

1 Exploring the Spatial Distribution of Persistent 2 SARS-CoV-2 Mutations - Leveraging mobility 3 data for targeted sampling

4 Riccardo Spott^{1,*}, Mathias W. Pletz^{1,2}, Carolin Fleischmann-Struzek^{1,2}, Aurelia Kimmig¹,
5 Christiane Hadlich³, Mathias Hauert³, Mara Lohde¹, Mateusz Jundzill¹, Mike Marquet¹, Petra
6 Dickmann⁴, Ruben Schüchner⁵, Martin Hölzer⁶, Denise Kühnert^{7,8}, Christian Brandt¹

7

8 ¹ Institute for Infectious Diseases and Infection Control, Jena University Hospital, Jena,
9 Germany

10 ² Center for Sepsis Control and Care, Jena University Hospital/Friedrich Schiller University
11 Jena, Jena, Germany

12 ³ SMA Development GmbH - epicinsights Agentur für Künstliche Intelligenz und Big Data
13 Analytics, Jena, Germany

14 ⁴ Department of Anaesthesiology and Intensive Care, Jena University Hospital, Jena,
15 Germany

16 ⁵ Thuringian State Authority for Consumer Protection, Department Health Protection

17 ⁶ Methodology and Research Infrastructure, Genome Competence Center (MF1), Robert
18 Koch Institute, Berlin, Germany

19 ⁷ Centre for Artificial Intelligence in Public Health Research, Robert Koch Institute, Berlin,
20 Germany

21 ⁸ Transmission, Infection, Diversification and Evolution Group, Max Planck Institute for
22 Geanthropology, 07745 Jena, Germany.

23 * Corresponding Author

24 Short Title

25 Mobility data for Surveillance

26 Keywords

27 Nanopore sequencing; SARS-CoV-2; WGS; mobility data; cluster tracking

28 Abstract

29 Given the rapid cross-country spread of SARS-CoV-2 and the resulting difficulty in tracking
30 lineage spread, we investigated the potential of combining mobile service data and fine-
31 granular metadata (such as postal codes and genomic data) to advance integrated genomic
32 surveillance of the pandemic in the federal state of Thuringia, Germany. We sequenced over
33 6,500 SARS-CoV-2 Alpha genomes (B.1.1.7) across seven months within Thuringia while
34 collecting patients' isolation dates and postal codes. Our dataset is complemented by over
35 66,000 publicly available German Alpha genomes and mobile service data for Thuringia. We
36 identified the existence and spread of nine persistent mutation variants within the Alpha
37 lineage, seven of which formed separate phylogenetic clusters with different spreading
38 patterns in Thuringia. The remaining two are sub-clusters. Mobile service data can indicate
39 these clusters' spread and highlight a potential sampling bias, especially of low-prevalence
40 variants. Thereby, mobile service data can be used either retrospectively to assess
41 surveillance coverage and efficiency from already collected data or to actively guide part of a
42 surveillance sampling process to districts where these variants are expected to emerge. The
43 latter concept proved successful as we introduced a mobility-guided sampling strategy for
44 the surveillance of Omicron sublineage BQ.1.1. The combination of mobile service data and
45 SARS-CoV-2 surveillance by genome sequencing is a valuable tool for more targeted and
46 responsive surveillance.

47 Introduction

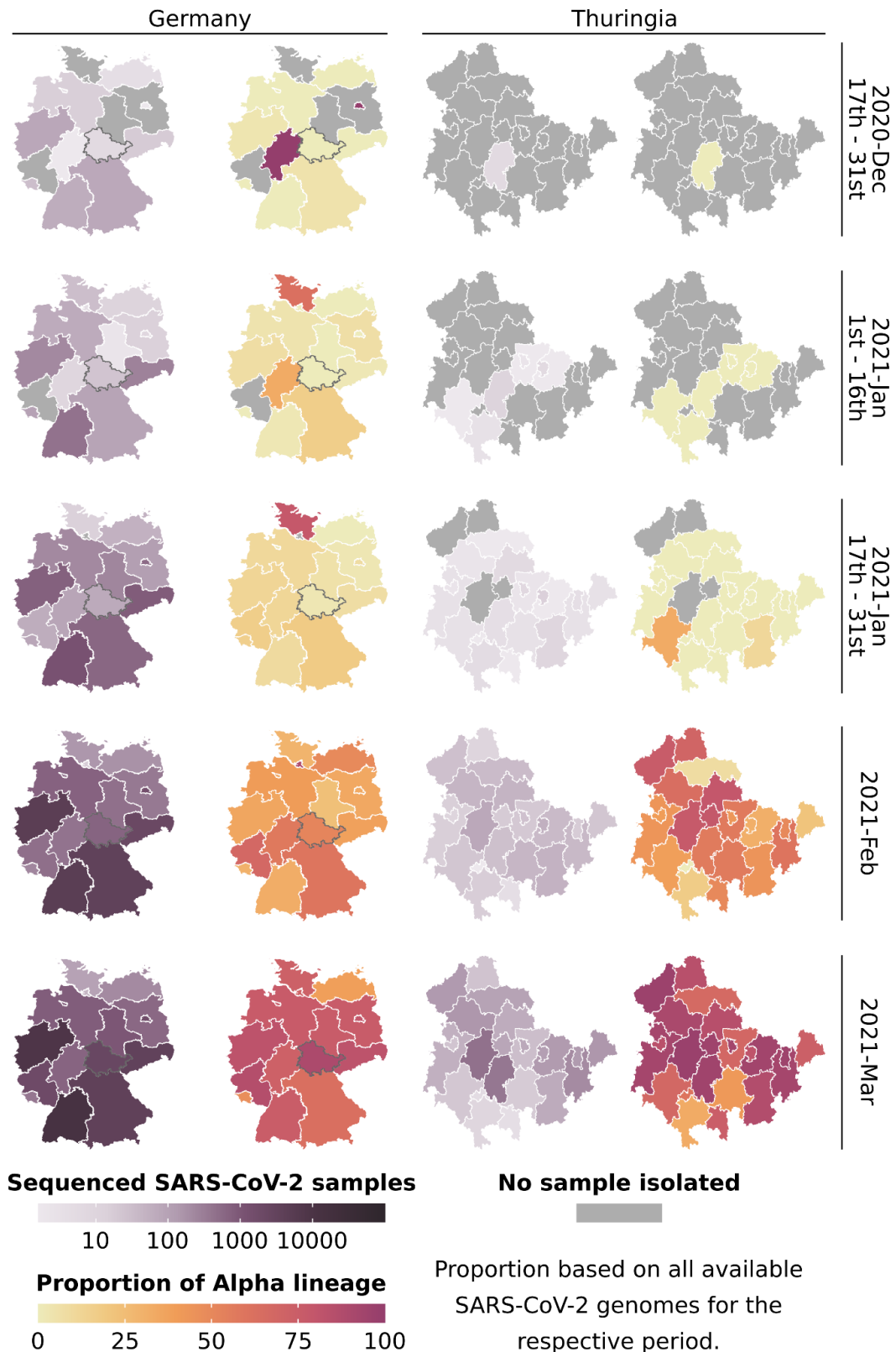
48 On March 11, 2020, the World Health Organization (WHO) classified the SARS-CoV-2 virus
49 (Severe Acute Respiratory Syndrome - Corona Virus - 2) as a global pandemic due to its
50 rapid spread and high infection rate.¹ The airborne virus has since caused significant
51 morbidity and mortality worldwide (<https://covid19.who.int/>). In an attempt to control its
52 spread, many countries initiated comprehensive surveillance efforts with molecular
53 techniques such as Polymerase Chain Reaction (PCR) and Whole Genome Sequencing
54 (WGS).^{2,3} Consequently, nearly 15.8 million SARS-CoV-2 sequences have been deposited
55 into the "Global Initiative on Sharing All Influenza Data" database (as of July 21, 2023,
56 GISAID). Many research groups have undertaken studies examining the viral spread by
57 integrating sequencing and epidemiological data to monitor the pandemic and investigate
58 local outbreaks.^{4,5} Most of these local projects are part of national surveillance programs
59 such as the UK's Genomics Consortium (COG-UK) or "national genomic surveillance" in the
60 USA.^{2,6} In Germany, the "Coronavirus-Surveillanceverordnung" (CorSurV) enacted by the
61 State Ministry of Health on January 19, 2021, mandated that laboratories with sequencing
62 capabilities process SARS-CoV-2-positive samples, offering financial compensation until
63 April 2023.³

64 Bioinformatics workflows developed in Germany, such as poreCov (for Oxford Nanopore
65 data) and CoVpipe2 (for Illumina data), reconstruct SARS-CoV-2 consensus genomes from
66 the sequencing data and prepare the results for upload and submission to the Robert Koch
67 Institute (RKI).^{7,8} As the German government's public health and biomedical research
68 institute responsible for disease control and prevention, the RKI collected the genomes via
69 the German Electronic Sequence Data Hub (DESH) and integrated them with additional
70 epidemiological information to provide an up-to-date overview of the ongoing viral spread.
71 For keeping track of the rapid SARS-CoV-2 evolution, PANGO (Phylogenetic Assignment of
72 Named Global Outbreak) provides a standard naming convention based on unique mutation
73 profiles and further criteria, resulting in the classification of over 3,660 lineages (as of August
74 2023).^{9,10} Additionally, the WHO classified important viral lineages as "Variants of Concern"
75 (VOC), "Variants of Interest" (VOI), or "Variants under Monitoring" (VUM), using Greek
76 designations in the past (e.g., "Alpha" (Pango lineage main designation B.1.1.7) or
77 "Omicron" (Pango lineage main designation B.1.1.529)). Further, the WHO also de-
78 escalated former VOCs to reflect the current SARS-CoV-2 variant landscape better. The first
79 defined VOC (now de-escalated), the Alpha lineage, rapidly replaced almost all previously
80 circulating lineages globally by the end of 2020 until the VOC Delta (main lineage B.1.617.2)
81 replaced it in mid-2021.¹¹⁻¹³

82 To predict or monitor the rapid viral spread throughout regions, various data types, like travel
83 data, passenger volumes, or passive wastewater monitoring were examined previously.^{14–16}
84 Furthermore, different studies explored mobility data with genomic data to retrace the origin
85 and spatial expanse of Alpha or utilized geolocation data to model the spread in metropolitan
86 areas to recreate case trajectories and the impact of mobility restrictions.^{17,18} Mobility data
87 was also used in Germany during the pandemic, revealing that lockdowns leave distant parts
88 of the country less connected due to the sharp decline in long-distance travel.¹⁹ These
89 studies focused on analyzing residential movement and contact tracing to evaluate and
90 inform health policies but were not applied to active molecular surveillance.
91 Here, we investigated whether mobile service data and fine-granular metadata (such as
92 postal codes and genomic data) can help predict the spread of the Alpha lineage or guide
93 the sampling for more targeted genomic surveillance with a focus on the German federal
94 state of Thuringia.

95 Results

96 The Alpha lineage spread rapidly through Thuringia, showing a pattern
97 similar to its nationwide spread
98 Thuringia is a rural federal state in central Germany with a population of 2.1 million and no
99 major airports (overview of Thuringia’s population density in Supplementary Figure S1). We
100 investigated if the spread of the Alpha lineage of SARS-CoV-2 behaved differently compared
101 to the whole of Germany. To understand its spread, we used 257,721 public SARS-CoV-2
102 genomes from Germany (excluding Thuringia; including 136,099 Alpha genomes) and 7,404
103 genomes from our own sequencing data for Thuringia (including 6,522 Alpha genomes) from
104 December 2020 to August 2021 (see Figure 1, Supplementary Figure S2, and
105 Supplementary Tables S1 and S2; for details, see Methods section "Alpha spread
106 datasets"). For Thuringia, district-level data (full postal code) per genome were available,
107 whereas, for Germany, only postal code data of the sending laboratories (referred to as
108 “primary diagnostic laboratory” by the RKI where the SARS-CoV-2 positive sample was
109 detected) and sequencing laboratories were publicly available.



110 **Figure 1: Total number of all sequenced SARS-CoV-2 samples (purple) and the proportion of**
 111 **the Alpha lineage for all sequenced samples (yellow-red) for each state of Germany and each**
 112 **district of Thuringia.** 257,721 publicly available German SARS-CoV-2 genomes and their metadata
 113 were used for the general German maps excluding data from Thuringia. For Thuringia, we always
 114 used 7,404 genomes and their metadata from our database for both the German and the Thuringian
 115 maps. Please note that for all states except Thuringia, we used the postal code of the sending
 116 laboratory as a proxy for the geographic location of a sample. Thuringia is highlighted by a grey
 117 border on the maps of Germany.
 118

119 In late December 2020, four federal states in Germany (from here on called states) reported
120 the first cases of the Alpha variant. Although sequencing was initially low, it gradually
121 increased in the following month. However, the Corona-Surveillance regulation was passed
122 at the end of January 2021, leading to a rapid increase in sampling and sequencing by
123 February since sequencing costs could be reimbursed. Even though Thuringia sequenced a
124 similar amount of SARS-CoV-2 samples compared to other German states (as shown in
125 Figure 1), the proportion of the Alpha variant to other lineages was relatively low. However,
126 the proportion of Alpha increased heavily in February.

127 By March, Alpha had spread to nearly all states and districts (districts are similar to counties
128 or provinces) in Germany (Median: 76.14 % compared to 35.92 % in February, excluding
129 Thuringia) and Thuringia (Median: 85.29 %, up from 50.00 % in February). So, there was no
130 noticeable difference in the Alpha proportions between Germany and Thuringia after
131 February. During the summer of June and July 2021, sequencing declined in Germany
132 (including Thuringia; Supplementary Figure S2) due to the decrease in overall daily cases,
133 as reported by Meintrup *et al.* and Oh and Hölzer *et al.*.^{20,21}

134 In summary, the spread of the Alpha lineage in Thuringia lagged roughly two weeks behind
135 the general spread in the rest of Germany but showed similar proportions. This suggests that
136 Thuringia experienced a delay in the initial arrival of Alpha. However, we did not observe any
137 difference in the overall spread afterward. Thuringia was among the first states to adopt new
138 containment measures, including contact limitations, closure of retail shops, and prohibition
139 of tourist journeys (14th December 2020). Jena, a city in Thuringia, was also the first
140 German city to implement mandatory public masking in March 2020.²² Contacts were further
141 restricted on January 9th, and people were urged to restrict their movement radius to 15 km,
142 which might explain the delay besides the absence of major airports nearby.

143 All Thuringian genomes were evenly distributed between other German samples in the
144 phylogenetic time tree (see Supplementary Figure S3). However, due to its rapid spread
145 from February onwards, it is difficult to accurately track how the Alpha lineage specifically
146 expanded (point of entries, exact origins, etc.). Consequently, we investigated whether
147 "sublineages" might be identifiable and trackable to address this.

148 Monitoring of Alpha subclusters in Thuringia reveals temporally and 149 regionally restricted distribution patterns

150 To identify possible clusters among the Alpha lineage spreading in Thuringia, we called each
151 Alpha genome's mutations via Nextclade by analyzing them using poreCov.^{7,23} We identified
152 nine clusters out of 70,429 Alpha genomes, based on their mutation profile, time period, and
153 phylogenetic distance (from here on called Alpha subclusters; for details, see Methods
154 "Subcluster identification"). All subclusters, their time period, and sample size in Thuringia

155 are summarized in Table 1. An overview of each subcluster (phylogenetic time tree, location,
 156 and period) is also provided on microreact.org as interactive views (see Methods “Subcluster
 157 identification”). Note that our subcluster definition is similar to the definition of a sublineage.
 158 However, during the Alpha wave, PANGO sublineages were only rarely defined (PANGO
 159 designation: Q.1 to Q.8; compared to the Delta and Omicron waves).

160 **Table 1: Overview of nine Alpha subclusters in Thuringia, their sample count, their time period,**
 161 **and their specific mutations that are shared across all members of the subcluster (excluding**
 162 **characteristic Alpha mutations that are shared across all subclusters).** The mutation used to
 163 define the subcluster is highlighted in bold.

Designation	Mutations	Number of samples	Time period	Remarks
1	S:H49Y , ORF1a:I841V	44	Feb-May 2021	S:H49Y eases cell entry in S-pseudotyped lentiviral system. ²⁴
2	S:N354K	63	Feb-May 2021	S:N354K slightly impaired mAb h11B11. ²⁵
3	S:G496S , ORF1a:E1013K	12	Mar-May 2021	S:G496S: compromises BA.1 replication fitness and changed mAb sensitivities, reduces ACE2 binding affinity, and increases immune evasion. ²⁶⁻²⁸
4	S:N703D , ORF1a:D1228G, ORF1a:A2123V	51	Mar-May 2021	-
5	S:T716V , N:G204P, ORF1a:D1600N	22	Apr-May 2021	-
6	S:S939F	206	Feb-May 2021	S:S939F: modulates T-cell propensity. ²⁹
6.1 [#]	S:V90F , S:S939F	55	Feb-May 2021	-
7	ORF1b:A520V	811	Feb-Jun 2021*	-
7.1 [§]	S:N185D , ORF1b:A520V, ORF1b:L1504F	40	Feb-May 2021	-

164 * only one sample for June # branch from subcluster 6 § branch from subcluster 7

165 Eight of these subclusters are based around a specific spike protein mutation, while the
 166 other contains a mutation within the ORF1b region. The subcluster 7.1 “S:N185D” branched
 167 out from the subcluster 7 “ORF1b:A520V” and subcluster 6.1 “S:V90F” is a branch of
 168 subcluster 6 “S:S939F” (see Table 1). These two branched subclusters still carry the specific
 169 mutation of their originating subcluster. The subclusters 3, 4, and 5 were observable
 170 between two and three months, and the other subclusters over at least four months. To
 171 investigate these subclusters' regional spread, each sample was mapped to its Thuringian
 172 district based on the resident's postal code it was isolated from. We then sorted the samples

173 according to their subclusters and visualized them throughout the subcluster's observed
174 period. The spread of two representative subclusters is exemplary visualized in Figure 2a,
175 and all the subclusters are available via Supplementary Figure S4 and their data in
176 Supplementary Table S3. Additionally, all subclusters and their metadata are also available
177 via Microreact (see Methods "Subcluster identification").

178

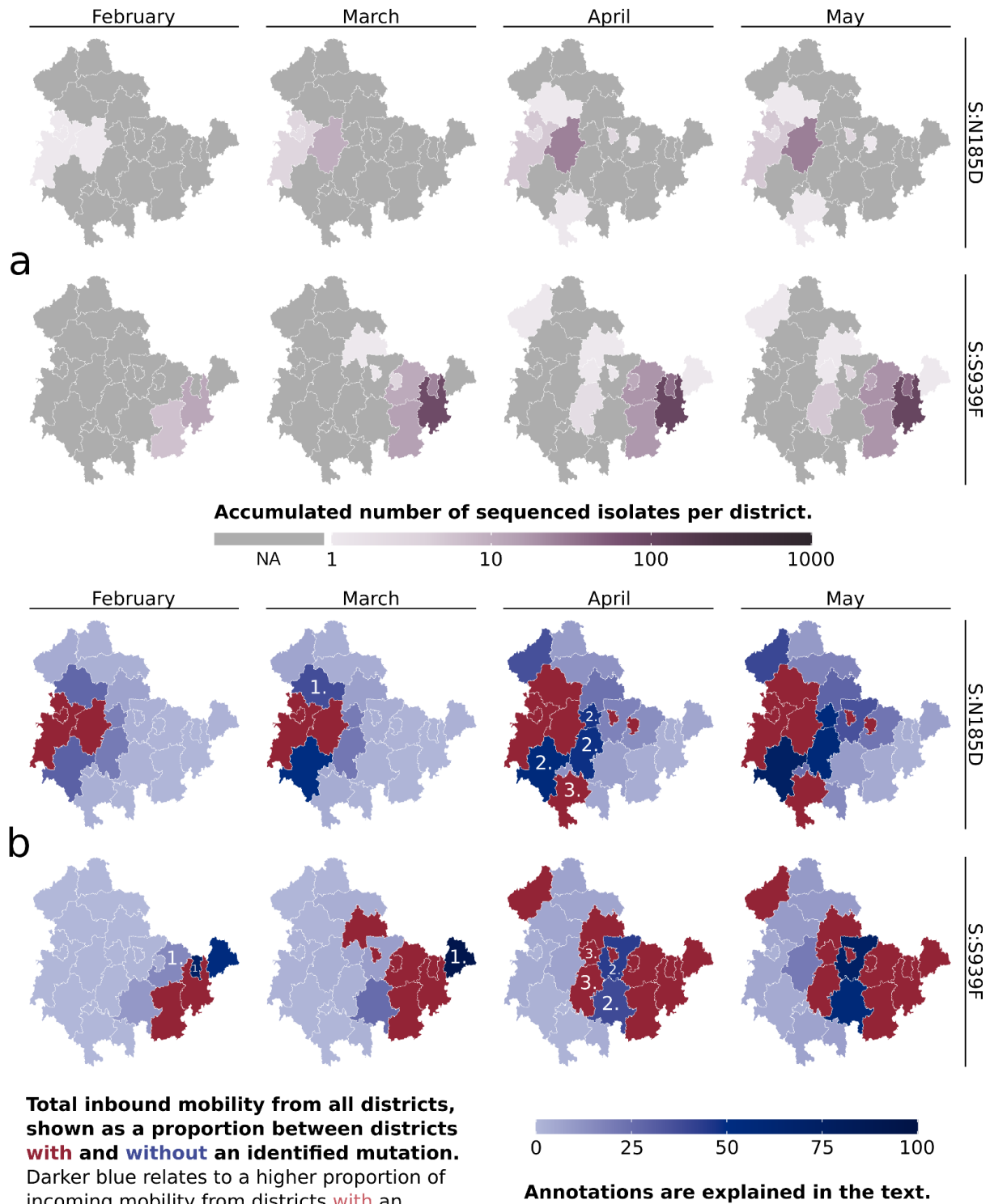
179 Each of the seven main mutation variant clusters originated from a different Thuringian
180 district. At the same time, two subclusters, 6.1 and 7.1, branched out from the same districts
181 as their original clusters (6 and 7), 12 or 13 days after their first emergence. The subclusters
182 mainly spread regionally confined and not across all of Thuringia (see Figure 2a) but were
183 also identified in other states of Germany (see "microreact.org"-project). For example, the
184 "S:S939F" subcluster spread across 15 states, with the first samples being isolated outside
185 of Thuringia. The eight Spike-mutation subclusters had expanded between four to twelve of
186 the 23 Thuringian districts within the observable time period of each subcluster. They
187 expanded by one to six districts per month, with a greater expansion mostly accompanied by
188 a larger increase in the subcluster sample number. In contrast, the ORF1b-variant even
189 comprised 21 districts and expanded between two to seven districts per month. Most of each
190 subcluster's samples were identified in their region of first occurrence, and no additional
191 samples were found after the given periods.

192 Several limitations need to be considered. The identified subclusters may have multiple
193 origins or may not originate from Thuringia. Due to the lack of precise zip codes (publicly
194 available German genomes are limited to postal codes of sending and sequencing
195 laboratories), monitoring the subclusters in other states on a district level was impossible.
196 Nevertheless, we could follow how the subclusters developed in Thuringia, even if multiple
197 origins may have affected the overall speed or length of each subcluster's occurrence.

198 Our surveillance sampling heavily relies on various institutions and partners, and only a
199 portion of the provided samples can be sequenced (see "Sampling" in Methods). For
200 example, the spread of subcluster "S:S939F" revealed two districts in April where no
201 respective samples were found (Figure 2a) despite being surrounded by districts with
202 "S:S939F"-samples present. This could be due to the lack of samples sent to sequencing
203 from those regions or the low prevalence. We, therefore, investigated if mobile service data
204 of residents, in addition to molecular surveillance, might be utilized to counteract this issue.

205 **Mobile service data indicates Alpha subcluster spread and sampling bias**
206 With the aim to predict the subcluster spread and, thereby, reduce surveillance-based
207 sampling bias, we utilized anonymized mobile service data from T-Systems International
208 GmbH. Around 200 million trips were used to determine the number of daily trips between

209 the Thuringian districts. We then combined this information with our fine-granular genomic
210 data to specify each district's monthly proportion of inbound mobility from subcluster-
211 affiliated districts (see Methods "Mobile service data"). The results are visualized in Figure
212 2b (complete overview in Supplementary Figure S5; data provided in Supplementary Table
213 S4).



214
 215 **Figure 2: Overview of the subclusters “S:N185D” and “S:S939F” in Thuringian districts.** a)
 216 Accumulated number of sequenced samples for each subcluster per district and per month. b)
 217 Combined visualization of each district’s “inbound mobility” from other districts (color intensity) and the
 218 occurrence of a subcluster sample (red = sample found, blue = no sample found). The inbound
 219 mobility of each district (color intensity) is shown as a proportion of incoming mobility from other
 220 districts with or without an identified sample. The darker the blue color of a district, the higher the
 221 proportion of inbound mobility from other districts with an identified subcluster sample (red districts).
 222 The light blue color describes that most of the inbound mobility of a district comes from other districts
 223 without an identified subcluster sample (blue districts). Numbers refer to district types 1, 2, and 3, as
 224 further defined in the main text.

225 The addition of mobile service data resulted in three different “types” of districts (see Figure
226 2b, annotated districts). Type 1 included districts with high inbound mobility from areas with
227 an identified variant, where the variant was eventually found afterward, while Type 2 were
228 districts with high inbound mobility from areas with an identified variant, where the variant
229 was never identified. Type 3 included districts not directly connected to a district with an
230 identified variant, but a variant was eventually identified while they border Type 2 district(s).
231 Our previous analysis of the subclusters' spreading pattern across the districts, based solely
232 on identified variants, indicated missed identifications in some districts due to the seemingly
233 illogical spread to districts without a connection to others (Figure 2a). The inclusion of mobile
234 service data revealed some of these districts to be Type 2 districts. This suggests that the
235 specific variant should be identifiable within these districts due to the observed high
236 incoming mobility from districts with identified variants. Type 2 districts were mainly observed
237 for subclusters with low prevalence and, consequently, low numbers of covered samples that
238 are usually more difficult to monitor. For example, we assumed missing identifications in
239 some districts of subclusters 1, 2, and 3, which through the mobile service data, are now
240 partially identified as Type 2 districts. In contrast, for fast-spreading, highly prevalent
241 subclusters, the regional coverage aligned well with the mobile service data, such as
242 subcluster 7 covering 811 samples (see Supplementary Figure S5).
243 Despite analyzing the mobile service data of districts from other federal states than
244 Thuringia, we could not apply them, as the lack of precise location data for samples outside
245 of Thuringia prevented the correct calculation of the incoming mobility. Based on the nine
246 observable clusters, we concluded that mobile service data might be a good prediction
247 marker for the spread of high-prevalence variants but, more importantly, a good indication of
248 districts that should have an identified low-prevalence variant. Next, we investigated if mobile
249 service data can improve active surveillance via guiding sample collection for genomic
250 sequencing.

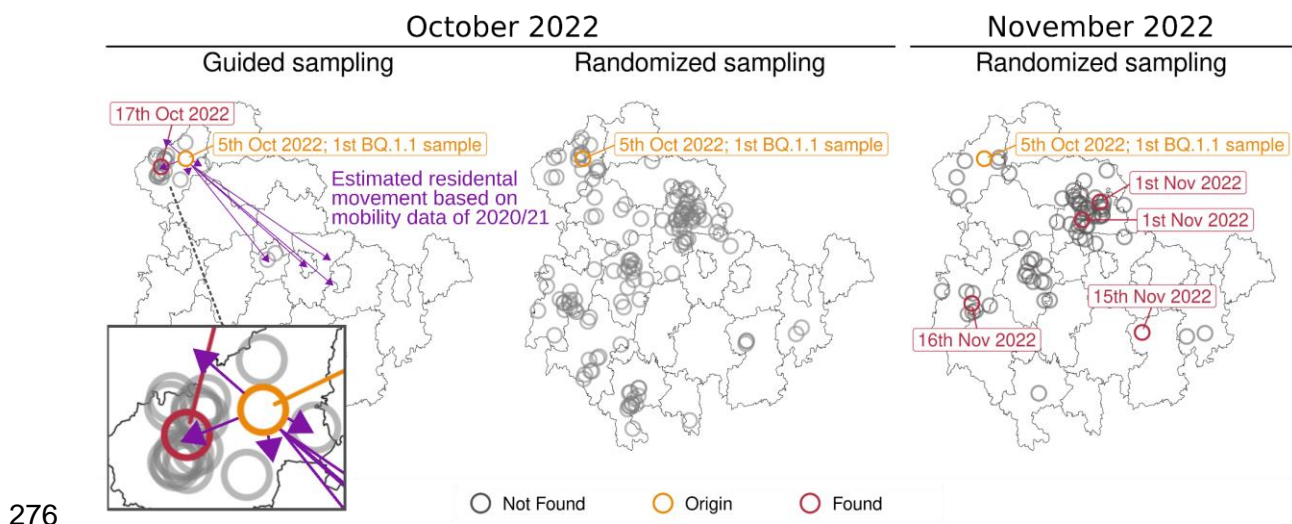
251 Proof of principle: Mobile service data-guided sampling for genomic 252 surveillance for Omicron BQ.1.1

253 Based on our previous findings, we implemented the “mobility-guided” sampling approach
254 under real pandemic circumstances over one month in addition to our active surveillance.

255 As the subject of investigation, we searched for a newly emerging (based on global news
256 reports) and ideally low prevalent SARS-CoV-2 lineage in Thuringia.

257 Among the various emerging Omicron sublineages during that time, sublineage BQ.1.1
258 fulfilled the defined criteria. First isolated in a northwestern-Thuringian community on
259 October 5, 2022, we identified this particular sublineage on October 14, 2022, among a
260 routine batch of 42 samples. BQ.1.1 was a low-prevalence sublineage that was identified

261 worldwide (<https://outbreak.info/situation-reports?pango=BQ.1.1>). Following its first
262 Thuringian identification, we utilized the past two years of mobile service data (2020/21) to
263 investigate the residential movements for the community of first detection. As a result, we
264 identified eight communities with the most residential movement from the originating
265 community (four in central and three in NW of Thuringia, one in NW-neighboring state
266 Saxony-Anhalt; purple arrows in Figure 3, “Guided sampling” in October 2022).
267 Subsequently, we specifically requested all the available samples from these communities
268 and collected 19 additional samples (isolated between the 17th and 25th of October 2022)
269 besides the randomized sampling strategy. More samples could not be included as sampling
270 was restricted to the submission of third parties we had no influence over. As part of the
271 general Thuringian surveillance, we collected 132 samples for October (covering dates
272 between the 5th and 31st) and 69 samples in November (covering dates between the 1st
273 and 25th; see “Randomized Sampling” for the according months in Figure 3). Randomized
274 sampling was not influenced or adjusted based on the mobility-guided sample collection. A
275 complete overview of all samples is provided in Supplementary Table S5.



276
277 **Figure 3: Overview of the mobility-guided sampling of the Omicron sublineage BQ.1.1 in**
278 **Thuringia compared to the default randomized sampling (surveillance) in October 2022.** For
279 clarity, the surveillance results in November 2022 were added to highlight the spreading progress of
280 BQ.1.1. Circles reflect the location of each sample (based on residents' zip codes). Orange circle:
281 First identified BQ.1.1 sample; Red circle: Additionally identified BQ.1.1 sample; Grey: Sample of
282 another lineage. Purple arrows show the eight residential movements with the most participants from
283 the originating community of the first identified BQ.1.1 sample (orange circle). The mobile service data
284 were extracted from 2020/21.

285 Among the 19 samples specifically collected based on mobile service data, we identified one
286 additional sample of the specific Omicron sublineage BQ.1.1 close to the originating
287 community. During the same period, our randomly sampled routine surveillance strategy did
288 not detect another sample (October 2022, Figure 3). Only in the one-month follow-up,
289 another four samples were identified across Thuringia through routine surveillance

290 (November 2022, Figure 3). During our attempt to implement the mobility-guided sampling
291 approach in real-time during the pandemic, we encountered three distinct limitations, some
292 of which are commonly observed in surveillance practices. The guided sampling depended
293 on the individual sample submitting institutions, affecting the availability of suitable samples,
294 especially for the communities of interest. By choosing a newly emerging Omicron
295 sublineage for our experiment, spread and, therefore, suitability were uncertain. In our case,
296 BQ.1.1's prevalence in Thuringia was even lower than expected, also remaining rare in
297 subsequent months, with only 42 samples found until June 2023, eight months after the first
298 occurrence in Thuringia. Due to the short preparation time, only mobile service data from the
299 past two years and no current data were available. Nevertheless, the available datasets still
300 reflect pandemic movement behavior since the pandemic was already ongoing for two years.
301 In summary, increasing the sampling depth in the suspected regions resulted in successfully
302 identifying the specified lineage using only a fraction of samples in contrast to the
303 randomized surveillance. Implementing such an approach effectively under pandemic
304 conditions poses difficult challenges due to the fluctuating sampling sizes. Although the
305 finding of the sample may have been coincidental, our proof of concept demonstrated how
306 we can leverage the potential of mobile service data for targeted surveillance sampling.

307 Conclusion

308 During the SARS-CoV-2 pandemic, diverse data sources like travel, wastewater, and
309 mobility data have been employed in surveillance and transmission tracking.¹⁴⁻¹⁹ In the
310 present study, we analyzed over 265,000 German SARS-CoV-2 genomes to examine
311 whether mobile service data can predict the spatial distribution of the Alpha lineage in the
312 German state of Thuringia and how they potentially benefit pandemic surveillance.
313 Our study shows that the absence of major transport hubs in Thuringia initially delayed the
314 spread of Alpha. However, its impact on the total distribution is limited, and the spread was
315 ultimately comparable between Germany and Thuringia. While our findings on mobile
316 service data may, therefore, also apply to Germany, we could not verify this because the
317 limited location data of publicly available German genomes prevented in-depth investigations
318 outside of Thuringia. Thus, precise patient location data are crucial to utilize mobile service
319 data in genomic surveillance, but privacy regulations may restrict access to this data. Shortly
320 after its emergence, Alpha formed mutation variants like the known sublineages Q.1 to Q.8
321 and the Thuringian subclusters identified by us. This reflects the ongoing evolution during
322 active circulation and indicates an even greater sublineage diversity, which has not been
323 surveyed as closely as in the subsequent Delta and Omicron waves. By monitoring the nine
324 Thuringian subclusters, rather than focusing solely on the parental lineage B.1.1.7, we were

325 again able to effectively track transmissions and gain a comprehensive understanding of the
326 regional spread. So, it underscores the importance of sequencing in pandemic surveillance
327 to explore such genomic changes and, thereby, keep track of the transmission chains and
328 potential outbreaks.

329 Mobile service data can support such surveillance in different ways. Previous studies
330 examined the capabilities of mobility data in the context of, e.g., case trajectories, but
331 retrospectively applied to already collected data, it can be used to examine surveillance
332 sampling coverage and possible sampling bias. We highlighted this approach exemplary
333 with the Alpha lineage, where mobile service data indicated an assumed sampling bias and
334 partially predicted the spread of our Thuringian subclusters. Another approach is to actively
335 guide the sampling process through the usage of mobile service data, which we
336 demonstrated with our proof of principle focusing on the Omicron-lineage BQ.1.1. The
337 allocation of surveillance resources for guided sampling can be flexibly adjusted to adapt to
338 specific circumstances and maximize efficiency. This allows for an increased regional
339 sampling coverage without increasing the general sampling volume, which especially helps
340 identify low-prevalence variants. We recommend using mobile service data as a supportive
341 element for general randomized surveillance to retrospectively evaluate its efficiency and
342 actively help screen for low-prevalence variants. These findings should also apply to other
343 surveillance efforts. Yet the feasibility depends on the availability and cost of such mobile
344 service data. Alternatively, financial resources could also be invested directly in increasing
345 sampling capacity and coverage, which ultimately depends on individual factors of the
346 respective surveillance. Mobile service data can also be used with other surveillance
347 approaches and elements. For example, wastewater surveillance can give further indications
348 to supplement guided sampling. At the same time, passenger data offers additional insights
349 into traffic hubs as sources of regional movement.

350 Methods

351 Sampling

352 Starting mid-2020, we initially sequenced hospital-intern samples, transitioning by January
353 2021 to approximately 43 PCR-positive samples per week: 20 from the hospital's
354 microbiology department and 23 randomly sourced by the Thuringian State Authority for
355 Consumer Protection ("Thüringer Landesamt für Verbraucherschutz"; TLV).
356 Until June 2023, our institute sequenced 3,770 SARS-CoV-2 samples, and SYNLAB Holding
357 Deutschland GmbH, Bioscientia Healthcare GmbH, and DIANOVIS GmbH provided
358 additional 7,800 Thuringian SARS-CoV-2 genomes and their metadata.

359 Sample preparation and sequencing

360 RNA isolation used the ZymoResearch “Quick-RNA Viral Kit” (Zymo Research Europe
361 GmbH, Germany, Product-ID: R1035), according to the manufacturer's instructions with
362 100 µl patient sample input and a centrifuge speed of 16,000 g.

363 The viral RNA underwent a Reverse Transcriptase (RT)-PCR followed by a multiplex-PCR
364 using the ARTIC V1200 primer set, according to Freed and Silander’s SARS-CoV-2
365 sequencing protocol (version 4, updating to version 5 by March 2021).³⁰ Subsequent DNA
366 quantification utilized the Qubit dsDNA HS assay (Invitrogen, USA).

367 From the amplified DNA, a sequencing library was prepared using the Nanopore SQK-
368 LSK109 and SQK-RBK004 kits (Oxford Nanopore Technologies, Oxford, UK), sequenced for
369 a maximum of 72 h utilizing an Oxford Nanopore MinION Mk1b sequencer with R.9-flowcells
370 and the MinKNOW software (versions MKE_1013_v1_revBC_11Apr2016 to
371 MKE_1013_v1_revBR_11Apr2016 in the respective period), and analyzed with the software
372 pipeline poreCov (versions 0.3.5 to 0.11.7; including basecalling, demultiplexing, adapter
373 removal, quality filtering, and genome alignment) to reconstruct consensus genomes.⁷

374 Sequencing data and the respective metadata (e.g., isolation date, sending laboratory
375 details) were submitted to the RKI through DESH. We also collected the postal code of the
376 isolation location or at least of the sending local health authority, storing all data additionally
377 in a local database.³¹ Due to data protection, such data is limited on the RKI's public GitHub
378 repository ([https://github.com/robert-koch-institut/SARS-CoV-2-](https://github.com/robert-koch-institut/SARS-CoV-2-Sequenzdaten-aus-Deutschland)
379 [Sequenzdaten aus Deutschland](https://github.com/robert-koch-institut/SARS-CoV-2-Sequenzdaten-aus-Deutschland)), providing instead postal codes of the sequencing and
380 sending laboratories.

381 Alpha spread datasets

382 From our local database, we extracted 8,397 samples with isolation dates before Oct 1st,
383 2021. After adding federal state and district information, 993 entries with non-Thuringian
384 locations were excluded, yielding 7,404 samples (including 6,522 Alpha genomes).

385 The publicly available RKI SARS-CoV-2 dataset was downloaded, containing 1,091,655
386 genomes with the respective metadata (17th Oct 2022; Zenodo-version 2022-10-16).³²
387 789,405 entries, isolated after Sep’21, and 64 entries without “sending laboratory”
388 information were removed. For the resulting 302,186 entries, location information (location,
389 federal state, district, longitude, latitude) were added based on the sending laboratory postal
390 code. Five entries with a non-existing postal code and all 44,465 Thuringian samples were
391 removed from the dataset, resulting in 257,721 samples (including 136,099 Alpha genomes).
392 Analyzing both datasets, we calculated the monthly proportion of Alpha lineage samples in
393 Thuringia and Germany per state/district, dividing Dec’20 and Jan’21 into first and second
394 halves.

395 Subcluster identification

396 Using a total of 70,429 German and Thuringian Alpha genomes, a phylogenetic time tree
397 was created (see Supplementary Method “Phylogenetic time tree construction” and
398 Supplementary Figure S3). We determined the frequency of all non-Alpha-specific mutations
399 among the 6,522 Thuringian Alpha genomes. We then screened for mutations present in at
400 least 20 genomes with a small phylogenetic distance and a time occurrence of at least two
401 months. This led to nine mutations, each of them creating a defined cluster covering
402 between 12 and 811 closely related genomes. We only kept mutation information of these
403 nine subclusters in the respective metadata, which, together with the tree file of the
404 phylogenetic time tree, was uploaded to a “microreact.org”-project, provided as
405 Supplementary File 1 and found under the following link:

406 <https://microreact.org/project/ftR2GfjF6iXtSwbmN4ARTx-thuringianalpha-linclusters#76ir->
407 [complete-overview](#).

408 Mobile service data

409 T-Systems International GmbH collected and aggregated mobile service data via the Cell ID
410 method, dividing a geographical area into so-called traffic cells. Each cell is assigned to
411 exactly one transmitter mast, with a spatial resolution from 500 m x 500 m up to 8 km x 8 km
412 (depending on the transmitter mast network density). Cell phones always register to the
413 closest traffic cell, which is recorded and stored in an Origin-Destination Matrix (ODM). For
414 population representation, the data was extrapolated with Deutsche Telekom's market
415 share. Due to data privacy, the registration data is combined into movement streams
416 between traffic cells, the status resolution is reduced to one hour (greater time intervals =
417 less resolution), and individual traffic cells are grouped into districts. The degree of
418 anonymization (k-value = 30, data privacy regulation) removed movement streams with less
419 than 30 participants, resulting in approximately 200 Mio trips in the ODM. SMA Development
420 GmbH analyzed all movements between the single Thuringian districts, adding each Alpha
421 sample's isolation time and location data (per subcluster). The movements were further
422 divided by months and originating district (subcluster-affiliated vs. -unaffiliated), determining
423 each district's monthly inbound mobility proportion from cluster-affiliated districts.

424 Research in Context

425 Evidence before this study

426 We searched Pubmed for studies about the use of mobile service data for surveillance
427 written in English. For the broadest possible search, we included any publication covering
428 mobile data and surveillance aspects, using the following search string: ("cellular data" OR

429 "cell phone data" OR "mobility data" OR "movement data" OR "migration data" OR "phone
430 data") AND ("Surveillance" OR "Monitoring" OR "Survey" OR "Pandemic" OR "Disease" OR
431 "Epidemic" OR "Outbreak"). Our search yielded 1,285 publications published between 1966
432 and 2023. We manually screened all these publications but found no study that applied
433 mobile service data for active, targeted surveillance. Across all studies, the general focus
434 was on tracking contacts or analyzing movements to assess, for instance, the efficiency of
435 non-pharmaceutical interventions or generate prediction models. Some studies suggested
436 targeted surveillance based on their results, but it was not yet applied. Additionally, we used
437 "suite.ai" and "chatGPT" (with BING-search access) to let them search for "studies that
438 utilize mobile service data to guide the sampling process for infectious disease surveillance".
439 While "suite.ai" found two studies and "chatGPT" found another ten studies and reviews,
440 none covered the direct application of the mobility data in active surveillance.

441 Added value of this study

442 This study highlights the value of combining mobile service data with fine-granular metadata
443 for integrated genomic surveillance during the SARS-CoV-2 pandemic in a German federal
444 state. We illustrated this strategy with the Omicron sublineage BQ.1.1 and how to guide the
445 sampling processes toward areas where the new variant was expected to emerge.
446 Additionally, we used mobile service data during the pandemic to assess our sampling
447 coverage. Our study is the first to actively guide part of the genomic surveillance process
448 during a pandemic.

449 Implications of all the available evidence

450 Efficient molecular surveillance setups are crucial in managing outbreaks from the local to
451 the global scale. Different data sources are investigated to increase this efficiency,
452 addressing factors like the more efficient usage of scarce surveillance resources and the
453 prediction of spread. Extending molecular surveillance with such data should improve the
454 future management of pandemics and outbreaks.

455 Author contributions

456 Sample collection and preparation, R.Sp, R.Sc., M.L., and M.M.; sequencing, R.Sp., M.L.,
457 and M.M.; database setup and maintenance, M.J.; software, C.B.; bioinformatic analysis,
458 R.Sp., M.L., M.M, and C.B.; mobile service data analysis, C.H., and M.Ha.; literature
459 research, R.Sp.; writing first draft, R.Sp.; reviewing and editing manuscript, R.Sp., C.B.,
460 M.M., C.F.-S., A.K., C.H., M.Ha., M.Hö., D.K., R.Sc., P.D., and M.W.P.; supervision, C.B.;

461 project administration, C.B.; funding acquisition, M.W.P.. All authors have read and agreed
462 to the published version of the manuscript.

463 Funding

464 This work was supported by grants from the Federal Ministry of Education and Research
465 (project “SARS-CoV-2Dx“), [grant number 13N15745] and the Thüringer Aufbaubank
466 (project “Pandemie Analyse mittels Advanced Analytics Methoden“), [grant number 2021
467 VF 0035]. We acknowledge support by the German Research Foundation Projekt-Nr.
468 512648189 and the Open Access Publication Fund of the Thueringer Universitaets- und
469 Landesbibliothek Jena.

470 Conflicts of interest

471 The authors have declared that no competing interests exist.

472 Data sharing statement

473 All genomic data (genomes and respective metadata) are available in the provided microreact
474 project (project file available under <https://osf.io/n5qj6/>;
475 [https://microreact.org/project/ftR2GfjF6iXtSwbmN4ARTx-thuringianalpha-linclusters#76ir-](https://microreact.org/project/ftR2GfjF6iXtSwbmN4ARTx-thuringianalpha-linclusters#76ir-complete-overview)
476 [complete-overview](https://microreact.org/project/ftR2GfjF6iXtSwbmN4ARTx-thuringianalpha-linclusters#76ir-complete-overview)). The mobile service data used in this study can only be published in
477 processed form (available under <https://osf.io/n5qj6/>). The original mobile service data can
478 not be made public due to legal reasons/ownership.

479 Declaration of generative AI and AI-assisted technologies in the 480 writing process

481 During the preparation of this work, the author(s) used Grammarly in order to correct general
482 English and improve readability. After using this tool/service, the author(s) reviewed and
483 edited the content as needed and take(s) full responsibility for the content of the publication.

484 References

- 485 1 Zhu N, Zhang D, Wang W, *et al.* A Novel Coronavirus from Patients with Pneumonia in
486 China, 2019. *N Engl J Med* 2020; **382**: 727–33.
- 487 2 An integrated national scale SARS-CoV-2 genomic surveillance network. *Lancet Microbe*
488 2020; **1**: e99–100.
- 489 3 German ‘Corona-Surveillanceverordnung’ as issued on the 18th of January 2021 and
490 updated on the 27th of June 2022.
491 https://www.bundesgesundheitsministerium.de/fileadmin/Dateien/3_Downloads/C/Corona

- 492 virus/Verordnungen/CorSurV_BAnz_AT_19.01.2021_V2.pdf;
493 <https://www.bundesanzeiger.de/pub/publication/hEG1y5vUCsSa43vJZSP/content/hEG1y5vUCsSa43vJZSP/BAnz%20AT%2028.06.2022%20V1.pdf?inline>.
494
495 4 Meredith LW, Hamilton WL, Warne B, *et al.* Rapid implementation of SARS-CoV-2
496 sequencing to investigate cases of health-care associated COVID-19: a prospective
497 genomic surveillance study. *Lancet Infect Dis* 2020; **20**: 1263–71.
498 5 Page AJ, Mather AE, Le-Viet T, *et al.* Large-scale sequencing of SARS-CoV-2 genomes
499 from one region allows detailed epidemiology and enables local outbreak management.
500 *Microb Genomics* 2021; **7**. DOI:10.1099/mgen.0.000589.
501 6 Lambrou AS, Shirk P, Steele MK, *et al.* Genomic Surveillance for SARS-CoV-2 Variants:
502 Predominance of the Delta (B.1.617.2) and Omicron (B.1.1.529) Variants — United
503 States, June 2021–January 2022. *MMWR Morb Mortal Wkly Rep* 2022; **71**: 206–11.
504 7 Brandt C, Krautwurst S, Spott R, *et al.* poreCov-An Easy to Use, Fast, and Robust
505 Workflow for SARS-CoV-2 Genome Reconstruction via Nanopore Sequencing. *Front*
506 *Genet* 2021; **12**: 711437.
507 8 Lataretu M, Drechsel O, Kmiecinski R, Trappe K, Hölzer M, Fuchs S. Lessons learned:
508 overcoming common challenges in reconstructing the SARS-CoV-2 genome from short-
509 read sequencing data via CoVpipe2. *F1000Research* 2023; **12**: 1091.
510 9 Rambaut A, Holmes EC, O’Toole Á, *et al.* A dynamic nomenclature proposal for SARS-
511 CoV-2 lineages to assist genomic epidemiology. *Nat Microbiol* 2020; **5**: 1403–7.
512 10 PANGO SARS-CoV-2 lineage overview; accessed 18th August 2023.
513 https://github.com/cov-lineages/lineages-website/blob/master/_data/lineage_data.full.json.
514 11 Washington NL, Gangavarapu K, Zeller M, *et al.* Emergence and rapid transmission of
515 SARS-CoV-2 B.1.1.7 in the United States. *Cell* 2021; **184**: 2587-2594.e7.
516 12 Walker AS, Vihta K-D, Gethings O, *et al.* Tracking the Emergence of SARS-CoV-2 Alpha
517 Variant in the United Kingdom. *N Engl J Med* 2021; **385**: 2582–5.
518 13 Michaelsen TY, Bennedbæk M, Christiansen LE, *et al.* Introduction and transmission of
519 SARS-CoV-2 lineage B.1.1.7, Alpha variant, in Denmark. *Genome Med* 2022; **14**: 47.
520 14 Alpert T, Brito AF, Lasek-Nesselquist E, *et al.* Early introductions and transmission of
521 SARS-CoV-2 variant B.1.1.7 in the United States. *Cell* 2021; **184**: 2595-2604.e13.
522 15 O’Toole Á, Hill V, Pybus OG, *et al.* Tracking the international spread of SARS-CoV-2
523 lineages B.1.1.7 and B.1.351/501Y-V2 with grinch. *Wellcome Open Res* 2021; **6**: 121.
524 16 Li J, Ahmed W, Metcalfe S, *et al.* Monitoring of SARS-CoV-2 in sewersheds with low
525 COVID-19 cases using a passive sampling technique. *Water Res* 2022; **218**: 118481.
526 17 Kraemer MUG, Hill V, Ruis C, *et al.* Spatiotemporal invasion dynamics of SARS-CoV-2
527 lineage B.1.1.7 emergence. *Science* 2021; **373**: 889–95.
528 18 Chang S, Pierson E, Koh PW, *et al.* Mobility network models of COVID-19 explain
529 inequities and inform reopening. *Nature* 2021; **589**: 82–7.
530 19 Schlosser F, Maier BF, Jack O, Hinrichs D, Zachariae A, Brockmann D. COVID-19
531 lockdown induces disease-mitigating structural changes in mobility networks. *Proc Natl*
532 *Acad Sci* 2020; **117**: 32883–90.
533 20 Meintrup D, Nowak-Machen M, Borgmann S. A Comparison of Germany and the United
534 Kingdom Indicates That More SARS-CoV-2 Circulation and Less Restrictions in the Warm
535 Season Might Reduce Overall COVID-19 Burden. *Life* 2022; **12**: 953.
536 21 Oh DY, Hölzer M, Paraskevopoulou S, *et al.* Advancing Precision Vaccinology by
537 Molecular and Genomic Surveillance of Severe Acute Respiratory Syndrome Coronavirus
538 2 in Germany, 2021. *Clin Infect Dis* 2022; **75**: S110–20.
539 22 Pletz MW, Steiner A, Kesselmeier M, *et al.* Introduction of mandatory masking in health
540 care and community: experience from Jena, Germany. *Infection* 2023; published online
541 April 3. DOI:10.1007/s15010-023-02015-w.
542 23 Aksamentov I, Roemer C, Hodcroft E, Neher R. Nextclade: clade assignment, mutation
543 calling and quality control for viral genomes. *J Open Source Softw* 2021; **6**: 3773.
544 24 Ozono S, Zhang Y, Ode H, *et al.* SARS-CoV-2 D614G spike mutation increases entry
545 efficiency with enhanced ACE2-binding affinity. *Nat Commun* 2021; **12**: 848.
546 25 Du Y, Shi R, Zhang Y, *et al.* A broadly neutralizing humanized ACE2-targeting antibody

- 547 against SARS-CoV-2 variants. *Nat Commun* 2021; **12**: 5000.
- 548 26Liang R, Ye Z-W, Ong CP, *et al*. The spike receptor-binding motif G496S substitution
549 determines the replication fitness of SARS-CoV-2 Omicron sublineage. *Emerg Microbes*
550 *Infect* 2022; **11**: 2093–101.
- 551 27Kimura I, Yamasoba D, Nasser H, *et al*. The SARS-CoV-2 spike S375F mutation
552 characterizes the Omicron BA.1 variant. *iScience* 2022; **25**: 105720.
- 553 28Asif A, Ilyas I, Abdullah M, Sarfraz S, Mustafa M, Mahmood A. The Comparison of
554 Mutational Progression in SARS-CoV-2: A Short Updated Overview. *J Mol Pathol* 2022; **3**:
555 201–18.
- 556 29Donzelli S, Spinella F, di Domenico EG, *et al*. Evidence of a SARS-CoV-2 double Spike
557 mutation D614G/S939F potentially affecting immune response of infected subjects.
558 *Comput Struct Biotechnol J* 2022; **20**: 733–44.
- 559 30Freed N, Silander O. SARS-CoV2 genome sequencing protocol (1200bp amplicon
560 "midnight" primer set, using Nanopore Rapid kit) v5. 2021
561 DOI:10.17504/protocols.io.btsrind6.
- 562 31Jundzill M, Spott R, Lohde M, Hölzer M, Viehweger A, Brandt C. Managing and
563 monitoring a pandemic: showcasing a practical approach for the genomic surveillance of
564 SARS-CoV-2. *Database* 2023; **2023**: baad071.
- 565 32Koch-Institut R. SARS-CoV-2 Sequenzdaten aus Deutschland. 2022; published online
566 Oct 16. DOI:10.5281/ZENODO.7212725.