

1 **Title:** SARS-CoV-2 diversity and transmission on a university campus across two academic years
2 during the pandemic

3
4 **Authors:** Casto AM^{1,2*}, Paredes MI^{2,3}, Bennett JC^{1,3}, Luiten KG¹, O’Hanlon J¹, Han PD⁴, Gamboa
5 LS^{4,5}, McDermot E^{4,5}, Truong M^{4,5}, Gottlieb GS^{1,6,7}, Acker Z^{4,5}, Wolf CR¹, Magedson A¹, Lo NK¹,
6 McDonald D¹, Wright TC¹, McCaffrey KM⁴, Figgins MD², Englund JA⁸, Boeckh M², Lockwood
7 CM^{4,9}, Nickerson DA⁵, Shendure J^{4,5,10}, Uyeki TM¹¹, Starita LM^{4,5}, Bedford T^{2,3, 4,5,10}, Chu HY^{1,3#},
8 Weil AA^{1,7#}

9
10 **Affiliations**

11 ¹Division of Allergy and Infectious Diseases, Department of Medicine, University of Washington

12 ²Vaccine and Infectious Diseases Division, Fred Hutchinson Cancer Center

13 ³Department of Epidemiology, University of Washington

14 ⁴Brotman Baty Institute for Precision Medicine

15 ⁵Department of Genome Sciences, University of Washington

16 ⁶Environmental Health and Safety Department, University of Washington

17 ⁷Department of Global Health, University of Washington

18 ⁸Seattle Children’s Research Institute, Department of Pediatrics, University of Washington

19 ⁹Department of Laboratory Medicine and Pathology, University of Washington

20 ¹⁰Howard Hughes Medical Institute

21 ¹¹Centers for Disease Control

22 *Corresponding Author

23 #These authors contributed equally to this work.

24
25 **Disclaimer:** The findings and conclusions in this report are those of the authors and do not
26 necessarily represent the official position of the Centers for Disease Control and Prevention.

27
28 **Keywords:** SARS-CoV-2, university, college, genome, transmission, clade, lineage, cluster

29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44

45 **Abstract**

46 Institutions of higher education (IHEs) have been a focus of SARS-CoV-2 transmission
47 studies but there is limited information on how viral diversity and transmission at IHEs changed
48 as the pandemic progressed. Here we analyze 3606 viral genomes from unique COVID-19
49 episodes collected at a public university in Seattle, Washington (WA) from September 2020 to
50 September 2022. Across the study period, we found evidence of frequent viral transmission
51 among university affiliates with 60% (n=2153) of viral genomes from campus specimens
52 genetically identical to at least one other campus specimen. Moreover, viruses from students
53 were observed in transmission clusters at a higher frequency than in the overall dataset while
54 viruses from symptomatic infections were observed in transmission clusters at a lower
55 frequency. Though only a small percentage of community viruses were identified as possible
56 descendants of viruses isolated in university study specimens, phylodynamic modelling
57 suggested a high rate of transmission events from campus into the local community,
58 particularly during the 2021-2022 academic year. We conclude that viral transmission was
59 common within the university population throughout the study period but that not all
60 university affiliates were equally likely to be involved. In addition, the transmission rate from
61 campus into the surrounding community may have increased during the second year of the
62 study, possibly due to return to in-person instruction.

63

64

65

66

67 Introduction

68 The Covid-19 pandemic was marked by serial emergence of new SARS-CoV-2 variants
69 which outcompeted and replaced older variants in the general population¹. Geographical and
70 chronologic variation in circulating variants was revealed through viral genome sequencing
71 surveillance². These efforts were integral for identifying functional differences among variants
72 with regards to transmissibility, symptom profile and severity, and degree of neutralization by
73 antibodies resulting from previous infection or vaccination³. Genomic surveillance also helped
74 to predict fluctuations in incidence of COVID-19 and has aided public health authorities in
75 devising mitigation strategies and informing new vaccine formulations and recommendations⁴.

76 Educational institutions, including institutions of higher education (IHE), have been a
77 major focus of surveillance efforts during the pandemic. These efforts have been aimed at
78 improving understanding of transmission and informing infection control strategies. Though
79 IHEs vary widely in size, demographics, and setting, most IHE populations are predominantly
80 made up of young, healthy adults, that are at relatively low risk of severe disease from SARS-
81 CoV-2^{5,6}. Risk may be further mitigated by high rates of vaccine uptake in the setting of vaccine
82 mandates at some IHEs⁷. However, features of IHE environments, such as communal housing
83 and frequent social events, may promote the spread of respiratory viruses, including SARS-CoV-
84 2^{8,9}. IHE populations tend to be highly mobile, with students travelling to domestic or
85 international locations during academic breaks, which may contribute to the dispersal of new
86 variants⁹.

87 SARS-CoV-2 epidemiology and transmission studies on IHE campuses have estimated
88 varying viral transmission rates between IHE populations and their surrounding

89 communities^{8,10–15}. Some of this variation is likely due to fixed differences among IHEs, as well
90 as temporal variation in viral transmission dynamics across the pandemic. Most studies of
91 SARS-CoV-2 in IHE populations to date have analyzed data collected during relatively short time
92 frames (such as single academic quarters or semesters), preventing an accurate assessment of
93 how transmission changed over time and how new viral variants and changes in mitigation
94 strategies and human behavior shaped transmission as the pandemic progressed.

95 Here we examine SARS-CoV-2 genomic sequence data and demographic, epidemiologic,
96 and clinical data collected across from September 2020 to September 2022 as part of a
97 university testing program, the Husky Coronavirus Testing (HCT) study. We used three different
98 means to assess viral diversity and transmission. First, we characterized diversity of viruses seen
99 on campus and compared this to viral diversity within the state. Second, we identified clusters
100 of closely related HCT sequences and assessed the impact of demographic, epidemiologic, and
101 clinical factors of study participants on cluster membership. Third, we examined phylogenetic
102 relationships between university and community viruses.

103

104 **Methods and Materials**

105 *Study Overview*

106 Data analyzed in this study were collected as part of the Husky Coronavirus Testing
107 (HCT) study. HCT provided SARS-CoV-2 testing for affiliates (students, faculty, and staff) for the
108 University of Washington (UW) main campus and two satellite campuses from September 2020
109 to July 2023. Details of specimen and demographic, epidemiologic, and clinical data collection
110 are described in prior studies^{16,17}. Briefly, participation was open to all English-speaking

111 university affiliates. Participants completed an electronic questionnaire at enrollment and were
112 sent daily attestation surveys. Throughout the study, participants were invited to test if they
113 reported new symptoms, have a known exposure to a SARS-COV-2 case, or were members of a
114 group experiencing an outbreak. Walk-in testing was also available for any reason. In addition,
115 participants were invited to test from September 2020 to August 2021 following attendance at
116 gatherings with >10 people, from September 2021 to July 2022 following out of state travel,
117 and from August – September 2022 following report of a positive rapid test. The UW IRB
118 approved this study (#00011148). All participants gave informed consent or assent and
119 parent/guardian consent for participants under 18 years of age.

120 *Swab collection*

121 SARS-CoV-2 testing was performed via (1) participant self-swab collection observed by
122 study personnel at testing sites on campus, (2) unobserved self-swab collection returned to an
123 on-campus drop box, and (3) unobserved self-swab collection picked up by courier. Two swab
124 types were used for observed swab collection at kiosks; US Cotton #3 Steripack Polyester Spun
125 Swabs placed in a 10 mL tube were used at the beginning of the study with transition to
126 RHINOstic™ RH_S000001 Automated Nasal Swabs placed in a Matrix™ 1.0 mL Thermo Fisher
127 3741 ScrewTop tube in 2021. Swabs returned to drop boxes were RHINOstic™ while those
128 returned via courier were US Cotton #3 swabs.

129 *Specimen testing for SARS-CoV-2*

130 Specimens collected prior to November 18, 2020 were placed in universal transport
131 media and viral nucleic acid was extracted using two commercial kits as described previously¹⁸.
132 Specimens collected after November 18, 2020 were stored without preservatives or media.

133 Specimens were prepped for testing using an extraction-free protocol¹⁸. Aliquots of 5µL were
134 used in four multiplexed RT-qPCR reactions. Two reactions used custom probe sets that
135 targeted Orf1b and two used probes targeted the S-gene. Viral probe sets were multiplexed
136 with a probe set for human RNase P. Specimens were positive for SARS-CoV-2 if viral gene
137 targets and human RNase P were detected in at least three of four reactions.

138 *Genome Sequencing*

139 Genome sequencing was attempted on all specimens that tested positive for the
140 presence of SARS-CoV-2 with an average cycle threshold of 30 or less. Magna Pure 96 kits
141 (Roche) were used to extract nucleic acids from specimens and sequencing libraries were
142 prepared using COVIDSeq kits (Illumina). Sequencing primers were updated at several points
143 during the study period to account for emergence of new viral variants. Sequencing was
144 performed using NextSeq2000 P200 kits (Illumina). Processing of raw sequence data and
145 generation of consensus genomes was performed using a publicly available bioinformatic
146 pipeline (<https://github.com/seattleflu/assembly>). All genome sequences used in this study
147 were submitted to Global Initiative on Sharing All Influenza Data (GISAID)¹⁹.

148 *Genomic Analyses*

149 SARS-CoV-2 genomes used in analyses that were generated outside of the HCT study
150 were downloaded from GISAID. Both HCT and GISAID genomes were screened for quality using
151 Nextclade CLI²⁰. Nextclade was also used to assign sequences to Nextstrain clades and Pango
152 lineages. Sequences given an unfavorable quality rating by Nextclade (based on a sequence's
153 complement of missing data, mixed sites, private mutations, mutation clusters, frameshifts, and
154 premature stop codons), with a missing collection date, or for which Nextclade was unable to

155 make a clade and/or lineage assignment were excluded from further analysis. Sequence
156 alignment, masking of problematic loci, and phylogenetic tree generation were performed
157 using Nextstrain²¹. Trees were visualized using Auspice. Groups of identical SARS-CoV-2
158 sequences were identified using a previously described R package
159 (<https://github.com/blab/size-genetic-clusters>)²². Each group of identical sequences were
160 characterized by a set of mutations relative to the reference genome (Wuhan/Hu-1/2019,
161 GenBank Accession MN908947). Any genome that carried the same mutations as a group of
162 identical genomes plus additional mutations relative to the reference was categorized as a
163 descendant of that group. Phylogenetic groups were identified as follows: a phylogenetic tree
164 was constructed for each Nextstrain clade, which included all HCT and GISAID sequences from
165 Washington State (WA) belonging to that clade during the study period. All terminal nodes in
166 these trees were designated as HCT or non-HCT. Augur trait was used to assign HCT versus non-
167 HCT states for all internal nodes and to provide a likelihood of each state assignment²³. All HCT
168 terminal nodes that descended from the same internal node assigned a state of HCT with a
169 likelihood of 95% or greater were grouped into a single phylogenetic cluster.

170 *Transmission Modeling*

171 Our transmission modeling analysis included SARS-CoV-2 sequences divided into 3
172 regions of origin: HCT, consisting of HCT sequences; KC, consisting of sequences from King
173 County (KC), WA; and other, which consisted of contextual sequences from around the world to
174 account for outside viral introductions. We employed an equal temporal subsampling scheme
175 to enrich for under sampled time periods by randomly choosing a maximum of 400 total
176 sequences per region (HCT, KC, and other) sampled equally per each calendar month via Augur

177 filter²³, resulting in a set of 1137 total sequences, which were input into the model. Given the
178 differential number of specimens in each region-year-month combination, not all demes
179 included a total number of 400 sequences. We chose an equal temporal subsampling scheme
180 based on recent work showing that maximizing spatiotemporal diversity reduces bias in
181 MASCOT²⁴.

182 Using the compiled input sequence set, we employed a MASCOT-Skyline approach,
183 which approximates the structured coalescent, to predict when the most recent common
184 ancestor for each sequence pair in our input set existed and which of the three regions this
185 ancestor would have existed in. To generate these predictions, we made assumptions about
186 effective population sizes of the three regions and migration among the regions. To allow for
187 population sizes to change over time, we modeled effective population sizes similar to the
188 Skygrid approach for unstructured populations²⁵. We estimated the effective population size for
189 each location between time $t=0 \times \text{tree height}$, ..., $t=1 \times \text{tree height}$. Between each time point
190 where we estimated N_e , we assumed exponential growth. *A priori*, we assumed that the
191 effective population size at time $t+1$ is normally distributed with mean 0 and standard deviation
192 σ , with σ being estimated. We assumed the migration rate to be constant forward-in-time,
193 $m \frac{f}{zy}$, between states y and z . As the structured coalescent assumes backwards-in-time
194 migration rates, we assumed that backwards-in-time rate of migration between state y and z ,
195 $m \frac{b}{yz} = m \frac{f}{zy} \times \frac{N_e(t)_z}{N_e(t)_y}$. To infer effective population sizes and migration rates over time, we
196 employed an adaptable multivariate gaussian operator²⁶.

197 Parameter traces were visually evaluated for convergence using Tracer (v1.7.1)²⁷ and
198 30% burn-in was applied for all phylodynamic analyses. Output from our modeling analysis was

199 a phylogenetic tree with internal nodes representing common ancestors of input sequence
200 pairs. Tree plotting was performed with baltic (<https://github.com/evogytis/baltic>) and data
201 visualizations were done using Altair²⁸. We summarized trees as maximum clade credibility
202 trees using TreeAnnotator and visually inspected posterior tree distributions using IcyTree²⁹.
203 Transmission between regions was calculated by measuring the number of migration jumps
204 from HCT to KC and vice versa walking from tips to root in the posterior set of trees. Persistence
205 time was measured by calculating the average number of days for a tip to leave its sampled
206 location (HCT, KC, other), walking backwards up the phylogeny from tip up until node location
207 was different from tip location³⁰.

208

209 **Results**

210 *Viral lineages and clades common in Washington State were observed among HCT specimens*

211 We sequenced 3,855 of 6,485 SARS-CoV-2 positive specimens collected by HCT from
212 September 2020 to September 2022. These sequences represent 3% of all SARS-CoV-2 genomes
213 generated from specimens collected in Washington State (WA) during this time. From this raw
214 sequence set, we retained only one sequence per person per infection and filtered out poor
215 quality sequences, resulting in 3,606 sequences (Figure 1) in the final dataset; 3195 of these
216 sequences were collected during academic year 2 (September 1, 2021 – September 30, 2022;
217 hereafter referred to as year 2) with 1813 collected between December 1, 2021 and February
218 28, 2022 (Supplementary Figure S1). The final sequence set contained sequences from 19
219 different Nextstrain clades and 115 Pango lineages (Supplementary Note S1; Supplementary
220 Tables S1, S2).

221 To provide context for diversity seen among SARS-CoV-2 genomes from HCT specimens,
222 we downloaded all SARS-CoV-2 genomes from specimens collected in WA outside of the HCT
223 study from September 2020 to September 2022 from GISAID EpiCoV database¹⁹. After filtering
224 out poor quality and duplicate sequences, a total of 119,215 WA genomes remained,
225 representing 27 different clades and 333 lineages. All clades with a frequency of >0.2% and all
226 lineages with a frequency of >0.4% among WA genomes were represented by at least one HCT
227 genome. Most lineages in WA were rare (<0.4% of all WA genomes) and so more than half
228 (n=224, 67.3%) of all WA lineages were not observed among HCT genomes. There were 6
229 lineages that were represented in HCT but not the WA sequence set. These were all from
230 samples collected in early January or late March 2022 (Supplementary Table S3). The percent of
231 WA clades and lineages observed in HCT fluctuated over time; in year 2, these percentages
232 appeared to spike at the beginning of academic quarters (Supplementary Figure S2;
233 Supplementary Table S4).

234 *Average delay of one month between variant observation in WA and in HCT*

235 The prevalence of clades among HCT and WA sequences over time is shown in Figure 2.
236 For clades and lineages seen among both HCT and WA genomes, we determined the date of
237 first observation of a lineage and clade in each group. The average number of days from
238 observation in WA to observation among HCT specimens was 35.1 days (median 24, range 0 –
239 116) for clades and 35.5 days (median 28, range -56 to 170) for lineages (Supplementary Figure
240 S3). Ten lineages were observed among HCT specimens before WA specimens. Notably, the
241 BA.2 lineage, which was represented by 428 (11.9%) HCT sequences and 6,704 (5.6%) WA
242 sequences and from which all currently circulating SARS-CoV-2 are descended, was among

243 these and was first observed on campus on January 3, 2022. Three of the other lineages first
244 observed in HCT were also collected in January 2022.

245 *Most HCT SARS-CoV-2 specimens were closely related to at least one other HCT specimen*

246 We used two different approaches to identify groups of closely related HCT SARS-CoV-2
247 genomes. First, we identified groups of identical genomes (which we refer to as “zero distance
248 clusters”). There were 1730 unique haplotypes among HCT sequences, including 2153
249 sequences that were identical to at least one other HCT sequence and 277 haplotypes
250 represented by more than one HCT sequence (Figure 3A, Supplementary Table S5,
251 Supplementary Figure S4). A single Omicron haplotype (clade 21K, lineage BA.1.1) was observed
252 for 655 different sequenced specimens collected from December 17, 2021 until March 8, 2022
253 (18.2% of all HCT genomes). Of the 277 zero distance clusters, 26 included 10 or more
254 sequences. For each clade, the average size of zero distance clusters decreased with time since
255 clade introduction (Supplementary Figure S5)²² consistent with declining transmission rate over
256 time following variant introduction. The longest period over which a single haplotype was
257 observed was 153 days (clade 21L, lineage BA.2).

258 We also identified groups of identical sequences among a combined HCT and WA
259 dataset. Of 277 HCT zero distance clusters, 133 (48%) represented a haplotype not observed
260 among WA genomes (we refer to these as “HCT-only zero distance clusters”). The largest HCT-
261 only zero distance cluster (clade 21K, BA.1.20) included 12 sequences and the most persistent
262 cluster (longest period from collection of first to last specimen) was observed over a period of
263 35 days (clade 21K, BA.1.1). To assess for possible “spill-over” of virus from university affiliates
264 into other populations, we looked for WA viruses that appeared to be descendants of one of

265 133 HCT-only zero distance clusters (see Methods). We found a total of 81 such non-HCT
266 viruses, associated with 19 clusters (Supplementary Tables S6, S7). Over half (n=42, 51.9%) of
267 these 81 viruses were of the BA.2 lineage (clade 21L). The largest number of non-HCT
268 descendants of a single cluster was 37 (clade 21L, BA.2, Supplementary Figure S6).

269 We created a phylogenetic tree for each clade that included all HCT and WA genomes.
270 We used these trees to identify clusters (which we refer to as “phylogenetic clusters”) of HCT
271 genomes that descend from a single introduction event (see Methods). These clusters ranged in
272 size from 2 to 70 sequences with 19 clusters including more than 10 sequenced specimens
273 (Figure 3, Supplementary Table S8). Most (n=198, 84.6%) of the 234 HCT phylogenetic clusters
274 included only HCT sequences. However, a total of 218 WA sequences were part of an HCT
275 sequence cluster. The largest number of non-HCT sequences in a single HCT phylogenetic
276 cluster was 37 (clade 20I, lineage B.1.1.7; Supplementary Tables S9, S10).

277 *Model suggests high transmission rate from the university into the surrounding community*

278 To further explore the relationship between SARS-CoV-2 in university affiliates and the
279 surrounding community, we modeled transmission dynamics to and from the HCT population.
280 We limited the WA sequences in this analysis to those from King County (KC) to more
281 accurately reflect the community immediately surrounding the university. In addition to
282 including an HCT and KC region in the model, we also included an “other” region representing
283 sequences from outside KC in WA and the rest of the world. After subsampling (see methods), a
284 total of 1137 genomes were used as input for the model. Results suggested a higher forward
285 migration rate from HCT into KC than vice versa (Figure 4; 10.8 migration events/lineage/year
286 [95% highest posterior density (HPD) 4.3-19.9] vs 0.13 migration events/lineage/year [95% HPD

287 0.068-0.179]). We estimated that KC had at least 433 (IQR: 415-444) viral introduction events
288 during the study period with at least 130 events (IQR: 126 – 137) coming from the HCT
289 population. These numbers represent the lower bound of the number of introductions as the
290 absolute number is constrained by the total number of sequences in our specimen set. Our
291 model indicated that viral lineages are more likely to circulate longer in the larger KC region
292 (92.6 days, IQR: 86.4-101.1 days) than in the HCT population (77.2 days, IQR: 71.5 – 82.6 days).
293 When analyzing transmission patterns across time, we find that viral flow between HCT and KC
294 was dominated by spread from KC to HCT during academic year 1 (September 1, 2020 – August
295 31, 2021, hereafter referred to as year 1), and from HCT to KC during year 2 (Figure 5).

296 *Participants with a sequenced viral genome were representative of those testing positive for*
297 *SARS-CoV-2*

298 The 3,606 HCT genomes were from 3,560 unique individuals (Supplementary Table S11).
299 Most (85.4%) were students, 57.5% identified as female, 8.6% were Latinx, and the majority
300 were White (50.4%) or Asian (32.0%). Average age at the time of infection was 25.1 years
301 (median 21.3 years, range 17.4 – 78.7). HCT participants with a sequenced specimen were
302 overall demographically representative of all HCT participants with a positive test
303 (Supplementary Table S12). 3,514 individuals had only one sequence in the dataset while 46
304 individuals, who experienced infection with more than one clade/lineage of SARS-CoV-2 during
305 the study period, had two sequences included in the dataset (Supplementary Table S13,
306 Supplementary Note S2). Of these, 36 were first infected by a non-Omicron SARS-CoV-2 variant
307 followed by an Omicron variant, 9 were infected by two different Omicron variants, and one
308 was infected with two different non-Omicron variants (Supplementary Note S3).

309 *Sequences from students, younger persons more likely to cluster with other HCT sequences*

310 We examined the impact of epidemiologic, demographic, and clinical factors of infected
311 persons on the relationship among viruses in the HCT population. Students were
312 overrepresented among re-infected individuals relative to their frequency in the complete
313 dataset (one proportion z-test, 95.7% versus 85.4%, $p = 0.048$), while symptomatic infections
314 were underrepresented (47.8% versus 75.5%, $p < 0.0001$, Figure 6, Supplementary Table S14).
315 Additionally, average age at the time of infection was lower for those who experienced re-
316 infection than for all participants with a sequenced virus (21.1 versus 25.1, $p < 0.0001$). This was
317 also observed when sequences from students and those from faculty/staff were considered
318 separately (20.5 versus 21.7, $p < 0.0001$ and 33.0 versus 44.8, $p = 0.0096$).

319 Sequences from students were overrepresented among those in zero distance, HCT-only
320 zero distance, and phylogenetic clusters (90.9%, 97.0%, 97.6% versus 85.5%, $p < 0.0001$ for all)
321 while sequences from non-students (faculty/staff/other) were underrepresented in all three
322 cluster types (9.1%, 3.0%, 2.4% versus 14.5%, $p < 0.0001$ for all). Sorority/fraternity members
323 were also overrepresented in all three cluster types (21.5%, 23.5%, 31.3% versus 18.6%, $p =$
324 $0.00022, 0.0276, <0.0001$), while sequences from symptomatic infections were
325 underrepresented among those in HCT-only zero distance and phylogenetic clusters (64.4%,
326 65.9% versus 75.5%, $p < 0.0001$ for both). Finally, average age at the time of infection was
327 lower for sequences in all three cluster types compared to average age for all sequences (23.4,
328 21.3, 21.1 versus 25.1, $p < 0.0001$ for all). This difference was also observed when sequences
329 from students were considered separately (21.3, 20.6, 20.6 versus 21.7, $p < 0.0001$ for all).
330 Average ages at infection for sequences from non-students in all 3 cluster types did not differ

331 from the overall average age for non-students (44.0, 39.2, 43.2 versus 44.8, $p = 0.3922$, 0.6292,
332 and 0.2305). Results of similar analyses for other demographic and epidemiologic variables are
333 shown in Supplementary Figure S8.

334

335 **Discussion**

336 We studied SARS-CoV-2 cases with associated viral genomic data in a large, public
337 university population over the first two academic years of the pandemic with a focus on
338 characterizing viral diversity and transmission dynamics. To our knowledge, this represents the
339 largest survey to date of IHE SARS-CoV-2 cases with viral sequence data and one of the few
340 based on data collected for more than one year. Our results provide an in-depth analysis of
341 SARS-CoV-2 in an IHE population during many changes in viral epidemiology, transmission
342 mitigation strategies, and human behavior. We found that some measures of viral diversity and
343 transmission dynamics differed between year 1, during which online-only instruction occurred,
344 and year 2, during which classes were conducted mostly in-person, or appeared to be impacted
345 by the academic calendar; other measures were stable throughout the study period. We also
346 observed that not all university affiliates were equally likely to be involved in campus-related
347 viral transmission. Students and sorority/fraternity members were overrepresented in
348 transmission clusters while faculty/staff and those with symptomatic infections were under-
349 represented relative to their representation in the full dataset. In addition, the average age at
350 the time of infection of those in transmission clusters was lower than the overall average age
351 and the average age of students in clusters was lower than the average age for all students.
352 These findings provide context for and aid in interpretation of other studies of SARS-CoV-2 at

353 IHEs, particularly those with shorter study periods. They also can help administrators of IHEs
354 target mitigation strategies toward affiliates at highest risk of involvement in campus-related
355 SARS-CoV-2 transmission.

356 Our study has several unique features, including evaluation of viral diversity on campus
357 compared to the state over time. We found that clades/lineages common in the state were
358 reliably observed within HCT even during year 1, when the number of samples collected was
359 relatively small. There was a notable difference between year 1 and year 2 in the similarity
360 between clade frequency on campus and in the state, with HCT clade frequencies observed
361 during study year two almost identical to statewide frequencies. We also noted that the
362 average delay between first observation of a clade or lineage on campus relative to first
363 observation in the state was shorter during year 2 relative to year 1. These differences may be
364 explained by differences in sample size, though we hypothesize that changes in infection
365 control measures and in the virus, such as a return to in-person instruction, the introduction of
366 more transmissible variants, and overall higher infection rates in year 2, also played a role. We
367 observed evidence that the academic calendar influenced the relationship between viral
368 diversity on campus and in the state. We saw spikes in the percent of WA clades and/or
369 lineages represented on campus at the beginning of academic quarters in year 2. Lineages
370 unique to HCT or seen first in HCT were also mostly collected at the beginning of academic
371 quarters in year 2. These observations are consistent with increased importation of viral
372 diversity into the HCT population when students, 10% of whom are international and 15% of
373 whom are out-of-state residents, were returning to campus for the start of in-person

374 instruction from diverse geographic locations. These patterns were not observed in year 1 of
375 the study when classes were held exclusively online.

376 We used two different methodologies to define putative campus transmission clusters,
377 one based on genetic distance and one based on phylogenetic relationships. This was done to
378 mitigate disadvantages of each method and assess robustness of results. When comparing zero
379 distance and phylogenetic clusters, we noted similarities that are likely to be reflective of
380 campus transmission dynamics. Both methodologies suggested that campus related
381 transmission was common throughout the study period as most specimens were closely related
382 to at least one other HCT specimen. Most clusters were confined to a single academic quarter
383 as expected given drop off in campus population during breaks, though interestingly, there
384 were some exceptions to this. Two phylogenetic clusters persisted through spring break 2021;
385 one large Delta phylogenetic cluster started in early September 2021 and persisted into fall
386 quarter 2021, and several Omicron clusters started in late December 2021 and persisted into
387 winter quarter 2022, suggesting some on-going transmission among university affiliates even
388 during academic breaks. Additionally, both cluster types indicated that campus transmission
389 chains could persist for weeks to months despite frequent lineage replacement within the
390 SARS-CoV-2 population. The number of clusters per academic quarter was stable during year 1
391 and increased during year 2. This seems to be due to the increase in the number of specimens
392 collected in year 2 relative to year 1, rather than a change in transmission dynamics, as the
393 percent of specimens falling into zero distance, HCT-only zero distance, and phylogenetic
394 clusters remained stable across the study period. Finally, we observed that average cluster size
395 (for both zero distance and phylogenetic clusters) for a particular variant decreased with

396 increasing time since emergence of that variant. This is consistent with prior observations in
397 WA²² and is thought to be indicative of a decrease in the effective reproduction number for
398 viral variants the longer they have been circulating in a population.

399 We also conducted a modelling analysis to assess concordance with our other results.
400 The average number of days that viral lineages circulated on campus was estimated to be 77.2.
401 This was surprisingly high given that longest persistence times of the three cluster types were
402 153 days (zero distance), 35 days (HCT-only zero distance), and 60 days (phylogenetic), but does
403 provide further support for the hypothesis that viral transmission among university affiliates
404 was common during the study period. Interestingly, model results also suggested that viral flow
405 between KC and HCT was mostly into the campus population during year 1 and then from the
406 campus population during year 2. This could be the result of spread of new viral variants (Delta,
407 Omicron), the return to in-person instruction in year 2 with an increase in campus population
408 relative to year 1, or some combination of these. Given that we found few WA sequences that
409 appeared to descend from HCT clusters, we were surprised that the model estimated a
410 substantial number of transmission events from HCT to KC (estimated minimum of 130) and a
411 significantly higher forward migration rate than from KC to HCT. Most of this difference is likely
412 attributable to the vastly different population sizes of the two regions; KC has population of
413 2.252 million while HCT enrolled 37,360 participants. If we imagine a SARS-CoV-2 transmission
414 chain starting in KC, we expect a 1.7% chance of that transmission chain jumping into HCT due
415 to differences in population size alone. Conversely, if the transmission chain started in HCT, we
416 estimate a 98.3% chance of this chain infecting an individual in KC. This asymmetry corresponds
417 to a 59.2 fold larger viral migration rate from HCT to KC than vice versa. The fact that this

418 magnitude increase is similar but still less than the 85.7 fold difference in estimated forward
419 migration rates suggests that the difference in migration rates can largely be explained by the
420 difference in population sizes, further augmented by the presence of population structure. In
421 summary, the results of the modeling analysis do not suggest that SARS-CoV-2 cases in the HCT
422 population had a disproportionate impact on KC, but also that nearly all HCT transmission
423 chains resulted in an infection in the KC population.

424 We used clinical, epidemiologic, and demographic data for HCT participants to assess
425 the relationship of these factors to viral clustering. While findings that viruses from students
426 were disproportionately represented in clusters and that average age of those with viruses in
427 clusters was younger than the population average are not surprising, these results provide
428 evidence to support the direction of limited infection control resources at IHEs to those most
429 likely to be involved in transmission chains. Studies with more detailed information about
430 participant housing, activities, behaviors, and vaccination status could help to further delineate
431 drivers of this association between student status, age and, cluster membership to more
432 strategically target infection control resources. Due to outbreaks associated with
433 sorority/fraternity membership at the study university in 2020^{16,31}, the HCT study did collect
434 data on participant involvement with these social groups and our results indicated that
435 sorority/fraternity members were disproportionally represented in transmission clusters. Data
436 on membership in other social groups, such as sport teams or clubs, was not collected by HCT
437 and we are unable to comment on the impact of participation in these activities on involvement
438 in campus-related viral transmission. We note, though, that unlike most social groups, sorority
439 and fraternity members frequently live together in communal housing, which could be a major

440 driver of their risk of involvement in viral transmission chains. Finally, differences in clustering
441 observed for symptomatic versus asymptomatic cases was an unanticipated result. One
442 possible explanation is differences in behavior of the two groups, such as increased social
443 distancing and isolation by symptomatic individuals.

444 Limitations of our study included incomplete case identification and sampling on the
445 university campus and in WA state during the study period. Additionally, sequence data could
446 not be generated for all HCT cases, particularly those involving specimens with low amounts of
447 viral RNA (high cycle thresholds). University affiliates could also test outside of HCT and
448 sequenced specimens from affiliates collected outside HCT were classified as non-HCT WA
449 sequences. This reflects the broader challenge of the lack of associated demographic and
450 clinical data for most SARS-CoV-2 genomes in GISAID. This limits our understanding of
451 relationships between viral transmission and factors such as gender, age, race/ethnicity,
452 symptoms, and place of residence below the level of state. In addition, changing availability of
453 genomic surveillance data over time and unequal sampling across WA and the world impacted
454 the probability that a case was represented by a sequence in our dataset. We attempted to
455 mitigate this potential bias in our modeling analysis by using spatiotemporal subsampling,
456 which has been shown to improve inferential power of similar models²⁴. However, conclusions
457 of our modeling analysis, and all similar modeling analyses, are limited by the fact that results
458 are based on assumptions about population sizes and migration rates, which may be
459 inaccurate, and on the input sequence set, which represents a small fraction of all SARS-CoV-2
460 cases occurring in the three regions during the study period.

461 Populations of IHEs have been and will continue to be a focus of SARS-CoV-2 research
462 out of concern that these populations are prone to frequent transmission, which may have
463 significant impacts on IHEs and surrounding communities. Varying results have made it
464 challenging to derive generalizable lessons from studies conducted in IHEs during the last
465 several years. Here we have characterized viral diversity and transmission at a single IHE over
466 two years to gain an understanding of how viral diversity and transmission dynamics at a single
467 institution can vary over time and to aid in the synthesis of the data and results from previous
468 studies into a cohesive knowledge base. Such knowledge is vital to future optimization of
469 interventions to limit spread of SARS-CoV-2 and other respiratory viruses in IHE populations.

470

471 **Acknowledgements**

472 We would first like to thank all the study participants and the HCT study team. We would also
473 like to thank the University of Washington, including UW Environment Health and Safety team
474 (Katia Harb, Sheryl Schwartz, Natalie Thiel, Kim Baker, and Julie Skene) and the UW Covid
475 Incident Command team (Margaret Shepherd, Josh Gana, Pamela Schreiber, and Jack Martin).
476 Finally, we would like to thank all contributors of data to GISAID. We have included a GISAID
477 acknowledgements table in the Supplementary Material (Supplementary Table S15).

478

479 **Funding**

480 This work was supported by a Howard Hughes Medical Institute Covid Supplement Award to TB
481 and by the United States Senate and House of Representatives, Bill 748, Coronavirus Aid, Relief,

482 and Economic Security Act. MIP is an ARCS Foundation scholar. TB is a Howard Hughes Medical
483 Institute Investigator.

484

485 **Data Availability Statement**

486 All SARS-CoV-2 genomes used in this study have been deposited to GISAID (<https://gisaid.org/>).

487 Additional data and software files are available on github

488 (https://github.com/amcasto/huskytesting_SARSCoV2genomics_First2Years).

489

490 **Authorship**

491 Conceptualization: AMC, MIP, TB, HYC, AAW; Methodology: AMC, MIP, TB; Software: AMC,

492 MIP, KGL, TB; Validation: AMC, MIP; Formal Analysis: AMC, MIP; Investigation: JCB, KGL, JO,

493 PDH, LSG, EM, MT, ZA, CRW, AM, NKL, DM, TCW, KMM; Resources: GSG, JAE, MB, CML, DAN,

494 JS, LMS, TB, HYC, AAW; Data Curation: JCB, KGL, JO, PDH, LSG, EM, MT, ZA, CRW, AM, NKL, DM,

495 TCW, KMM; Writing – Original Draft Preparation: AMC, MIP, TB, HYC, AAW; Writing – Review &

496 Editing: All Authors; Visualization: AMC, MIP; Supervision: GSG, JAE, MB, CML, DAN, JS, TMU,

497 LMS, TB, HYC, AAW; Project Administration: JCB, GSG, ZA, CRW, NKL, JAE, MB, CML, DAN, JS,

498 LMS, TB, HYC, AAW; Funding Acquisition: JAE, MB, CML, DAN, JS, LMS, TB, HYC, AAW

499

500 **Competing Interests**

501 GSG has received research grants and/or research support from the US National Institutes of

502 Health, the University of Washington, the Bill & Melinda Gates Foundation, Gilead Sciences,

503 Alere Technologies, Merck & Co., Janssen Pharmaceutica, Cerus Corporation, ViiV Healthcare,

504 Bristol-Myers Squibb, Roche Molecular Systems, Abbott Molecular Diagnostics, and THERA
505 Technologies/TaiMed Biologics, Inc, all outside of the submitted work. JAE reports grants to her
506 institution from AstraZeneca, GlaxoSmithKline, Merck, and Pfizer and acts as a consultant for
507 Abbvie, Ark Biopharma, AstraZeneca, GlaxoSmithKline, Meissa Vaccines, Merck, Moderna,
508 Pfizer, and Sanofi Pasteur. MB has performed consulting for Allovir, Symbio, and Evrys Bio and
509 has received research support from Merck. CML reports that her spouse is an employee of
510 Bayer. HYC reports consulting for Ellume, Pfizer, and the Bill and Melinda Gates Foundation; has
511 served on advisory boards for Vir, Merck and Abbvie; has conducted CME teaching with
512 Medscape, Vindico, and Clinical Care Options; and has received research funding from Gates
513 Ventures, and support and reagents from Ellume and Cepheid outside of the submitted
514 work. All other authors have no competing interests to report.

515

516

517 **References**

518

519 1. Balloux, F. *et al.* The past, current and future epidemiological dynamic of SARS-CoV-2. *Oxf*

520 *Open Immunol* **3**, iqac003 (2022).

521 2. Tosta, S. *et al.* Global SARS-CoV-2 genomic surveillance: What we have learned (so far).

522 *Infect Genet Evol* **108**, 105405 (2023).

523 3. Markov, P. V. *et al.* The evolution of SARS-CoV-2. *Nat Rev Microbiol* **21**, 361–379 (2023).

524 4. Volz, E. Fitness, growth and transmissibility of SARS-CoV-2 genetic variants. *Nat Rev Genet*

525 **24**, 724–734 (2023).

- 526 5. Dudley, J. COVID-19 Transmission Under the Public Health Radar: High Prevalence in Young
527 Adults for COVID-19 Pandemic Wave 1. *International Journal of Infectious Diseases* **116**, S29
528 (2022).
- 529 6. Romero Starke, K. *et al.* The isolated effect of age on the risk of COVID-19 severe outcomes:
530 a systematic review with meta-analysis. *BMJ Glob Health* **6**, e006434 (2021).
- 531 7. Petros, B. A. *et al.* Early Introduction and Rise of the Omicron Severe Acute Respiratory
532 Syndrome Coronavirus 2 (SARS-CoV-2) Variant in Highly Vaccinated University Populations.
533 *Clin Infect Dis* **76**, e400–e408 (2023).
- 534 8. Valesano, A. L. *et al.* SARS-CoV-2 Genomic Surveillance Reveals Little Spread From a Large
535 University Campus to the Surrounding Community. *Open Forum Infect Dis* **8**, ofab518
536 (2021).
- 537 9. Nickbakhsh, S. *et al.* Genomic epidemiology of SARS-CoV-2 in a university outbreak setting
538 and implications for public health planning. *Sci Rep* **12**, 11735 (2022).
- 539 10. Richmond, C. S., Sabin, A. P., Jobe, D. A., Lovrich, S. D. & Kenny, P. A. *SARS-CoV-2*
540 *sequencing reveals rapid transmission from college student clusters resulting in morbidity*
541 *and deaths in vulnerable populations.*
542 <http://medrxiv.org/lookup/doi/10.1101/2020.10.12.20210294> (2020)
543 doi:10.1101/2020.10.12.20210294.
- 544 11. Leidner, A. J. *et al.* Opening of Large Institutions of Higher Education and County-Level
545 COVID-19 Incidence - United States, July 6-September 17, 2020. *MMWR Morb Mortal Wkly*
546 *Rep* **70**, 14–19 (2021).

- 547 12. Andersen, M. S., Bento, A. I., Basu, A., Marsicano, C. R. & Simon, K. I. College openings in
548 the United States increase mobility and COVID-19 incidence. *PLoS One* **17**, e0272820 (2022).
- 549 13. Srinivasa, V. R. *et al.* Genomic Epidemiology of Severe Acute Respiratory Syndrome
550 Coronavirus 2 Transmission Among University Students in Western Pennsylvania. *J Infect*
551 *Dis* **228**, 37–45 (2023).
- 552 14. Turcinovic, J. *et al.* Transmission Dynamics and Rare Clustered Transmission Within an
553 Urban University Population Before Widespread Vaccination. *J Infect Dis* **jjad397** (2023)
554 doi:10.1093/infdis/jiad397.
- 555 15. Turcinovic, J. *et al.* Linking contact tracing with genomic surveillance to deconvolute SARS-
556 CoV-2 transmission on a university campus. *iScience* **25**, 105337 (2022).
- 557 16. Weil, A. A. *et al.* SARS-CoV-2 Epidemiology on a Public University Campus in Washington
558 State. *Open Forum Infect Dis* **8**, ofab464 (2021).
- 559 17. Weil, A. A. *et al.* Genomic surveillance of SARS-CoV-2 Omicron variants on a university
560 campus. *Nat Commun* **13**, 5240 (2022).
- 561 18. Srivatsan, S. *et al.* Preliminary support for a “dry swab, extraction free” protocol for SARS-
562 CoV-2 testing via RT-qPCR. <http://biorxiv.org/lookup/doi/10.1101/2020.04.22.056283>
563 (2020) doi:10.1101/2020.04.22.056283.
- 564 19. Shu, Y. & McCauley, J. GISAID: Global initiative on sharing all influenza data - from vision to
565 reality. *Euro Surveill* **22**, 30494 (2017).
- 566 20. Aksamentov, I., Roemer, C., Hodcroft, E. & Neher, R. Nextclade: clade assignment, mutation
567 calling and quality control for viral genomes. *JOSS* **6**, 3773 (2021).

- 568 21. Hadfield, J. *et al.* Nextstrain: real-time tracking of pathogen evolution. *Bioinformatics* **34**,
569 4121–4123 (2018).
- 570 22. Tran-Kiem, C. & Bedford, T. *Estimating the reproduction number and transmission*
571 *heterogeneity from the size distribution of clusters of identical pathogen sequences.*
572 <http://medrxiv.org/lookup/doi/10.1101/2023.04.05.23287263> (2023)
573 doi:10.1101/2023.04.05.23287263.
- 574 23. Huddleston, J. *et al.* Augur: a bioinformatics toolkit for phylogenetic analyses of human
575 pathogens. *J Open Source Softw* **6**, 2906 (2021).
- 576 24. Layan, M. *et al.* Impact and mitigation of sampling bias to determine viral spread: Evaluating
577 discrete phylogeography through CTMC modeling and structured coalescent model
578 approximations. *Virus Evol* **9**, vead010 (2023).
- 579 25. Gill, M. S. *et al.* Improving Bayesian population dynamics inference: a coalescent-based
580 model for multiple loci. *Mol Biol Evol* **30**, 713–724 (2013).
- 581 26. Baele, G., Lemey, P., Rambaut, A. & Suchard, M. A. Adaptive MCMC in Bayesian
582 phylogenetics: an application to analyzing partitioned data in BEAST. *Bioinformatics* **33**,
583 1798–1805 (2017).
- 584 27. Rambaut, A., Drummond, A. J., Xie, D., Baele, G. & Suchard, M. A. Posterior Summarization
585 in Bayesian Phylogenetics Using Tracer 1.7. *Syst Biol* **67**, 901–904 (2018).
- 586 28. VanderPlas, J. *et al.* Altair: Interactive Statistical Visualizations for Python. *JOSS* **3**, 1057
587 (2018).
- 588 29. Vaughan, T. G. IcyTree: rapid browser-based visualization for phylogenetic trees and
589 networks. *Bioinformatics* **33**, 2392–2394 (2017).

590 30. Bedford, T., Cobey, S., Beerli, P. & Pascual, M. Global migration dynamics underlie evolution
591 and persistence of human influenza A (H3N2). *PLoS Pathog* **6**, e1000918 (2010).

592 31. UW News Staff. UW, Public Health - Seattle & King County responding to coronavirus cases
593 in Greek system. *UW News* (2020).

594

595

596

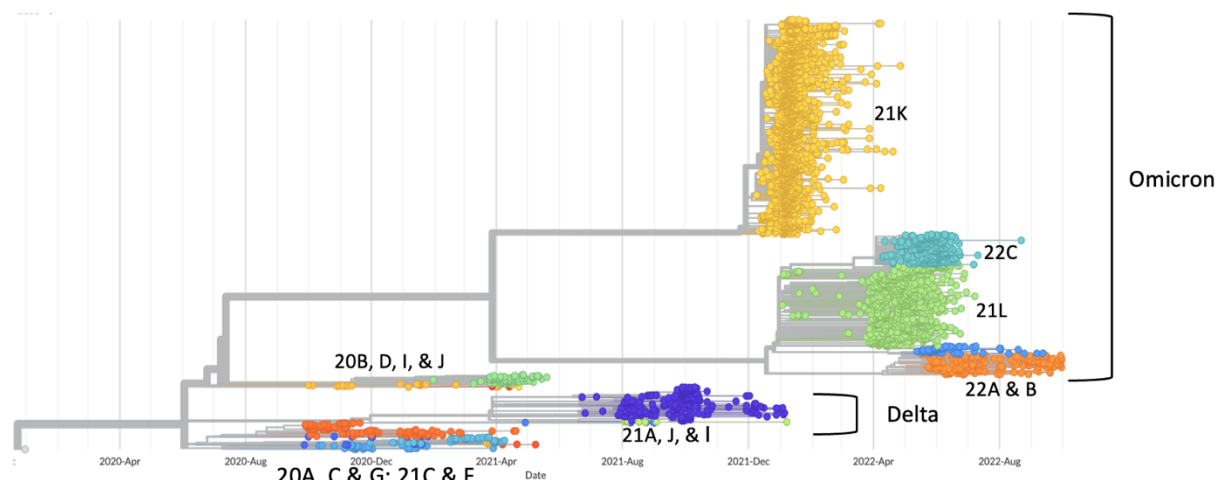
597

598

599

600 **Figures**

601



602

603 **Figure 1: Phylogenetic tree of all 3606 HCT sequences.** Tips are colored by Nextstrain clade.

604 Date of specimen collection is on the x-axis.

605

606

607

608

609

610

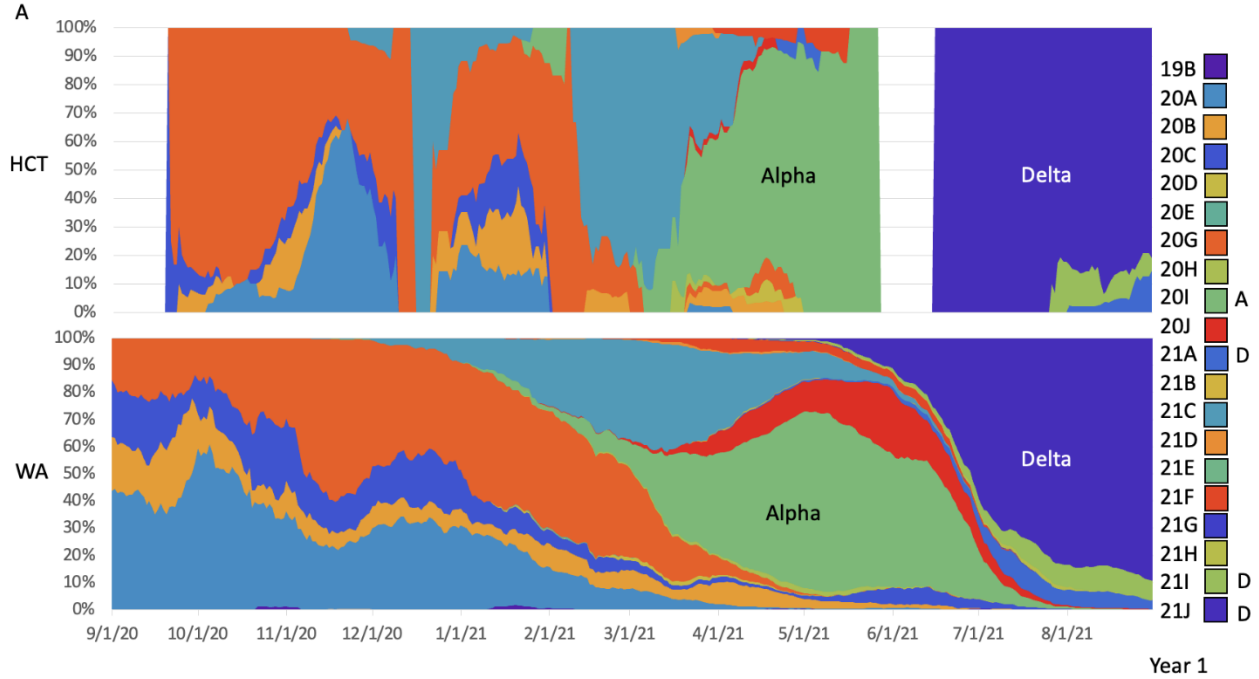
611

612

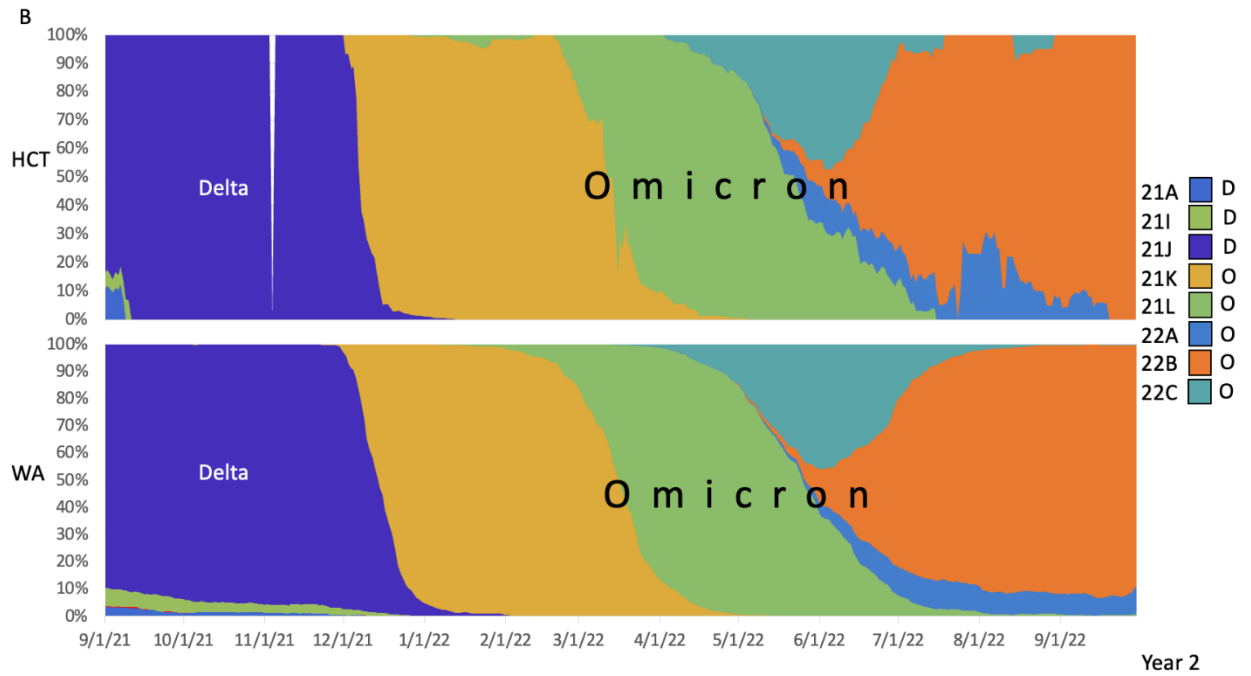
613

614

615
616
617
618
619



620



621
622
623
624

625 **Figure 2: Frequency of Nextstrain clades among HCT and WA SARS-CoV-2 genomes over time.**

626 Each time point on the x-axis represents a two-week sliding window centered on that date. Y-
627 axis shows the distribution of genomes collected in that window among different clades. Blank
628 sections for HCT represent two-week windows during which no sequenced specimens were
629 collected. Chart colors correspond to Nextstrain clades as shown in the legend. Alpha variant
630 clades are labeled “A”, Delta variant clades are labeled “D”, and Omicron variant clades are
631 labeled “O” in the legend. A) Specimens collected during year 1 (September 1, 2020 to August
632 31, 2021). B) Specimens collected during year 2 (September 1, 2021 to September 30, 2022).

633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658

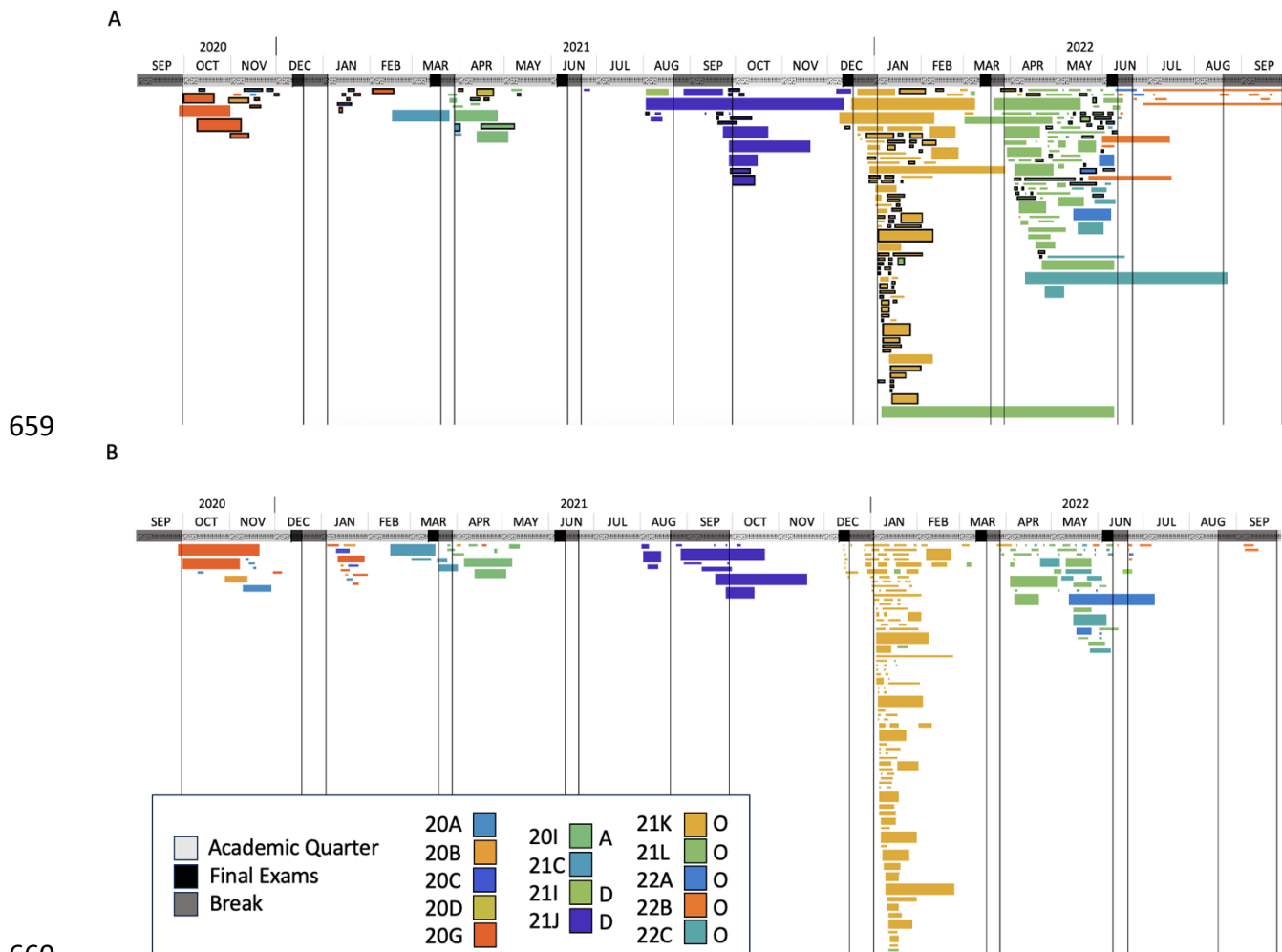
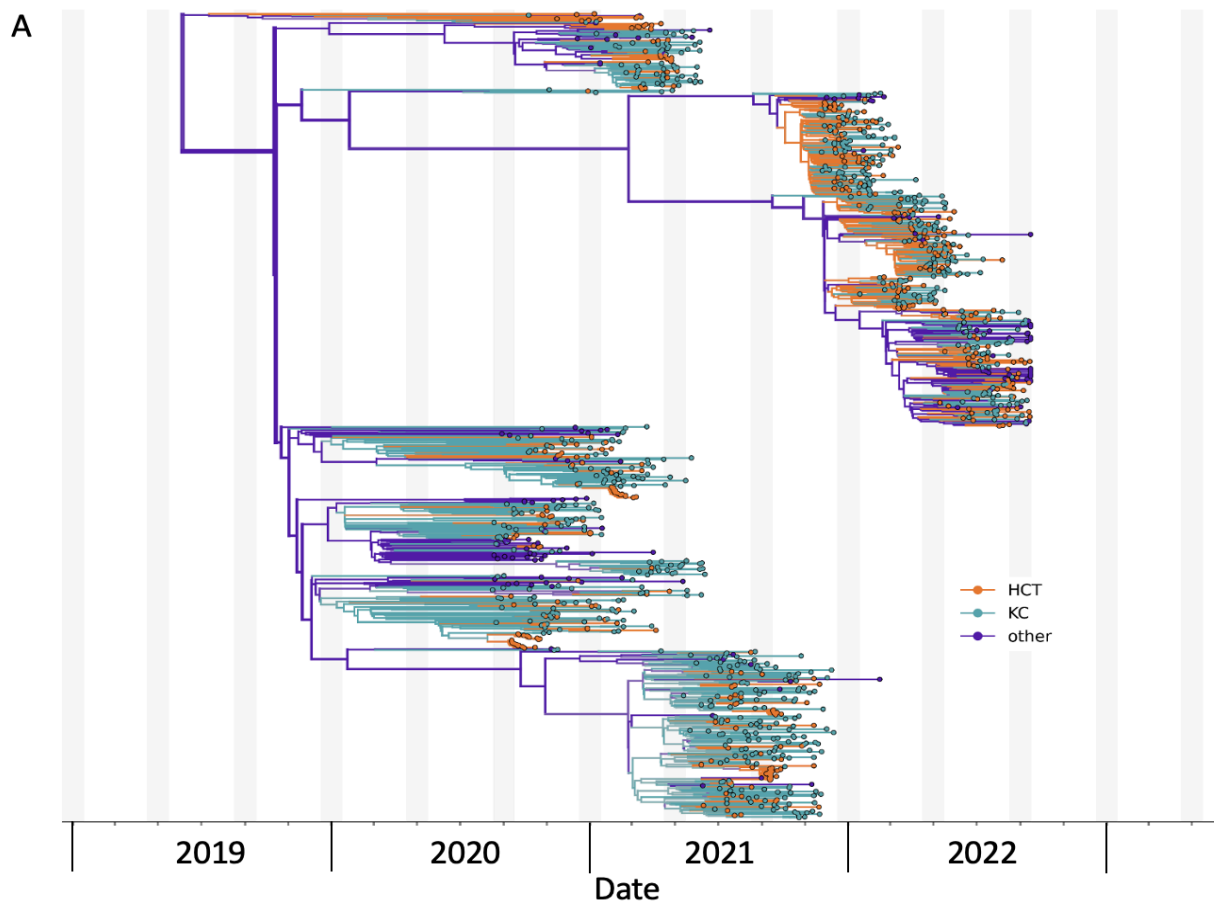
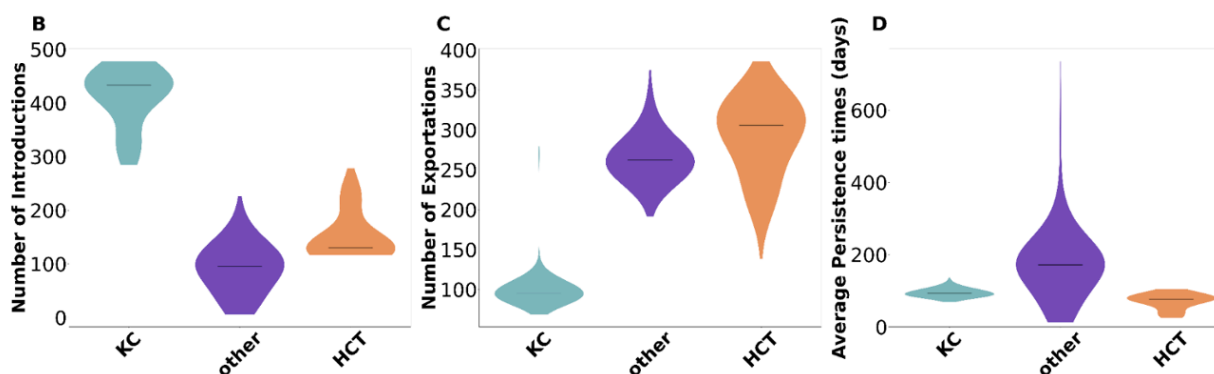


Figure 3: Chronology of clusters of HCT sequences. Each bar represents a sequence cluster which extends from the collection date of the first specimen in the cluster to the collection date of the last specimen. Month and year of specimen collection is denoted at the top with dates of academic quarters in light gray, final exams in black, and university breaks in dark gray. Color of cluster bars corresponds to viral clade as indicated in the legend. Alpha variant clades are labeled “A”, Delta variant clades are labeled “D”, and Omicron variant clades are labeled “O” in the legend. Height of cluster bar is proportional to the number of sequences in the cluster. A) Zero distance clusters (groups of identical sequences). Bars with black outlines represent HCT-only zero distance clusters (groups of identical sequences with haplotype unique to HCT). B) Phylogenetic clusters (groups of sequences that cluster phylogenetically).



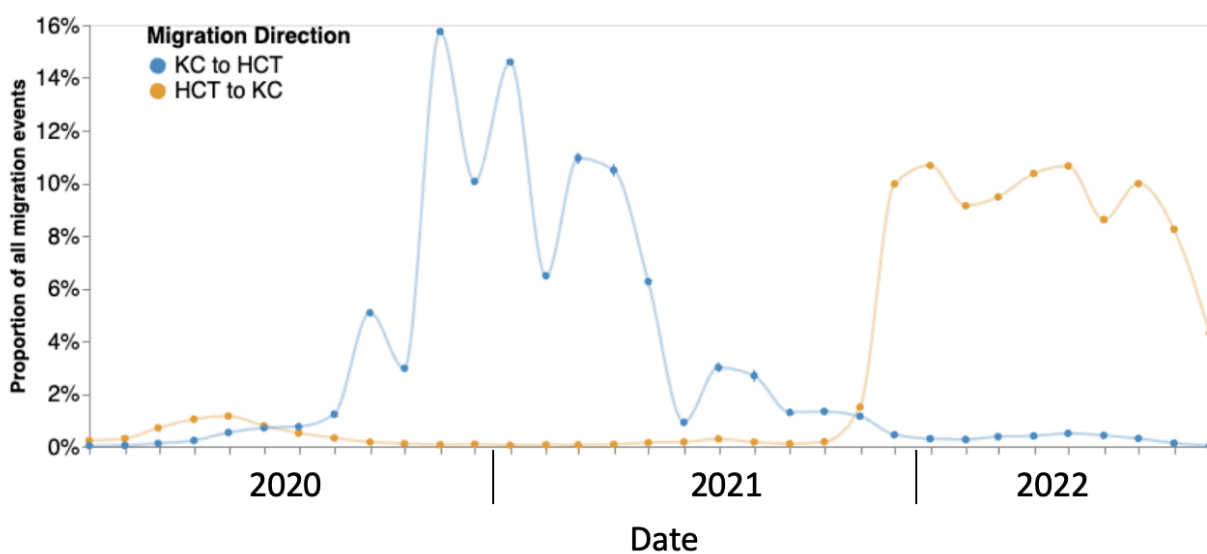
675
676



677
678
679
680
681
682
683
684
685

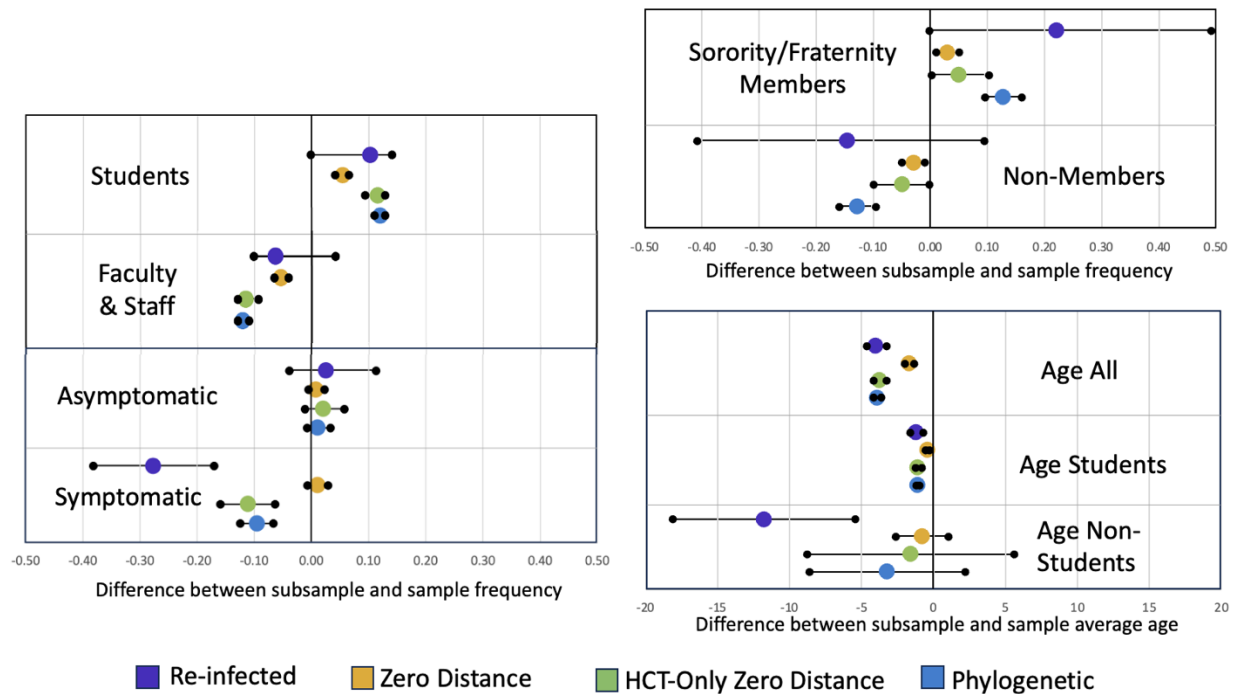
Figure 4: Phylodynamic analysis of SARS-CoV-2 transmission between KC and HCT

populations. A) Maximum clade credibility tree summary of the Bayesian inference conducted using MASCOT-Skyline on 1137 sequences. Colors correspond to the locations in the legend. KC = Sequences from King County (excluding those from HCT), HCT = Sequences from specimens collected by HCT, other = global contextual sequences from outside of King County sampled to increase spatiotemporal diversity. Estimated number of introductions B), exports C), and average time of local persistence in days D) for each region. Horizontal black line denotes median estimates.



686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706

Figure 5: Proportion of all migration events between HCT and KC. KC to HCT migration is in blue and HCT to KC migration is in orange. Proportions do not add up to 100% as migration events including the “other” region were excluded.



707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732

Figure 6: Demographic and epidemiologic characteristics of participants who experienced re-infection and of those with sequenced specimens in clusters. For student, faculty & staff, asymptomatic, symptomatic, sorority/fraternity members, and non-members categories, the mid-line ($x = 0$) represents the percent representation of that category in the whole dataset. The purple, yellow, green, and blue dots give the percent difference between the frequency of that category among those who experienced re-infection and those with sequences in zero distance (groups of identical sequences), HCT-only zero distance (groups of identical sequences with haplotype unique to HCT), and phylogenetic clusters (groups of sequences that cluster phylogenetically), respectively, and its frequency in the whole dataset. Intervals marked by black dots connected by a bar indicate the 95% confidence interval. For age all, age students, and age non-students categories, the mid-line ($x = 0$) represents the average age at time of infection for all sequences, sequences from students, and sequences from non-students, respectively. The purple, yellow, green, and blue dots give the difference between the above averages and the average age of those who experienced re-infection and for sequences in zero distance, HCT-only zero distance, and phylogenetic clusters, respectively. Intervals marked by black dots connected by a bar indicate the 95% confidence interval.

733 **Supplements**

734

735 **Supplementary Note S1**

736

737 Nextstrain clades and Pango lineages are two of the most commonly used methods for
738 classifying SARS-CoV-2. The Pango lineage system names lineages using a series of letters
739 separated from a series of numbers by a period (for instance, A.1 and B.2.75). The system is
740 hierarchical such that B.2.75 is a sublineage of the B.2 lineage. If a sublineage becomes
741 common enough, it is assigned its own alphabetic name. For instance, B.1.1.529 is equivalent to
742 BA. The Nextstrain clade system assigns names of two numbers corresponding to the year
743 when the viral group was first observed and then a letter (for instance, 19A and 22C). In
744 general, Nextstrain clades are usually larger (more inclusive) than Pango lineages, such that a
745 clade will be composed of multiple lineages.

746

747 **Supplementary Note S2**

748

749 For purposes of this study, we deemed an HCT participant to have been re-infected or
750 superinfected with SARS-CoV-2 during the study period if they had at least two sequenced
751 SARS-CoV-2 specimens that were determined to be of different clades and/or lineages. We did
752 not observe any participants with 3 or more sequenced specimens each of different clades
753 and/or lineages. It is possible that we excluded from the re-infected/superinfected participant
754 group participants that experienced re-infection/superinfection with two closely related viruses
755 (ie of the same clade and lineage). However, given the difficulty in distinguishing these
756 situations from on-going shedding of the same virus, we chose to limit our list of re-
757 infected/superinfected participants to those that had viruses of different clades and/or
758 lineages.

759

760 **Supplementary Note S3**

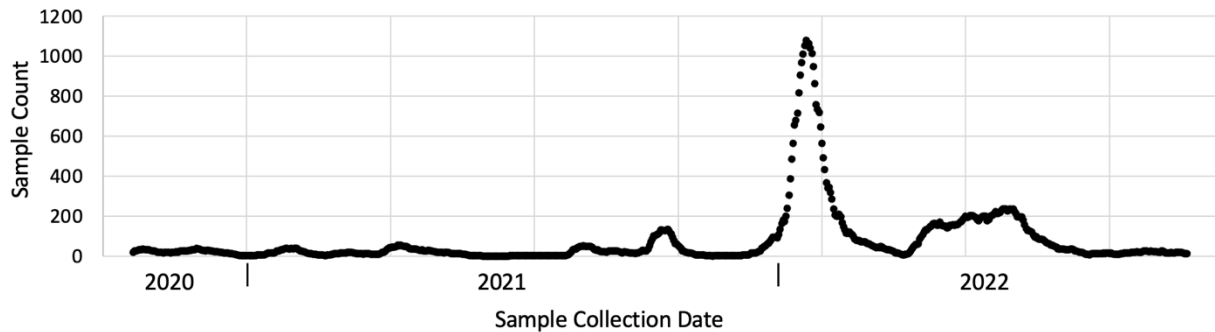
761

762 One of the 46 individuals identified as re-infected or superinfected had two sequenced
763 specimens of two different clades/lineages collected only 6 days apart (first specimen is
764 22B/BA.5 and second is 21L/BA.2.3). The sequencing reads for both specimens were reviewed
765 and there was no evidence in either specimen of mixed infection. These two specimens were
766 the participant's only SARS-CoV-2 positive specimens collected by HCT. Prior to testing positive
767 the first time, the participant reported a known SARS-CoV-2 exposure but no recent travel. The
768 participant denied having any history of previous SARS-CoV-2 infection prior to this first positive
769 test. The participant reported experiencing symptoms for several days. The possibility of a
770 specimen swap cannot be completely excluded, though there are numerous measures built into
771 the study protocol to prevent such errors and a review of the documentation surrounding the
772 collection of these specimens found no anomalies. Individuals testing using someone else's
773 information also cannot be excluded as a possibility, though is less likely in this case as both
774 specimens were collected via staff observed in-person testing.

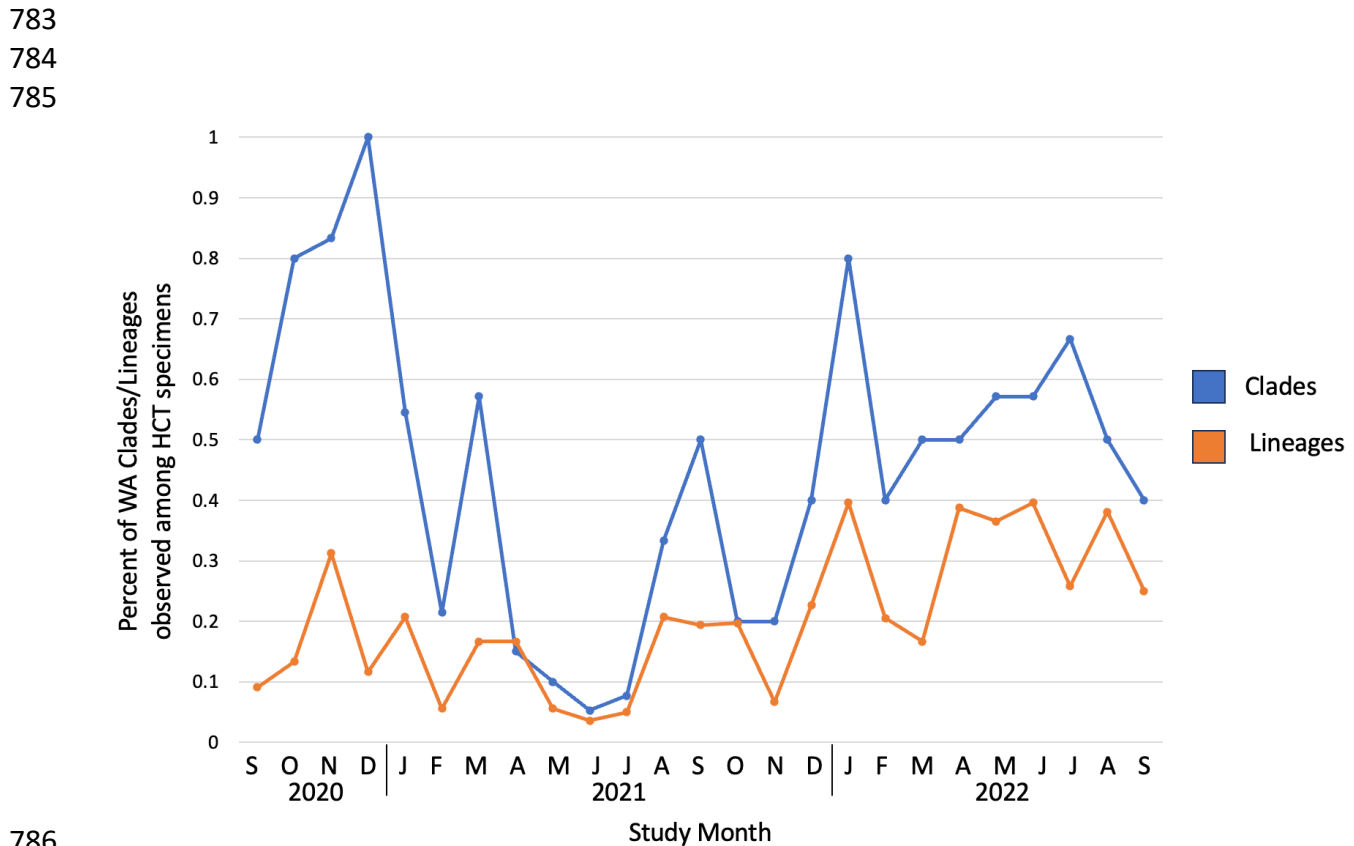
775

776

777
778 **Supplementary Figures**
779



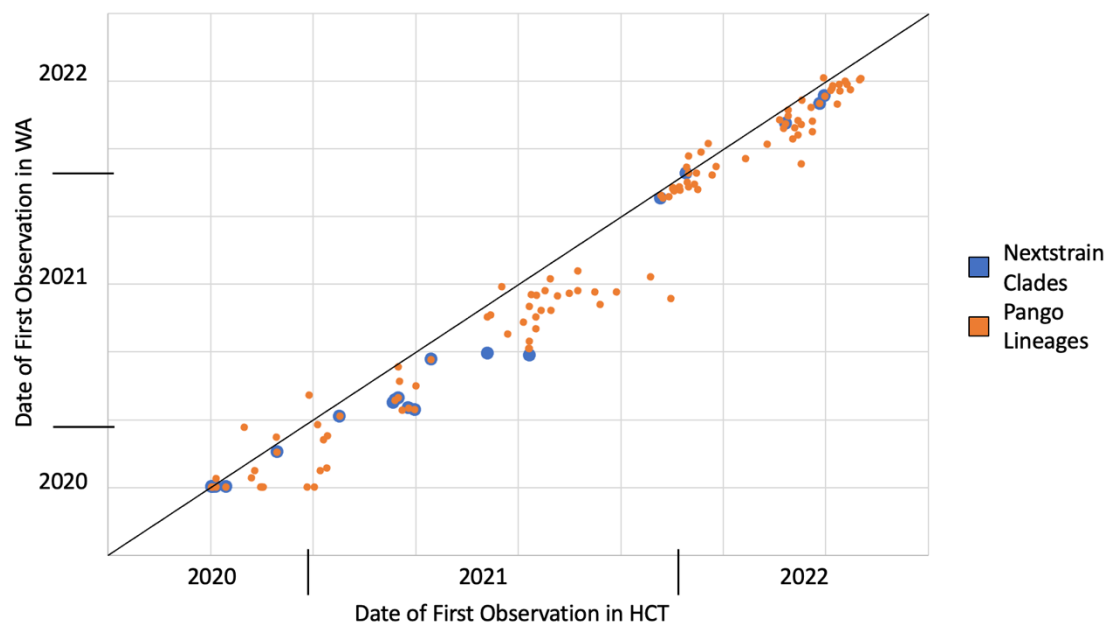
780
781 **Supplementary Figure S1: Total number of sequenced HCT specimens collected across the**
782 **study period by two week sliding window.**



786
787 **Supplementary Figure S2: Percent of WA clades/lineages observed among sequenced HCT**
788 **specimens by study month.**

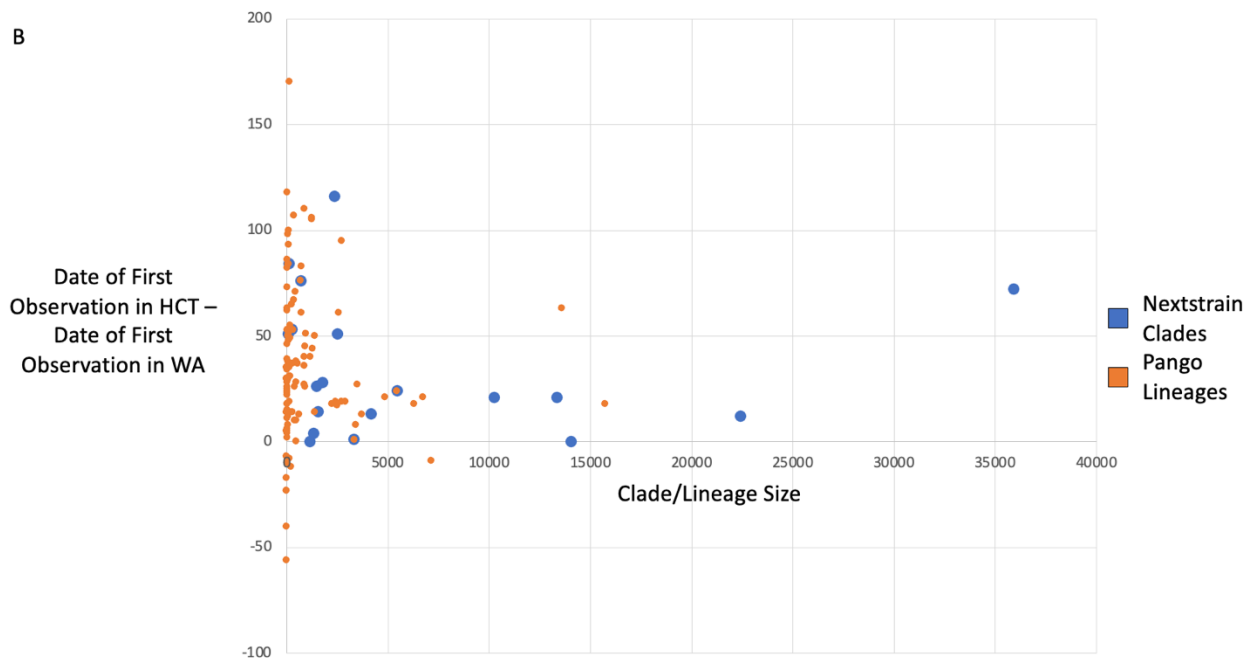
789
790
791

A



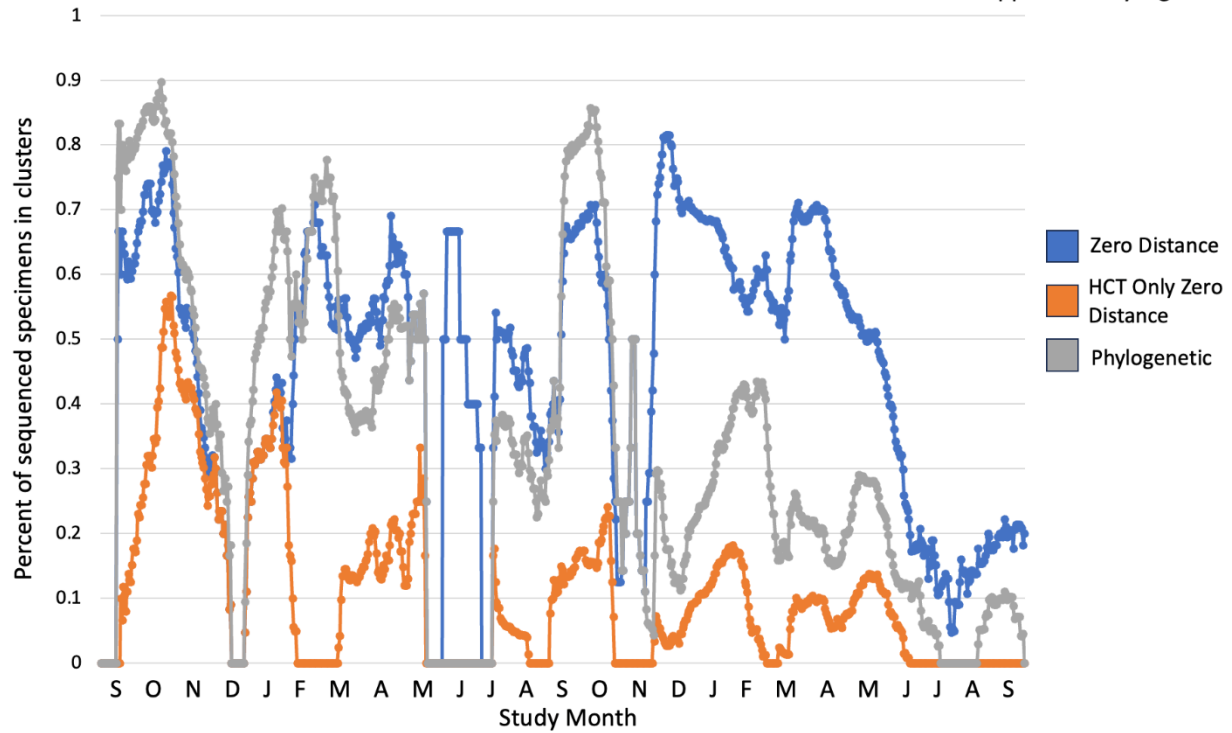
792
793

B



794
795
796
797
798
799
800
801

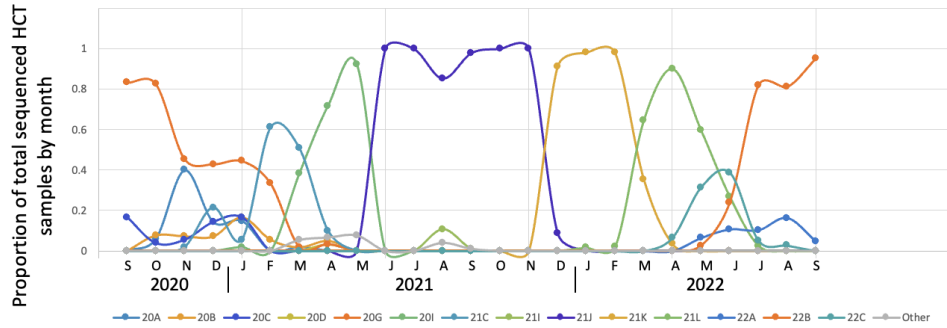
Supplementary Figure S3: Date of first observation of clades and lineages in WA and in HCT. Blue dots represent Nextstrain clades. Orange dots represent Pango lineages. A) Chart with date of first observation in HCT on the x-axis and first observation in WA on the y-axis. B) Chart with size of clade/lineages on the x-axis and date of first observation in HCT minus first observation in WA (in days) on the y-axis.



802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826

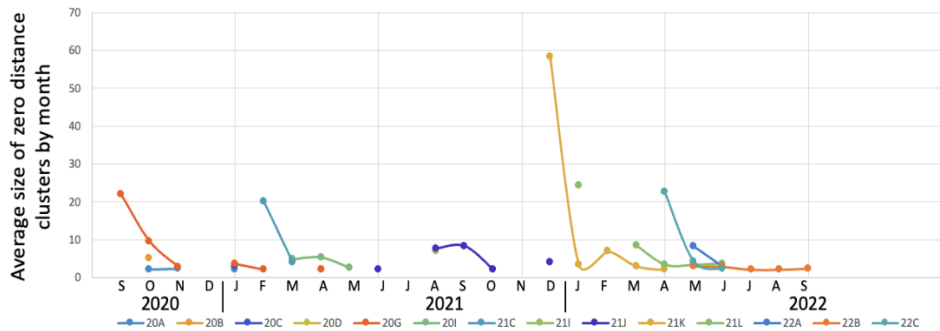
Supplementary Figure S4: Percent of sequenced specimens in zero distance, HCT-only zero distance, and phylogenetic clusters calculated in one month sliding window periods.

A



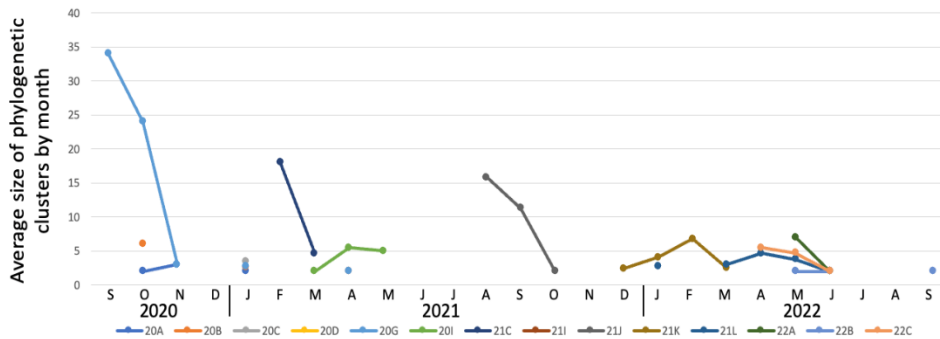
827
828

B



829

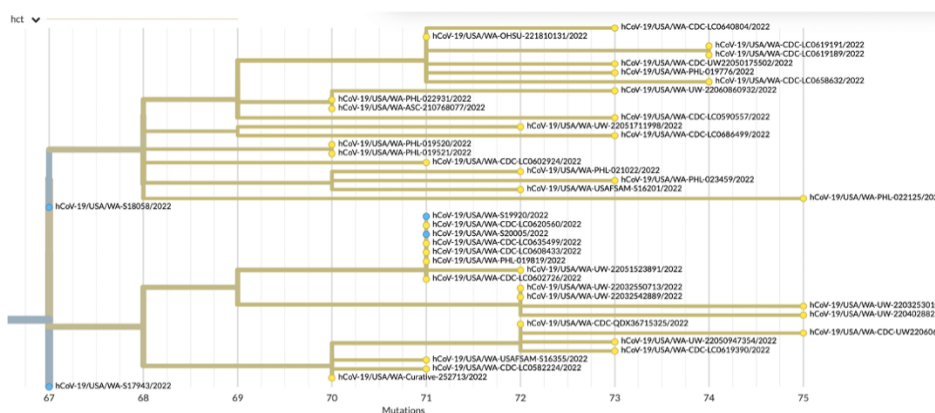
C



830
831
832
833

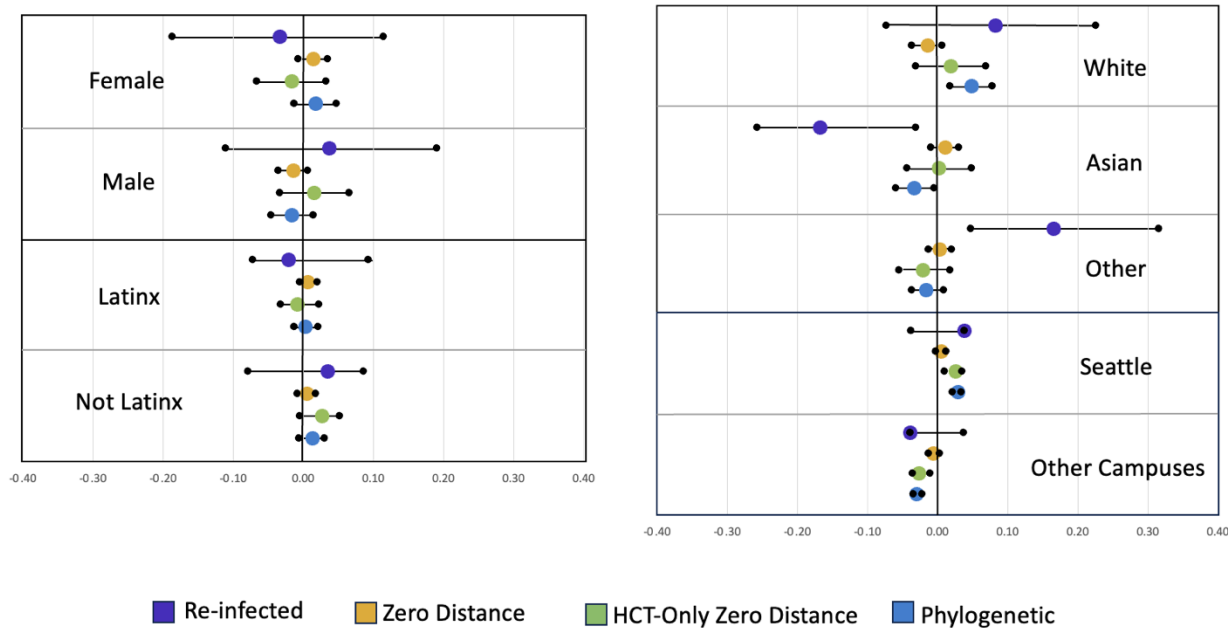
834 **Supplementary Figure S5: Average size of HCT sequence clusters by month.** A) Proportion of
835 total sequenced HCT specimens represented by each Nextstrain clade by month. B) Average
836 size of zero distance clusters of each Nextstrain clade by month. C) Average size of phylogenetic
837 clusters of each Nextstrain clade by month.

838
839
840
841
842
843
844
845
846
847
848



849
850 **Supplementary Figure S6: Phylogenetic tree showing HCT-only zero distance cluster with**
851 **largest number of non-HCT descendants.** HCT genomes are represented by blue nodes and
852 non-HCT genomes by yellow nodes.

853
854
855
856
857
858
859
860
861
862



863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878

Supplementary Figure S7: Demographic and epidemiologic characteristics of participants who experienced re-infection and of those with sequenced specimens in clusters. For each category (female, male, Latinx, Not Latinx, White, Asian, Other Race, Seattle campus, and other campuses), the mid-line ($x = 0$) represents the frequency of that category in the whole dataset. The purple, yellow, green, and blue dots give the percent difference between the frequency of that category among those who experienced re-infection and those with sequences in zero distance (groups of identical sequences), HCT-only zero distance (groups of identical sequences with haplotype unique to HCT), and phylogenetic clusters (groups of sequences that cluster phylogenetically), respectively, and the frequency of that category in the whole dataset. The intervals marked by black dots connected by a bar mark the 95% confidence interval of the values given by the purple, yellow, green, and blue dots.