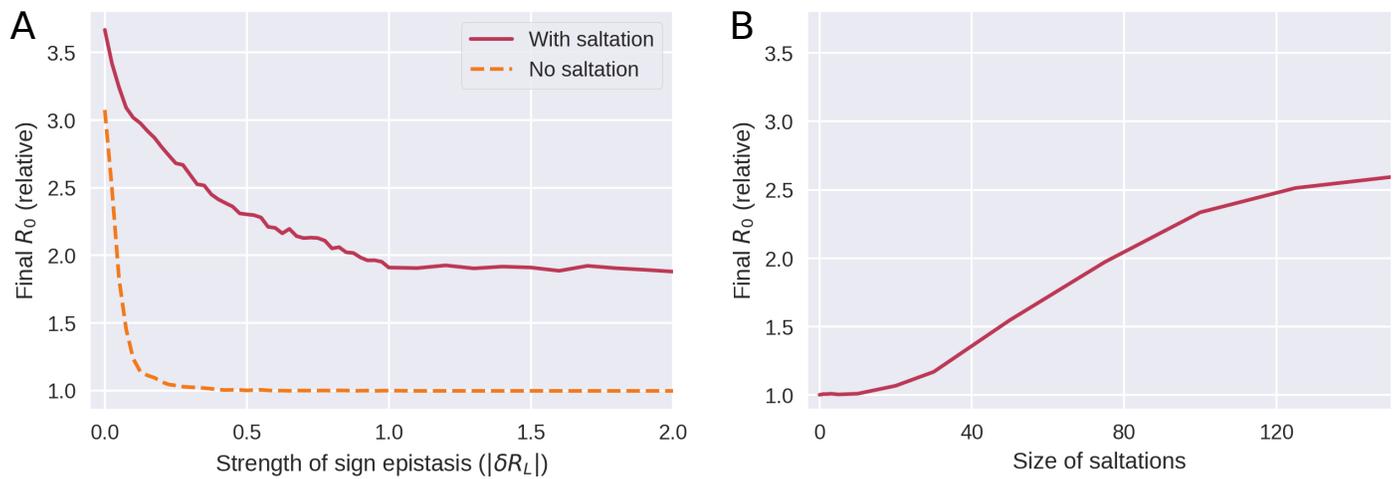


Figure 5. Saltation allows highly transmissible variants to evolve by facilitating evolution across fitness landscape troughs. **A)** Evolution under varying degrees of sign epistasis. The vertical axis indicates the final average reproductive number in the model population after 300 generations of the simulation, relative to (divided by) the reproductive number of the initial variant. The horizontal axis indicates the depth of a valley in the fitness landscape, $|\delta R_L|$, understood as the reduction in reproductive number suffered due to a deleterious configuration. Here, δR_L was distributed according to a Dirac δ distribution and as such its value was deterministic. This panel is based on 90,000 simulations and the parameters used were $d_0 = 3$ and $\delta R_H = 1$. **B)** Evolution with varying degrees of saltation. Moderate sign epistasis is assumed ($\delta R_L = -0.5$). All other parameters are as in panel A. This panel is based on 7600 simulations.

234 **2019; Smith and Ashby, 2022)**. However, such a fitness valley indicates the presence of deleterious
 235 mutations which impart a high probability of extinction of the lineage in question, preventing the
 236 fitness peak from being reached.

237 In Figure 5, we explore how the strength of epistasis and the size of saltations affect the ability
 238 of the pathogen to evolve new, highly transmissible strains. By the *strength* of epistasis, we mean
 239 the typical depth of a valley in the fitness landscape, $|\delta R_L|$, i.e. the loss in reproductive number
 240 suffered.

241 For a pathogen which does not undergo saltational evolution (Figure 5A, dashed curve), signif-
 242 icant sign epistasis ($|\delta R_L| \gtrsim 0.25$ at $d_0 = 3$ in our simulations) is a roadblock to evolution of high-
 243 fitness variants. However, a pathogen which undergoes saltation (fully drawn curve) can overcome
 244 this epistatic hindrance. Above a certain threshold (at $|\delta R_L| \approx 1$ in Figure 5A), stronger sign epista-
 245 sis ceases to further impede the emergence of high-fitness variants. The mechanism behind this is
 246 that sign epistasis becomes so strong that a fitness valley may be overcome only by pure saltation
 247 and is no longer traversable by gradual evolution or a combination of the two.

248 As shown in Figure 5B, large saltations are necessary to overcome even moderate sign epistasis,
 249 further explaining why the Hamming peaks seen in SARS-CoV-2 are so large.

250 Epidemic dynamics and spatial structure

251 In Figure 3, we made a number of simplifying assumptions, the major ones being constant preva-
 252 lence and absence of any spatial or population structure.

253 We first relax the former assumption by implementing susceptible-infected-recovered-susceptible
 254 (SIRS) dynamics. The infected individuals are now assumed to make up only a fraction of a larger
 255 population of total size N . Our aim is to ascertain whether the diversity dynamics observed in the
 256 previous section are fundamentally altered by allowing a variable number of infected individuals,
 257 $I(t)$, as well as susceptible depletion and waning immunity.

258 In Figure 6, a typical course of a simulation with SIRS dynamics, epistasis and saltation is shown.
 259 As shown in panel A, the number of recovered (immune) individuals varies non-monotonically over
 260 time, reflecting that individuals acquire immunity after being infected, and that the immunity even-

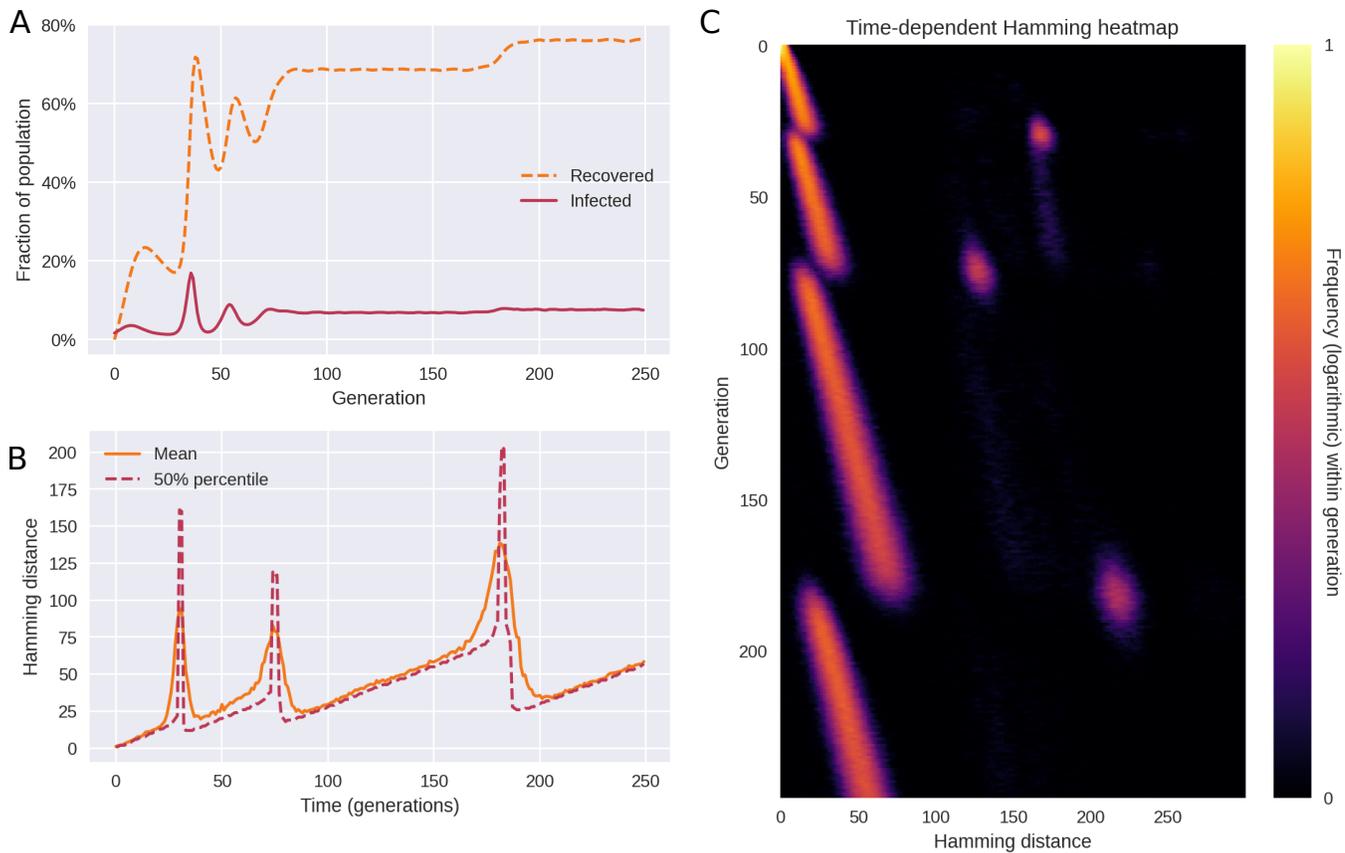


Figure 6. Saltational evolution under susceptible-infected-recovered-susceptible (SIRS) dynamics. The reproductive number of the initial variant is $R_0 = 1.2$ and immunity wanes at a rate of $\omega = 0.1$. Population size is $N = 6 \times 10^5$. The parameters of fitness-altering mutations are $\delta R_H = 1$ and $\delta R_L = -\infty$ (i.e. deleterious mutations are fatal to the organism or prevent transmission). **(A)** Time evolution of the recovered (or immune) fraction of the population. Successive variants have higher reproductive numbers (R_0), eventually leading to two different endemic plateaus. **(B)** Even under variable prevalence, the Hamming dynamics looks similar to that of Figure 3. **(C)** The full temporal Hamming distribution is characterized by the same kind of punctuated evolution as in the simpler constant-prevalence case of Figure 3.

Figure 6–Figure supplement 1. Simulations with spatial (metapopulation) structure.

261 tually wanes. However, as successive variants of greater fitness (greater reproductive number R_0)
 262 arise, an endemic plateau is eventually reached. While the epidemiology is very different from
 263 that of Figure 3, the Hamming distribution (Figure 6C) is remarkably similar. This indicates that
 264 the mechanism of saltational evolution in conjunction with sign epistasis robustly reproduces the
 265 punctuated evolutionary dynamics seen in Figure 1.

266 In simulations with variable incidence, higher incidence translates to an increased risk of emer-
 267 gence of new variants, all else being equal. Since saltations are simulated as a constant-rate (Pois-
 268 son) process for each infected individual, the risk of emergence scales with the number of infected.
 269 Since simulations are stochastic, this tendency is not necessarily clear from a single realization
 270 such as Figure 6. A similar frequency dependence is likely to hold for SARS-CoV-2, since rare occur-
 271 rences in terms of within-host evolution are proportionally more likely to be observed with higher
 272 incidence.

273 Next, we probe the impact of spatial separation on the diversity dynamics. Spatial structure is
 274 implemented by augmenting the model with a metapopulation element, see Materials and Meth-
 275 ods for details. We find that the main effect of space is to protract transitions, such that the
 276 transient multimodality of the Hamming distribution lasts longer. The duration of coexistence of
 277 strains with different fitnesses is observed to be determined by the transmission rate β_{ij} between
 278 different populations (i.e. with $i \neq j$). In Figure 6–Figure supplement 1, we probe three situations

279 where (relative) inter-population transmission rates are either 0, 10^{-4} or 10^{-3} (with intra-population
280 transmission rates $T_{ii} \approx 1$). We find that with very low ($< 10^{-3}$) transmission rates between popu-
281 lations, spatial structure leads to drawn-out transitions, but that this effect disappears as soon as
282 significant transmission between populations occur. This intuitively makes sense, since the within-
283 population transmission rate will dominate as soon as even a few cases of a new variant have
284 spilled into a population.

285 It is worth noting that spatial structure in itself cannot lead to the saltational signature seen in
286 Figure 1. By this we mean that a pathogen which evolves gradually (i.e. obeys the weak mutation
287 limit) in multiple spatial patches will not lead to a sudden spike in the Hamming distance once spill-
288 over happens. The reason for this is that two spatially remote lineages diverge from each other (in
289 terms of Hamming distance) at approximately the same rate as less geographically distant pairs of
290 lineages.

291 Discussion

292 The pattern of evolution observed in SARS-CoV-2 suggests that transmissibility of the pathogen has
293 mainly increased due to large evolutionary ‘jumps’, rather than due to gradual evolution, something
294 that may turn out to be a signature feature of the pathogen. Our model simulations highlight how
295 this preference for adaptation by saltation may be explained by an ability to overcome epistatic
296 ‘fitness valleys’. The implications for public health are clear; any situation which facilitates such
297 jumps should be treated with heightened awareness. They represent a high risk for the emer-
298 gence of new, concerning variants which could not have emerged through gradual evolution.

299 While much attention has been given to the role of immunocompromised individuals, for good
300 reasons, it is important to realize that other probable mechanisms of saltation exist. For instance,
301 consider reverse zoonosis – the transmission from humans to animals. The epistatic landscape
302 may be very different in animals, affording a way of bypassing what would otherwise be troughs in
303 the human SARS-CoV-2 fitness landscape. Reintroduction of the mutated lineage into the human
304 population would then constitute a ‘jump’ in terms of Hamming distance, and potentially also phe-
305 notypically. An example of such back-and-forth transmission between human and animal hosts
306 leading to a large number of novel mutations was the so-called Cluster 5 variant, which evolved in
307 mink (*Neovison vison*) in Denmark and subsequently spread to humans (*Hammer et al., 2021*). This
308 mink-derived variant, which was only one of several which escaped into the human population,
309 exhibited 35 substitutions and four deletions in the spike protein alone (*Larsen et al., 2021*).

310 From a public health perspective, these possible mechanisms have one important thing in com-
311 mon; they underscore the importance of widespread and equitable distribution of up-to-date vac-
312 cines, since saltational evolution in disadvantaged or remote populations carries a risk of emer-
313 gence of new, highly transmissible variants.

314 The plurality of potential etiologies highlights the need for comprehensive research into the mech-
315 anisms which may underlie the observed saltational evolution. Such studies would be most wel-
316 come and would have to consider multiple scales, from molecular mechanisms and within-host
317 evolution to the epidemiological dynamics which may contribute to saltations.

318 The type of analysis performed in this study requires large amounts of sequence data, beyond
319 what could usually be obtained for infectious diseases prior to COVID-19. As shown in Figure 1–
320 Figure supplement 2, a similarly clear and detailed distribution of nucleotide distances could not
321 be obtained for influenza H1N1 or H3N2. This is just one example of how incredibly useful the
322 high level of genomic surveillance achieved for SARS-CoV-2 is, and more generally highlights the
323 potential that extensive sequencing of pathogens holds for advancing phylodynamic understand-
324 ing across pathogens (*Grenfell et al., 2004*). While many countries have since scaled down the
325 level of testing and sequencing of SARS-CoV-2, scientific insights based on this data will no doubt
326 continue to emerge and have a lasting impact on our understanding of pandemics more broadly.

327 In our simulations, we have not explicitly modelled any effects of immune memory. We have
328 allowed for new variants with higher effective reproduction numbers to arise, but have not dis-

329 tingished between whether that advantage stemmed purely from higher infectiousness or from
330 some degree of immune escape. However, it is worth noting that the empirical pattern of punctuated evolution held for every major transition (Figure 1), including the transition to the Alpha variant. When the Alpha variant became dominant, only about 3 infections per 100 people had been recorded in the United Kingdom (*Ritchie et al., 2022*) and vaccinations had not yet begun in earnest. While this is surely an undercount, a general depletion of susceptibles was not a main driver for the success or emergence of the Alpha variant. As such, the punctuated evolutionary pattern does not seem to be hinged on a connection between Hamming distance and evasion of immunity. In the case of the transition to Omicron, immune escape certainly played a role (*Meng et al., 2022; Zhang et al., 2022*), but it would seem that the mechanism of punctuated evolution is more general than that. In *Meijers et al. (2022)*, the authors explicitly decompose fitness advantages into intrinsic and antigenic. Introducing a similar distinction in a genotypic fitness landscape model with saltation would be an interesting extension of the present work.

342 Even in the absence of explicit modeling of immune memory and partial cross-immunity, there are good reasons to believe that saltations will continue to play a role in facilitating the emergence of new variants. As described by *Plotkin et al. (2002)*, accumulating immunity changes the fitness landscape of a pathogen over time, lowering some fitness peaks while rendering other peaks relatively more advantageous for the virus. Saltations can then enable the pathogen to reach those fitness peaks. Indeed, it is plausible that high levels of (more or less strain-specific) immunity in a population may increase the rate at which new strains emerge by saltation. Such a connection would further underscore the importance of broadly effective and widely available vaccines as well as any measures which decrease the likelihood of accelerated evolution within hosts, with its risk of seeding saltation events.

352 As a consequence of the parsimony of our model, we have not explicitly modelled recombination events, but rather assume that each multi-site jump involves a random set of sites. Recombination has been reported in SARS-CoV-2, including – but not limited to – in conjunction with treatment of immunosuppressed patients (*Burel et al., 2022; Duerr et al., 2022; Varabyou et al., 2021; Focosi and Maggi, 2022*). Future work could explore the implications of allowing for recombination events in this type of model.

358 The influenza model of *Koelle et al. (2006)*, which gives rise to Hamming dynamics reminiscent of the saltation-free simulations of Figure 4, does so in a very different way. There, it is assumed that the pathogen explores a *neutral network* (a set of antigenically and fitness-wise equivalent genotypes which are connected by one-mutation neighbours (*Newman and Engelhardt, 1998*)) in the vicinity the prevailing strain. This goes on until the ‘random walk’ happens upon a configuration which is substantially antigenically different from the prevailing cluster, albeit connected to it by a single mutation. Once this happens, a new cluster emerges which has only limited cross-immunity with the prevailing strain. Since all steps along the way are small, the new variant will be very close (genotypically) to a member of the previous cluster. Consequently, this type of dynamics does not produce abrupt spikes in Hamming distance, such as the ones shown in Figures 1 and 3.

368 There are a few models in the literature that seek to address the connection between saltation, epistasis and the likelihood of emergence of new variants (*Katsnelson et al. (2019)* and *Smith and Ashby (2022)*, the latter of which is based on the model by *Gog and Grenfell (2002)*). However, in contrast to existing theoretical studies, we address the empirical temporal development of diversity and propose a model which can directly replicate the main features of that distribution.

373 We have focused on capturing the main features of the evolution of SARS-CoV-2 as parsimoniously as possible and although we have explored a number of biologically motivated extensions, our model still represents a theoretical foundation upon which more sophisticated models can be built. There is much to be done in terms of understanding and modelling the precise fitness landscape of SARS-CoV-2, including its dependence on host immunity history. More broadly, an increase in genomic surveillance across multiple pathogens will doubtlessly lead to new insights into the diversity dynamics of other pathogens. This would not only enable research into the evolution

380 of individual pathogens, but allow us to question how co-circulating pathogens affect the diversity
381 dynamics of one another.

382 **Methods and Materials**

383 **Temporally resolved Hamming distributions from sequence data**

384 In this section, we describe the data processing workflow which was used to generate the Hamming
385 distance plots of Figures 1 and 2. We have used the open, GenBank-derived dataset of aligned
386 SARS-CoV-2 sequences from Nextstrain (*Hadfield et al., 2018*). In our main analysis, we have used
387 UK sequences from the 1st of March 2020 onwards. See also Figure 1–Figure supplement 1 for an
388 illustration of the following workflow:

389

- 390 • For each day t in the interval:
 - 391 – Select 5,000 random pairs of whole-genome sequences (i.e. 10,000 sequences) obtained
 - 392 within a 1-week time window starting on day t .
 - 393 – For each pair of sequences s_i and s_j :
 - 394 * Go through both sequences, site by site, and record the number of differences be-
 - 395 tween them, H_{ij} . This is the pairwise Hamming distance.
 - 396 – Compute a probability density/histogram $p_t(H)$ based on the observed Hamming dis-
 - 397 tances $\{H_{ij}\}$.

398 It is then this function, $p_t(H)$, that is plotted in Figure 1A. In practice, we have used the metadata pro-
399 vided by Nextstrain, which contains fields describing the nucleotide differences relative to the ref-
400 erence strain Wuhan-Hu-1 (GenBank reference sequence accession number MN908947.3), rather
401 than operating directly on the whole-genome sequences. Numerically, this makes no difference,
402 but it affords a large increase in performance, since it allows us to avoid processing unchanged
403 regions of the genome, which do not contribute to the Hamming distance.

404 For Appendix 2 Figure 1, which instead shows the distance to the reference sequence (the 'ab-
405 solute' Hamming distance), the above workflow is slightly altered:

- 406 • For each day t in the interval:
 - 407 – Select 5,000 random whole-genome sequences obtained within a 1-week time window
 - 408 starting on day t .
 - 409 – For sequences s_i :
 - 410 * Go through the sequence, site by site, and record the number of differences H_i
 - 411 between s_i and the reference sequence. This is the absolute Hamming distance.
 - 412 – Compute a probability density/histogram $p_{0,t}(H)$ based on the observed Hamming dis-
 - 413 tances $\{H_i\}$.

414 It is then this function, $p_{0,t}(H)$, that is plotted in Appendix 2 Figure 1.

415 **Branching model with saltational evolution**

416 The mechanistic model developed for this study is a discrete-time branching model coupled to a
417 genotypic fitness landscape model.

418 In the simulations of Figures 3 and 4, we assumed a constant prevalence, for simplicity. This
419 amounts to keeping the mean effective reproductive number across the population at unity. In
420 Figure 6 we relax this assumption and explore a version of the model with epidemic dynamics.
421 We start by documenting the constant-prevalence version of the model, as well as the genotypic
422 fitness landscape element, before we go on to describe how we incorporate SIRS dynamics and
423 spatial structure.

424 Evolutionary branching model with constant prevalence

425 In the model, each new generation of infections consists of a fixed number of individuals, N , and
426 generations do not overlap. Consequently, there are N infected individuals at any given time. Each
427 infected individual i has an associated bit-string G_i of length L , representing the genome of the
428 pathogen. We do not explicitly model any within-host diversity, as we are only interested in the
429 genome of the pathogen that is eventually transferred during transmission.

430
431 At each time step (corresponding to one generation), a new random individual i is repeatedly
432 selected and allowed to infect a number z_i of new individuals, selected from a Poisson distribution
433 with mean R_i , i.e. $z_i \sim \text{Pois}(R_i)$. This continues until a total of N new transmissions have occurred
434 in that generation, ensuring that the prevalence is kept constant. At transmission, the pathogen
435 genome of the infector is copied to the infectee. The personal reproductive number R_i is deter-
436 mined by the fitness of the bit-string G_i , the details of which are discussed in the next subsection.
437 In each newly infected individual, there is a risk of mutation. The number of point mutations m_i
438 that occur within the i 'th host is drawn from a distribution. In the case of homogeneous mutation
439 (i.e. absence of saltation), m_i is drawn from a Poisson distribution characterized by a mutation
440 rate $\mu_0 < 1$. Saltation, on the other hand, is simulated by drawing m_i from a bimodal distribution
441 characterized by two different mutation rates/sizes μ_0 and μ_1 , ensuring that an outsized amount
442 of mutation can take place within a single host on rare occasions. Concretely, we have used the
443 distribution $P_s(m)$ given by:

$$P_s(m) = (1 - \epsilon)\text{Pois}(m; \mu_0) + \epsilon U(m; \mu_1 \pm \Delta\mu) \quad (1)$$

444 Where $U(m; \mu_1 \pm \Delta\mu)$ is the uniform distribution centered on μ_1 with half-width $\Delta\mu$, and $\text{Pois}(m; \mu_0)$
445 is a Poisson distribution with mean μ_0 . $\epsilon \ll 1$ is a small dimensionless quantity measuring the fre-
446 quency of saltational mutation. The parameter μ_0 gives the rate of non-saltational mutation while
447 μ_1 is the typical size of a saltation.

448 We use this simple bimodal distribution out of convenience, but our results do not change qualita-
449 tive if another bimodal distribution is used.

450 Once the quantity m_i has been drawn, a number m_i of random bit flips are then performed in the
451 genomic bitstring G_i , each flip corresponding to a point mutation.

452 Modelling sign epistasis

453 Before simulations start, a number of N_e of 'epitopes' (regions in the genome on which fitness
454 depends), each of length L_e , are designated. We assume non-overlapping epitopal regions and
455 thus require $L_e N_e \leq L$.

456 Within each epitope, a number N_H of highly fit combinations are assigned. We have assumed
457 $N_H = 1$ for all of our simulations, but since the general N_H case is no more complicated, we include
458 the parameter here. The fitness of each combination is measured in terms of its contribution δR_H
459 to the individual reproductive number. In general, δR_H for each combination may be drawn from
460 a distribution $P_H(\delta R_H)$ to allow for a variety of combinations with different fitness values.

461 Tunable sign epistasis is modeled by assigning a fitness contribution $\delta R_L \leq 0$ to each combination
462 which lies within a Hamming distance d_0 of a high-fitness combination. The overall fitness of a
463 given genotype is then obtained by adding up the contributions for each of the N_e epitopes:

$$R_0 = R_0^{\text{initial}} + \sum_{i=1}^{N_e} \delta R_i, \quad (2)$$

464 with the constraint that $R_0 \geq 0$. In practice this constraint is enforced by letting

$$R_0 = \max\left(0, R_0^{\text{initial}} + \sum_{i=1}^{N_e} \delta R_i\right), \quad (3)$$

Table 1. Model parameters and their values.

Parameter	Description	Value (base case)
L	Genome length (bits)	1000
L_e	Length of each epitopal sequence	5
N_H	Number of highly fit configurations of each epitope	1
N_e	Number of epitopes	5
d_0	Hamming distance within which genotype is deleterious	3
$\langle \delta R_H \rangle$	Avg. fitness advantage (change in R) due to beneficial genotype	1
$\langle \delta R_L \rangle$	Avg. fitness contribution (change in R) due to deleterious genotype	-1
μ_0	Base mutation rate (for whole genome)	0.3
ϵ	Frequency of saltation	0 or 0.0001
μ_1	Typical size of saltations	150
$\Delta\mu$	Half-width of saltation size distribution	50
T_{ij}	Relative transmission rate betw. populations i, j Subject to $\sum_j T_{ij} = 1$.	0-1
ω	Rate of waning of immunity in SIRS simulations	0.1

465 High sign epistasis is then achieved when $d_0 > 1$ and $\delta R_L \ll 0$. However, the model also allows
 466 for incomplete or partial sign epistasis: if δR_L for each combination is drawn from a distribution
 467 $P_L(\delta R_L)$ which has support at $\delta R_L = 0$, then each peak in the fitness landscape will not be com-
 468 pletely surrounded by troughs. In other words, in that case it may be possible to evolve to a highly
 469 fit variant through a series of single point mutations without suffering decreased fitness in the pro-
 470 cess.

471 Unless otherwise specified, we run our simulations with the parameter values given in Table 1

472 Incorporating SIRS dynamics

473 In the simulations of Figures 3 and 4 we assumed constant incidence, meaning that the number
 474 of infected within any given generation was $I(t) = I_0$ with I_0 a constant (thus, prevalence was con-
 475 stant as well). However, to relax this assumption we incorporate susceptible-infected-recovered-
 476 susceptible (SIRS) dynamics.

477 We continue to model the infected individuals by a branching process, but now also keep track of
 478 the number of susceptible ($S(t)$) and recovered ($R(t)$) individuals in each generation t of the out-
 479 break.

480 The simulations proceed as follows. At time $t = 0$, let:

$$S(t = 0) = N - I_0,$$

$$I(t = 0) = I_0,$$

$$R(t = 0) = 0.$$

481 Then, at each subsequent time step t , the transmission probability is modulated by a factor $S(t)/N$
 482 representing susceptible depletion. Whenever transmission does occur, the number of suscepti-
 483 bles decrease by one while the number of infected individuals increase correspondingly.

484 When an individual is recovered (i.e. after one generation of being infected), $I(t)$ is decreased by
 485 one while $R(t)$ increases by one.

486 Each recovered person has a constant probability rate ω for becoming susceptible once again. In
 487 other words, this is modeled as a Poisson process with rate ω . Note that this not only corresponds
 488 to waning of immunity, but also to any other mechanism by which a recovered individual may
 489 become replaced by a susceptible one (such as population turnover). However, we will refer to
 490 ω as the rate of waning. In our simulations (Figure 6 and Figure 6–Figure supplement 1), we set
 491 $1/\omega = 10$ meaning that duration of immunity averages 10 generations. This figure is not supposed

492 to reflect any particular value for SARS-CoV-2, but is rather used to illustrate the robustness of the
493 pattern of punctuated evolution to waning immunity. In the interest of simplicity, we have ignored
494 any seasonal effects on transmission. We consider this a reasonable simplification, both due to
495 the conceptual nature of our model and the understanding that susceptible dynamics rather than
496 seasonality is the major limiting factor in the pandemic phase (*Baker et al., 2020*).

497 **Spatial structure**

498 We implement a minimal model of spatial structure by incorporating a metapopulation element.
499 Let there be n_{pops} populations, each with total population N_i ($i \in \{1, \dots, n_{\text{pops}}\}$). At time $t = 0$, let the
500 number of susceptible, infected and recovered individuals in each population be given by:

$$S_i(t = 0) = N_i - I_{i,0},$$

$$I_i(t = 0) = I_{i,0},$$

$$R_i(t = 0) = 0.$$

501 In our simulations (Figure 6–Figure supplement 1) we assume identical population sizes, $N_i =$
502 N/n_{pops} , and an initial equipartitioning of infected individuals $I_{i,0} = I_0/n_{\text{pops}}$, where $N = \sum_i N_i$ and
503 $I_0 = \sum_i I_{i,0}$.

504 The transmission rate between populations is then determined by the matrix elements $\beta_{ij} = \beta T_{ij}$
505 where each element T_{ij} gives the relative transmission rate from population i to j and β represents
506 the transmissibility of the strain the infected individual carries. We assume that \mathbf{T} is a symmetric
507 matrix, $T_{ij} = T_{ji}$.

508 In Figure 6–Figure supplement 1 we took \mathbf{T} to have the following form:

$$T = \begin{bmatrix} 1 - \varepsilon & \varepsilon & 0 \\ \varepsilon & 1 - 2\varepsilon & \varepsilon \\ 0 & \varepsilon & 1 - \varepsilon \end{bmatrix} \quad (4)$$

509 with ε at either 0 (panel A), 10^{-4} (panel B) or 10^{-3} (panel C). This corresponds to a linear layout of
510 three populations, with transmission occurring only between adjacent compartments.

511 **Modelling decreasing absolute Hamming distance**

512 As described in Appendix 2, the typical Hamming distance between circulating genomes and the
513 ancestral variant is not necessarily monotonically increasing with time. We call this distance the
514 *absolute* Hamming distance, in contrast to the pairwise distance between concurrently circulating
515 genomes which we call the *relative* Hamming distance (to reflect that the absolute Hamming dis-
516 tance is measured with respect to a fixed point in the genomic space).

517 We begin by describing a very simple variation upon the model which has the effect of allowing the
518 absolute Hamming distance to decrease (as well as increase) at variant transitions. In this section,
519 we assume a constant prevalence of N infected individuals.

520 Assume that a fraction p_d of the first generation (i.e. $p_d N$ individuals) have prolonged infections,
521 lasting τ_d typical generations before onward transmission. Assume furthermore that mutations
522 happen at a rate μ_d for these individuals, such that a number of point mutations, $\mu_d \tau_d$, occurs
523 before onward transmission. Here, μ_d is the mutation rate associated with these prolonged infec-
524 tions. The rest of the population is assumed to be homogeneous with respect to the occurrence of
525 mutations, all possessing a mutation rate μ_0 . We draw τ_d from a uniform distribution with support
526 throughout the entire simulation (which is assumed to have duration t_f), $\tau_d \sim U(\tau_d; t_f/2 \pm t_f/2)$.
527 Furthermore, the fitness advantage δR_H of different epitope configurations was drawn from a uni-
528 form distribution as well, to avoid fitness degeneracy (multiple equally fit variants).

529 This simple modification of the model enables a non-monotonic time development of the absolute
530 Hamming distance, as shown in Appendix 2 Figure 2B, while preserving the dynamics of relative
531 Hamming distance shown in Figure 3.

532 This is of course a highly simplistic variation upon the base model, but it serves to show that pro-
533 longed infection or introduction of (mutated versions of) previous variants can account for absolute
534 Hamming distance sometimes decreasing at variant transitions.

535 Acknowledgments

536 We would like to thank the members of the Grenfell and Levin Labs at The Department of Ecology
537 and Evolutionary Biology, Princeton University, for fertile plenary discussions. We would also like
538 to thank Arne Traulsen and Chadi M. Saad-Roy for enlightening discussions pertaining to the for-
539 mulation of our model and Christian Berrig and Viggo Andreasen at the PandemiX Center, Roskilde
540 University, for much appreciated comments on data visualization.

541 Research Funding

542 BFN and LS received funding from the Carlsberg Foundation under its Semper Ardens programme
543 (grant #CF20-0046, BFN and LS). BTG received funding from the Flu Lab. SAL acknowledges the sup-
544 port of the the C3.ai Digital Transformation Institute and Microsoft Corporation, Gift from Google
545 and the National Science Foundation (CNS-2027908, CCF1917819).

546 References

- 547 **Avanzato VA**, Matson MJ, Seifert SN, Pryce R, Williamson BN, Anzick SL, Barbian K, Judson SD, Fischer ER,
548 Martens C, et al. Case study: prolonged infectious SARS-CoV-2 shedding from an asymptomatic immuno-
549 compromised individual with cancer. *Cell*. 2020; 183(7):1901–1912.
- 550 **Baker RE**, Yang W, Vecchi GA, Metcalf CJE, Grenfell BT. Susceptible supply limits the role of climate in the early
551 SARS-CoV-2 pandemic. *Science*. 2020; 369(6501):315–319.
- 552 **Blanquart F**, Bataillon T. Epistasis and the structure of fitness landscapes: are experimental fitness landscapes
553 compatible with Fisher's geometric model? *Genetics*. 2016; 203(2):847–862.
- 554 **Burel E**, Colson P, Lagier JC, Levasseur A, Bedotto M, Lavrard-Meyer P, Fournier PE, La Scola B, Raoult D. Se-
555 quential appearance and isolation of a SARS-CoV-2 recombinant between two major SARS-CoV-2 variants in
556 a chronically infected immunocompromised patient. *Viruses*. 2022; 14(6):1266.
- 557 **Choi B**, Choudhary MC, Regan J, Sparks JA, Padera RF, Qiu X, Solomon IH, Kuo HH, Boucau J, Bowman K, et al.
558 Persistence and evolution of SARS-CoV-2 in an immunocompromised host. *New England Journal of Medicine*.
559 2020; 383(23):2291–2293.
- 560 **Corey L**, Beyrer C, Cohen MS, Michael NL, Bedford T, Rolland M. SARS-CoV-2 variants in patients with immuno-
561 suppression. *New England Journal of Medicine*. 2021; 385(6):562–566.
- 562 **Duerr R**, Dimartino D, Marier C, Zappile P, Wang G, Plitnick J, Griesemer S, Lasek-Nesselquist E, Dittman M,
563 Ortigoza MB, et al. Delta-Omicron recombinant SARS-CoV-2 in a transplant patient treated with Sotrovimab.
564 *bioRxiv*. 2022; .
- 565 **Ferguson NM**, Galvani AP, Bush RM. Ecological and immunological determinants of influenza evolution. *Nature*.
566 2003; 422(6930):428–433.
- 567 **Fisher R**. The genetical theory of natural selection. The Clarendon Press; 1930.
- 568 **Focosi D**, Maggi F. Recombination in Coronaviruses, with a Focus on SARS-CoV-2. *Viruses*. 2022; 14(6):1239.
- 569 **Gog JR**, Grenfell BT. Dynamics and selection of many-strain pathogens. *Proceedings of the National Academy*
570 *of Sciences*. 2002; 99(26):17209–17214.
- 571 **Grenfell BT**, Pybus OG, Gog JR, Wood JL, Daly JM, Mumford JA, Holmes EC. Unifying the epidemiological and
572 evolutionary dynamics of pathogens. *science*. 2004; 303(5656):327–332.
- 573 **Hadfield J**, Megill C, Bell SM, Huddleston J, Potter B, Callender C, Sagulenko P, Bedford T, Neher RA. Nextstrain:
574 real-time tracking of pathogen evolution. *Bioinformatics*. 2018; 34(23):4121–4123.

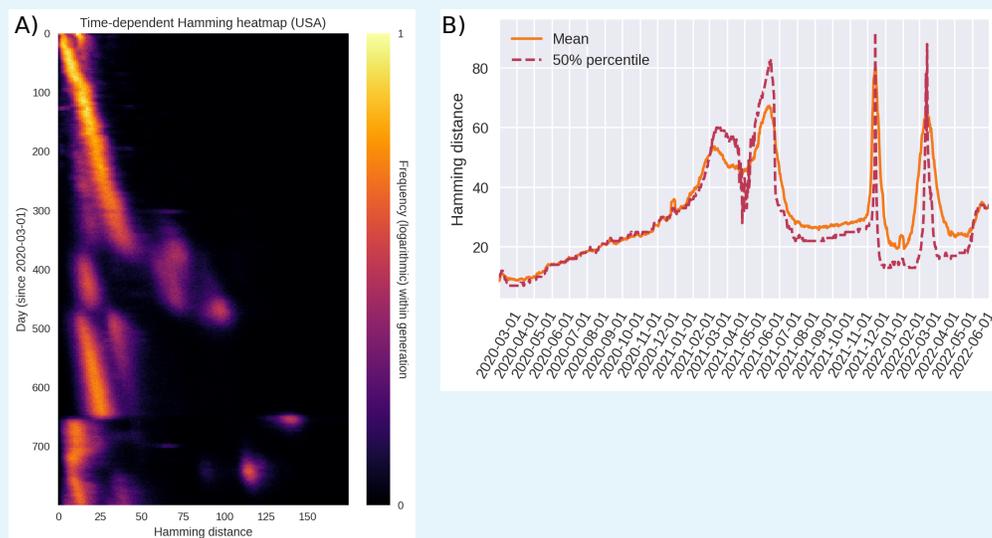
- 575 **Hammer AS**, Quaade ML, Rasmussen TB, Fonager J, Rasmussen M, Mundbjerg K, Lohse L, Strandbygaard B,
576 Jørgensen CS, Alfaro-Núñez A, et al. SARS-CoV-2 transmission between mink (Neovison vison) and humans,
577 Denmark. *Emerging infectious diseases*. 2021; 27(2):547.
- 578 **Harari S**, Tahor M, Rutsinsky N, Meijer S, Miller D, Henig O, Halutz O, Levytskyi K, Ben-Ami R, Adler A, Paran
579 Y, Stern A. Drivers of adaptive evolution during chronic SARS-CoV-2 infections. *Nature Medicine*. 2022 Jun;
580 p. 1–8. <https://www.nature.com/articles/s41591-022-01882-4>, doi: 10.1038/s41591-022-01882-4, publisher:
581 Nature Publishing Group.
- 582 **Hodcroft EB**, CoVariants: SARS-CoV-2 Mutations and Variants of Interest.; 2021. <https://covariants.org/>.
- 583 **Hou YJ**, Chiba S, Halfmann P, Ehre C, Kuroda M, Dinnon III KH, Leist SR, Schäfer A, Nakajima N, Takahashi K,
584 et al. SARS-CoV-2 D614G variant exhibits efficient replication ex vivo and transmission in vivo. *Science*. 2020;
585 370(6523):1464–1468.
- 586 **Isabel S**, Graña-Miraglia L, Gutierrez JM, Bundalovic-Torma C, Groves HE, Isabel MR, Eshaghi A, Patel SN, Gubbay
587 JB, Poutanen T, et al. Evolutionary and structural analyses of SARS-CoV-2 D614G spike protein mutation now
588 documented worldwide. *Scientific reports*. 2020; 10(1):1–9.
- 589 **Jackson B**, Boni MF, Bull MJ, Collieran A, Colquhoun RM, Darby AC, Haldenby S, Hill V, Lucaci A, McCrone JT,
590 et al. Generation and transmission of interlineage recombinants in the SARS-CoV-2 pandemic. *Cell*. 2021;
591 184(20):5179–5188.
- 592 **Katsnelson MI**, Wolf YI, Koonin EV. On the feasibility of saltational evolution. *Proceedings of the National
593 Academy of Sciences*. 2019; 116(42):21068–21075.
- 594 **Kauffman S**, Levin S. Towards a general theory of adaptive walks on rugged landscapes. *Journal of Theoretical
595 Biology*. 1987; 128(1):11–45.
- 596 **Kauffman SA**, Weinberger ED. The NK model of rugged fitness landscapes and its application to maturation
597 of the immune response. *Journal of theoretical biology*. 1989; 141(2):211–245.
- 598 **Kemp SA**, Collier DA, Datir RP, Ferreira IA, Gayed S, Jahun A, Hosmillo M, Rees-Spear C, Mlcochova P, Lumb IU,
599 et al. SARS-CoV-2 evolution during treatment of chronic infection. *Nature*. 2021; 592(7853):277–282.
- 600 **Koelle K**, Cobey S, Grenfell B, Pascual M. Epochal evolution shapes the phylodynamics of interpandemic in-
601 fluenza A (H3N2) in humans. *Science*. 2006; 314(5807):1898–1903.
- 602 **Kumata R**, Sasaki A. Antigenic escape accelerated by the presence of immunocompromised hosts. *bioRxiv*.
603 2022; <https://www.biorxiv.org/content/early/2022/06/14/2022.06.13.495792>, doi: 10.1101/2022.06.13.495792.
- 604 **Kupferschmidt K**. Where did 'weird' Omicron come from? *Science*. 2021; .
- 605 **Larsen HD**, Fonager J, Lomholt FK, Dalby T, Benedetti G, Kristensen B, Urth TR, Rasmussen M, Lassaunière R,
606 Rasmussen TB, et al. Preliminary report of an outbreak of SARS-CoV-2 in mink and mink farmers associated
607 with community spread, Denmark, June to November 2020. *Eurosurveillance*. 2021; 26(5):2100009.
- 608 **Martin G**. Fisher's geometrical model emerges as a property of complex integrated phenotypic networks.
609 *Genetics*. 2014; 197(1):237–255.
- 610 **Meijers M**, Ruchnewitz D, Łuksza M, Lässig M. Vaccination shapes evolutionary trajectories of SARS-CoV-2.
611 *arXiv*. 2022; <https://arxiv.org/abs/2207.09329>, doi: 10.48550/ARXIV.2207.09329.
- 612 **Meng B**, Abdullahi A, Ferreira IA, Goonawardane N, Saito A, Kimura I, Yamasoba D, Gerber PP, Fatihi S, Rathore
613 S, et al. Altered TMPRSS2 usage by SARS-CoV-2 Omicron impacts infectivity and fusogenicity. *Nature*. 2022;
614 603(7902):706–714.
- 615 **Newman ME**, Engelhardt R. Effects of selective neutrality on the evolution of molecular species. *Proceedings
616 of the Royal Society of London Series B: Biological Sciences*. 1998; 265(1403):1333–1338.
- 617 **Nussenblatt V**, Roder AE, Das S, de Wit E, Youn JH, Banakis S, Mushegian A, Mederos C, Wang W, Chung M, et al.
618 Yearlong COVID-19 Infection reveals within-host evolution of SARS-CoV-2 in a patient with B-cell depletion.
619 *The Journal of infectious diseases*. 2022; 225(7):1118–1123.
- 620 **Plotkin JB**, Dushoff J, Levin SA. Hemagglutinin sequence clusters and the antigenic evolution of influenza A
621 virus. *Proceedings of the National Academy of Sciences*. 2002; 99(9):6263–6268.

- 622 **Ritchie H**, Mathieu E, Rodés-Guirao L, Appel C, Giattino C, Ortiz-Ospina E, Hasell J, Macdon-
623 ald B, Beltekian D, Roser M. Coronavirus Pandemic (COVID-19). Our World in Data. 2022;
624 <https://ourworldindata.org/coronavirus>.
- 625 **Smith CA**, Ashby B. Antigenic evolution of SARS-CoV-2 in immunocompromised hosts. medRxiv. 2022; .
- 626 **Starr TN**, Greaney AJ, Hannon WW, Loes AN, Hauser K, Dillen JR, Ferri E, Farrell AG, Dadonaite B, McCallum
627 M, Matreyek KA, Corti D, Veessler D, Snell G, Bloom JD. Shifting mutational constraints in the SARS-CoV-2
628 receptor-binding domain during viral evolution. *Science*. 2022; 377(6604):420–424. [https://www.science.org/](https://www.science.org/doi/abs/10.1126/science.abo7896)
629 [doi/abs/10.1126/science.abo7896](https://www.science.org/doi/abs/10.1126/science.abo7896), doi: 10.1126/science.abo7896.
- 630 **Truong TT**, Ryutov A, Pandey U, Yee R, Goldberg L, Bhojwani D, Aguayo-Hiraldo P, Pinsky BA, Pekosz A, Shen L,
631 et al. Increased viral variants in children and young adults with impaired humoral immunity and persistent
632 SARS-CoV-2 infection: A consecutive case series. *EBioMedicine*. 2021; 67:103355.
- 633 **Varabyou A**, Pockrandt C, Salzberg SL, Perteau M. Rapid detection of inter-clade recombination in SARS-CoV-2
634 with Bolotie. *Genetics*. 2021; 218(3):iyab074.
- 635 **Volz E**, Hill V, McCrone JT, Price A, Jorgensen D, O'Toole Á, Southgate J, Johnson R, Jackson B, Nascimento FF,
636 et al. Evaluating the effects of SARS-CoV-2 spike mutation D614G on transmissibility and pathogenicity. *Cell*.
637 2021; 184(1):64–75.
- 638 **Weinreich DM**, Delaney NF, DePristo MA, Hartl DL. Darwinian evolution can follow only very few mutational
639 paths to fitter proteins. *science*. 2006; 312(5770):111–114.
- 640 **Weinreich DM**, Lan Y, Wylie CS, Heckendorn RB. Should evolutionary geneticists worry about higher-order
641 epistasis? *Current opinion in genetics & development*. 2013; 23(6):700–707.
- 642 **Zhang L**, Li Q, Liang Z, Li T, Liu S, Cui Q, Nie J, Wu Q, Qu X, Huang W, et al. The significant immune escape of
643 pseudotyped SARS-CoV-2 variant Omicron. *Emerging microbes & infections*. 2022; 11(1):1–5.

644 **Appendix 1**

645 **Diversity dynamics based on US sequences**

646 In the main text, we based our analysis on sequence data from the United Kingdom, since
647 the UK genomic surveillance was extensive. However, the overall results remain the same
648 if one instead considers United States sequence data. The exact timing of the variant tran-
649 sitions differ – most notably, the interval between the Alpha and Delta transitions is shorter
650 – but the qualitative features are similar. Variant transitions continue to be associated with
651 large jumps in Hamming distance, indicative of saltation. The temporal Hamming distri-
652 bution for the US is shown in Figure 1. Due to the shorter interval, the Alpha and Delta
653 transitions are not as clearly separated as in the UK data, and the Alpha transition appears
654 somewhat drawn out. As demonstrated in Figure 6–Figure supplement 1, this could be a
655 result of the presence of greater spatial effects.



656 **Appendix 1 Figure 1.** Hamming distance analysis analogous to that of Figure 1, except based on
657 United States sequence data, rather than British sequences. **A)** The full, time-dependent Hamming
658 distribution (US data, GenBank via Nextstrain). **B)** Time evolution of the mean and median pairwise
659 Hamming distances.
660

662 Appendix 2

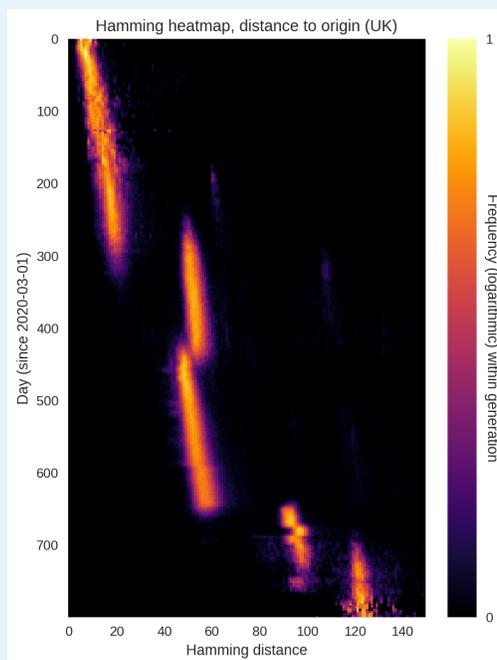
663 Fitting the origin-centered Hamming distribution

664 Our analysis is mainly centered around the dynamics of diversity, as measured by the dis-
665 similarity of SARS-CoV-2 genomes which are in circulation at a given point in time (Figure 1).
666 Our primary goal was to formulate a dynamical model which could qualitatively replicate
667 this pattern as parsimoniously as possible.

668 In this section, we additionally consider the distance between circulating genomes and the
669 *origin* (meaning the sequence Wuhan-Hu-1, GenBank reference sequence accession num-
670 ber MN908947.3). In Appendix 2 Figure 1, we show the distribution of distances between
671 circulating genomes and this reference sequence, which we will call the *absolute* Hamming
672 distance. The data analysis workflow in creating this plot is analogous to that detailed in
673 Figure 1–Figure supplement 1, but instead of continuously picking *pairs* of genomes, we
674 pick single genomes which are then compared with the reference genome. This also means
675 that, given N genomes within a given time window, there can be only N data points with
676 this method. This is in contrast to the pairwise comparison in Figure 1, where N genomes
677 give rise to $N(N - 1) \sim N^2$ possible pairings.

678 In Appendix 2 Figure 1, we see that the absolute Hamming distance at first increases ap-
679 proximately linearly with time until the Alpha transition, at which point a large jump in the
680 absolute Hamming distance is observed. Interestingly, the transition from Alpha to Delta
681 is associated with a slight decrease in this distance. The Omicron transition is again associ-
682 ated with a large increase in the absolute Hamming distance.

683



684

685

686

Appendix 2 Figure 1. Time evolution of the *absolute* Hamming distance, i.e. the distance to the origin (defined as the Wuhan-Hu-1 reference sequence). UK sequence data.

One possible explanation for such a decrease in absolute Hamming distance is that prolonged infections with an earlier variant have occurred in some individuals or a population, albeit without accompanying accelerated evolution. Once the unobserved lineage spills into

688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711

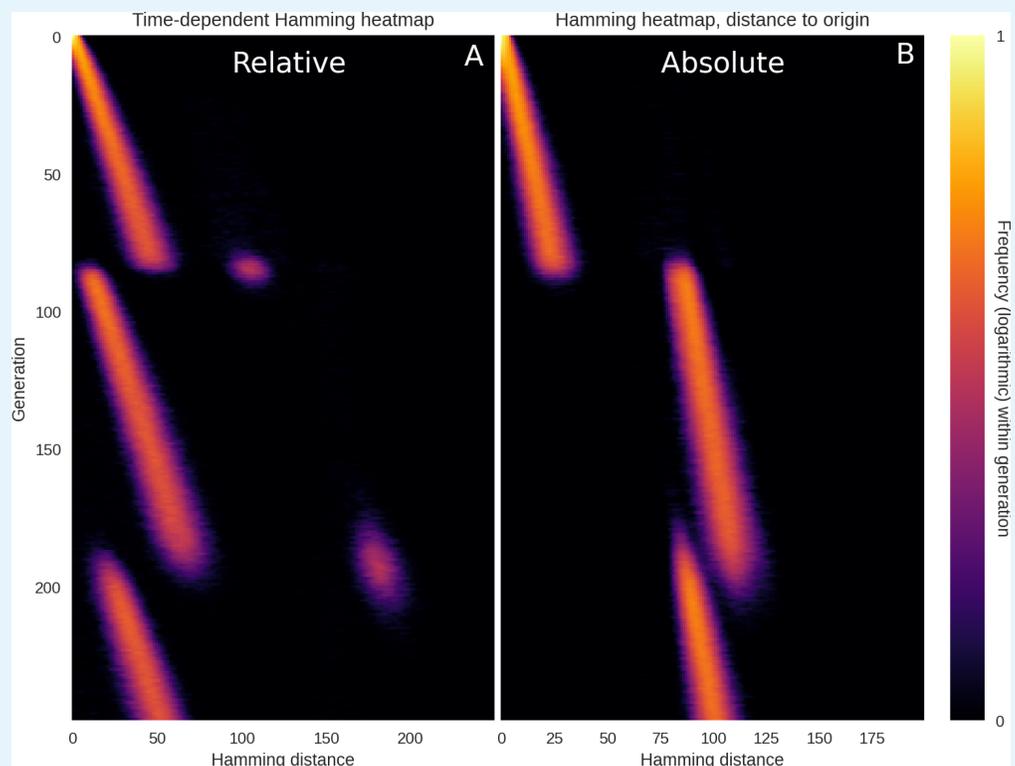
the sampled population, it may lead to a variant transition despite being closer to the ancestral variant than to the currently dominating one.

In terms of our model, this kind of dynamics can be quite simply incorporated. While it is not impossible to see decreases in absolute Hamming distance in the simple model formulation behind Figure 3, it is highly statistically unlikely, since each saltation happens on the basis of the genomes present in the previous generation.

The very simplest way to incorporate the possibility of decreasing absolute Hamming distance is to assume that a certain fraction of the population are initially infected with the ancestral strain, but that their infections are prolonged (and that they thus only transmit the disease much later). If the pathogen undergoes mutation within these hosts, it is possible to obtain a pattern such as the one shown in Appendix 2 Figure 2B. As evidenced by panel A of that same figure, the pairwise 'relative' Hamming distance between genomes in the same generation is not qualitatively affected by this addition to the model. Transitions are still driven by large saltations, and the (relative) Hamming distance is characterized by periods of linear growth punctuated by large increases and subsequent collapses in diversity.

The technical details of this addition to the model can be found in the Materials and Methods section.

While allowing for persistent infections with the initial variant is of course a very simple variation on the base model, it does show that prolonged infections with previous variants can account for occasional sudden decreases in the absolute Hamming distance.



712
713
714
715
716
717
718

Appendix 2 Figure 2. Simulation with persistent infections with original variant. When the original variant persists in a small fraction of the population, absolute Hamming distance can decrease at variant transitions. **A)** The 'relative' Hamming distance, the same quantity that was plotted in e.g. main text figures 1 and 3. It measures the dissimilarity between concurrently circulating pathogen genomes. **B)** The 'absolute' Hamming distribution, measuring the distance between circulating pathogen genomes and the reference sequence, namely that of the initial variant.

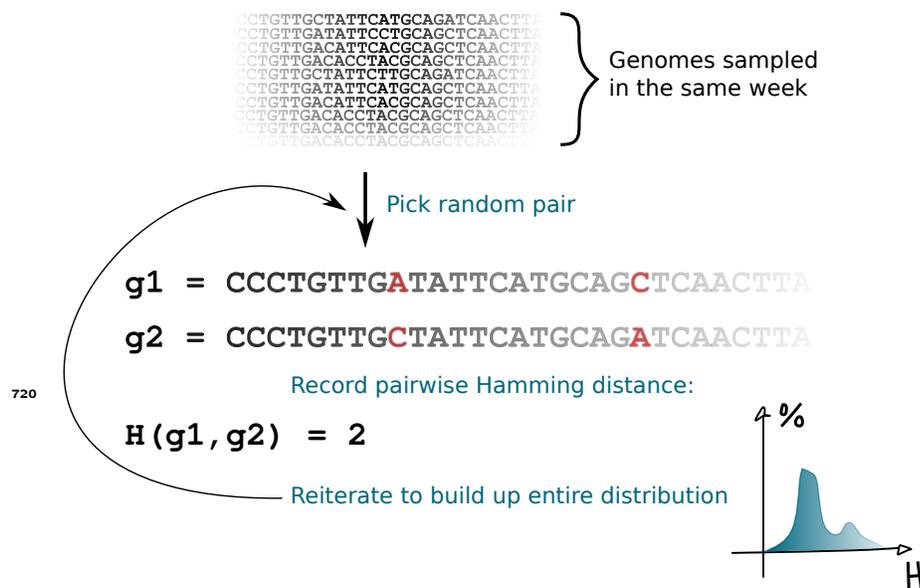


Figure 1–Figure supplement 1. Data analysis workflow. To generate the Hamming distribution for a given point in time, all sequences sampled within a week-long window starting on the given day are pooled. Then, pairs of sequences are repeatedly selected at random from this sequence pool, and the pairwise Hamming distance (number of sites which differ) is computed. All the computed Hamming distances are then pooled and a distribution (histogram) is generated.

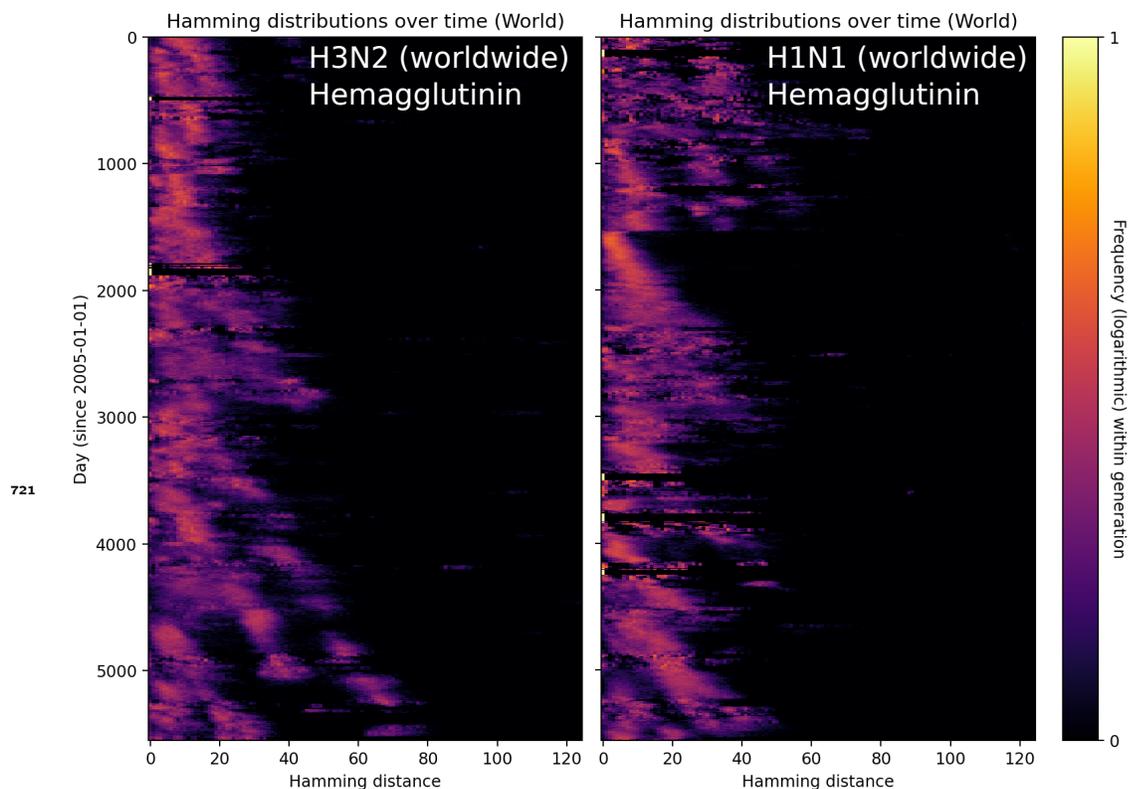
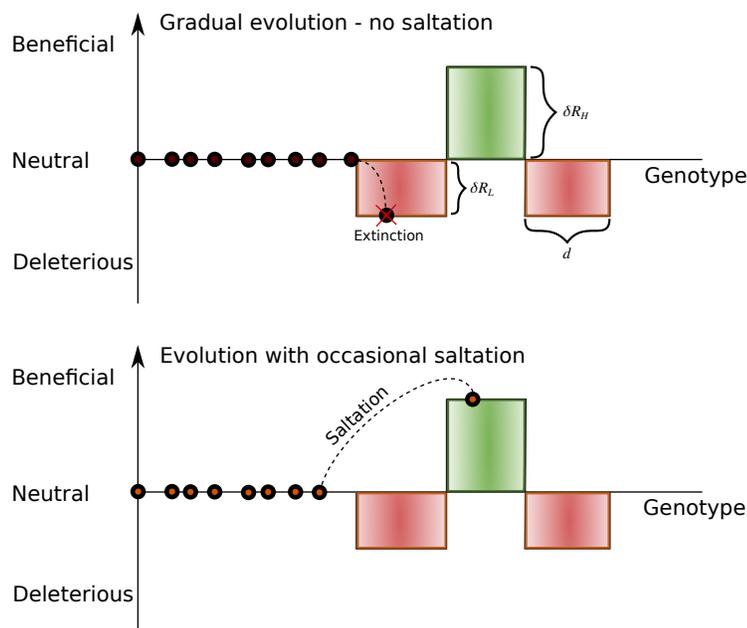


Figure 1-Figure supplement 2. Hamming distributions for influenza H3N2 and H1N1. With influenza, the amount of genomic surveillance data is much more limited and the temporal Hamming distributions are much less well-defined. In order to have enough data for each time point, a sampling window of 30 days was used here, as opposed to the 7 days used for SARS-CoV-2 in the main text.

A) Fitness landscape component



B) Branching model

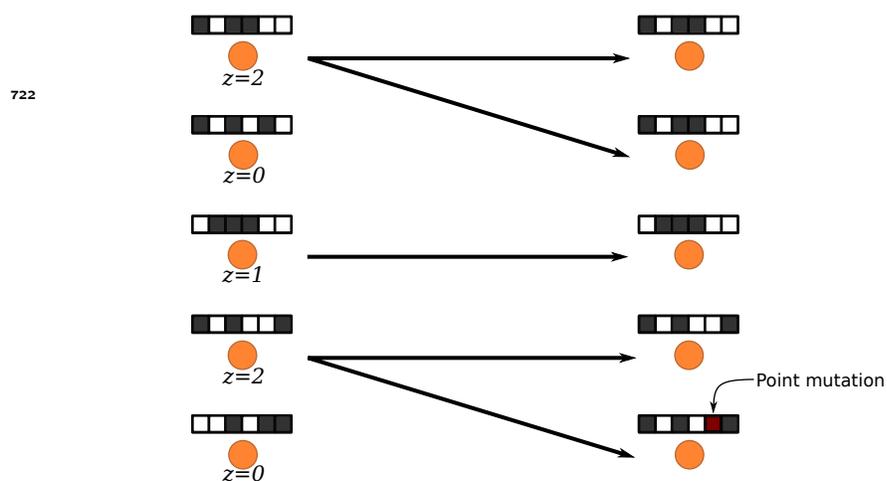


Figure 3-Figure supplement 1. Model schematics. A) The fitness landscape and epistasis components of the model. The majority of the fitness landscape is assumed neutral. In the case of gradual evolution devoid of saltation (top), the pathogen performs a random walk in this neutral space until it hits upon a deleterious configuration. As a model of sign epistasis, beneficial configurations are surrounded by deleterious ones. In the case of gradual evolution, the deleterious regions are unlikely to be traversed before the lineage dies out. However, in the case of saltational evolution (bottom), several point mutations may occasionally happen in the same genome within the same generation, leading to a jump which can enable the pathogen to bypass a deleterious region. Note that this is only a 1-dimensional conceptual representation of a highly multidimensional fitness landscape. **B)** In each generation of the branching model, each individual stochastically infects z new individuals. Upon infection, the pathogen genome (depicted as a string of black and white squares) is transferred. Occasionally a point mutation will occur, as indicated in the lower right genome. In the case of saltation (see panel A), multiple such point mutations can occur within the same genome in the same generation.

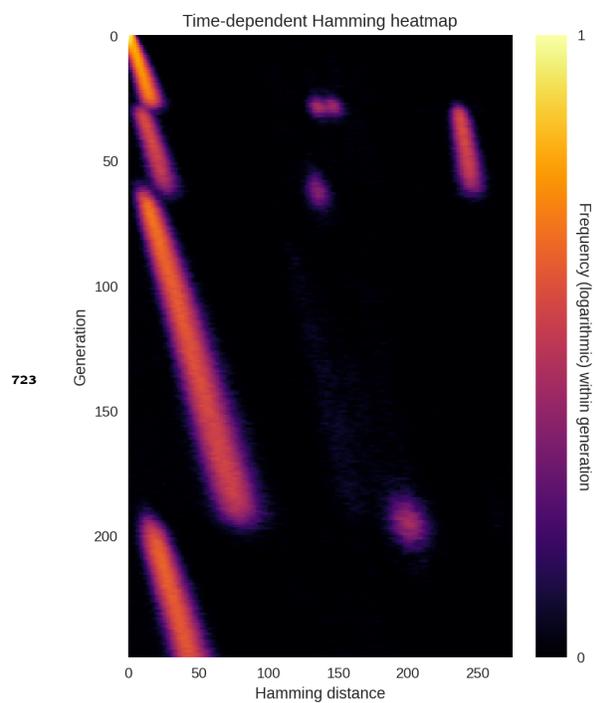


Figure 3–Figure supplement 2. Temporary coexistence of two equally fit variants.

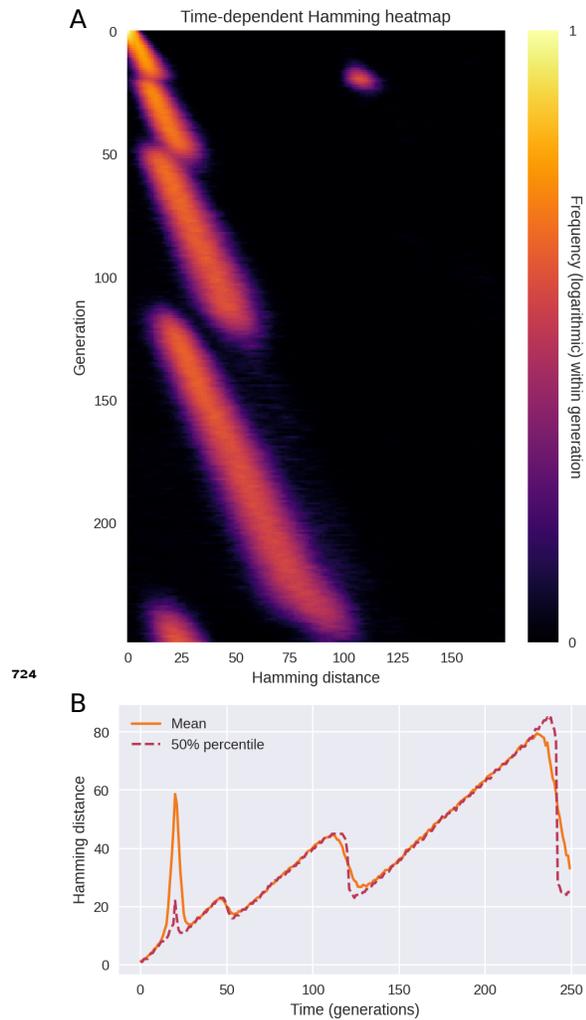


Figure 4-Figure supplement 1. Saltational evolution in the absence of sign epistasis. When saltational evolution is allowed, but epistasis is absent or very weak, a mixture of qualitatively different transitions occur. Some resemble the diversity spikes seen in Figure 3, but more commonly transitions will involve a gradual, linear increase in diversity followed by a collapse, as seen in Figure 4. **A)** Time evolution of the Hamming distance distribution. For each generation indicated on the vertical axis, the colour encodes the histogram of Hamming distances between genomes within that generation. **B)** Time evolution of the mean and median Hamming distance between genomes present in any given generation of the model simulation. In these simulations, $\delta R_L = 0$ (no epistasis) while saltations were of typical size $\mu_1 = 150$.

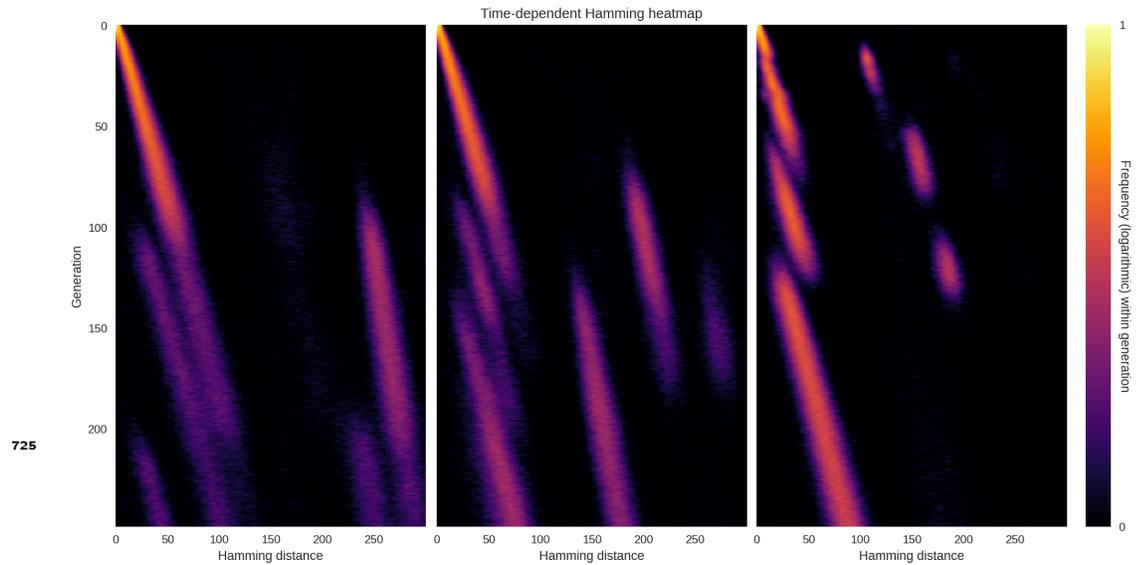


Figure 6-Figure supplement 1. Spatial structure leads to prolonged transition. Here we simulate the same SIRS dynamics as in Figure 6, but in a metapopulation consisting of three subpopulations. The within-population transmission rate $T_{ii} \approx 1$ ($i \in \{1, 2, 3\}$) is assumed much greater than the between-population transmission rate T_{ij} (with $j = i \pm 1$). **(A)** With inter-population transmission rate $\beta_{i,i\pm 1} = 0$, mutations never spread from one population to another and coexistence of variants with different fitness can last indefinitely. **(B)** With an inter-population transmission rate of 10^{-4} , transitions are severely prolonged but coexistence of variants with different fitness values does not last indefinitely. **(C)** At an inter-population transmission rate of 10^{-3} , transitions are only moderately prolonged compared to the non-spatial dynamics of Figure 6.