1    **SARS-CoV-2 shifting transmission dynamics and hidden reservoirs limited the efficacy**

2    **of public health interventions in Italy**

3

4    Marta Giovanetti[1,2,3]*, Eleonora Cella[4]*, Francesca Benedetti[5]*, Brittany Rife

5    Magalis[6]*, Vagner Fonseca[2,7,8], Silvia Fabris[3], Giovanni Campisi[9], Alessandra Ciccozzi[3],

6    Silvia Angeletti[10], Alessandra Borsetti[11], Vittoradolfo Tambone[12], Caterina Sagnelli[13],

7    Stefano Pascarella[14], Alberto Riva[15], Giancarlo Ceccarelli[16], Alessandro Marcello[17], Taj

8    Azarian[4], Eduan Wilkinson[7], Tulio de Oliveira[7], Luiz Carlos Junior Alcantara[1,2], Roberto

9    Cauda[18], Arnaldo Caruso[9], Natalie E Dean[19], Cameron Browne[20], Jose Lourenco[21], Marco

10    Salemi[6^], Davide Zella[5^], Massimo Ciccozzi[3^]

11

12    [1]Laboratório de Flavivírus, Instituto Oswaldo Cruz, Fundação Oswaldo Cruz, Rio de

13    Janeiro, Brazil; [2]Laboratório de Genética Celular e Molecular, ICB, Universidade Federal de

14    Minas Gerais, Belo Horizonte, Minas Gerais, Brazil; [3]Medical Statistic and Molecular

15    Epidemiology Unit, University of Biomedical Campus, Rome, Italy; [4]Burnett School of

16    Biomedical Sciences, University of Central Florida, Orlando (Florida, USA); [5]Institute of

17    Human Virology, Department of Biochemistry and Molecular Biology, School of Medicine,

18    University of Maryland (Baltimore, USA); [6]Emerging Pathogens Institute & Department of

19    Pathology, College of Medicine, University of Florida, Gainesville, FL 32610, USA;

20    [7]KwaZulu-Natal Research Innovation and Sequencing Platform (KRISP), School of

21    Laboratory Medicine and Medical Sciences, College of Health Sciences, University of

22    KwaZulu-Natal, Durban, South Africa; [8]Coordenação Geral dos Laboratórios de Saúde

23    Pública/Secretaria de Vigilância em Saúde, Ministério da Saúde, (CGLAB/SVS-MS)

24    Brasília, Distrito Federal 70719-040, Brazil; [9]Department of Molecular and Translational

25    Medicine, Section of Microbiology, University of Brescia, Brescia, Italy; [10]Unit of Clinical

26    Laboratory Science, University Campus Bio-Medico of Rome, Rome, Italy; [11]National

27    HIV/AIDS Researh Center, Istituto Superiore di Sanità, Rome, Italy; [12]Anthropology, and

28    Applied Ethics, Campus Bio-Medico University, Rome, Italy; [13]Department of Mental Health

29    and Public Medicine, University of Campania "Luigi Vanvitelli", Naples, Italy; [14]Department

30    of Biochemical Sciences "A. Rossi Fanelli", University of Rome "La Sapienza", Rome, Italy;

31    [15]ICBR, University of Florida, Gainesville, FL 32610, USA; [16] Department of Public Health

32    and Infectious Diseases, Policlinico Umberto I Università 'Sapienza', Rome; [17]Laboratory of

33    Molecular Virology, International Centre for Genetic Engineering and Biotechnology

34    (ICGEB), Trieste, Italy; [18]Department Infectious Diseases, – Fondazione Policlinico

35    Universitario "A. Gemelli" IRCCS, Rome, Italy; [19] Department of Epidemiology, College of

36    Public Health and Health Professions, University of Florida, Gainesville, FL 32610, USA; [20]

37    Department of Mathematics, University of Lafayette, LA, USA; [21]Department of Zoology,

38    University of Oxford, Oxford OX1 3PS, UK.

39

40        [*]These authors contributed equally to this article.

41    ^Correspondence and requests for materials should be addressed to

42    m.ciccozzi@unicampus.it; Dzella@ihv.umaryland.edu; salemi@pathology.ufl.edu

43

44    **Abstract**

45    We investigated SARS-CoV-2 transmission dynamics in Italy, one of the countries hit

46    hardest by the pandemic, using phylodynamic analysis of viral genetic and epidemiological

47    data. We observed the co-circulation of at least 13 different SARS-CoV-2 lineages over time,

48    which were linked to multiple importations and characterized by large transmission clusters

49    concomitant with a high number of infections. Subsequent implementation of a three-phase

50    nationwide lockdown strategy greatly reduced infection numbers and hospitalizations. Yet

51    we present evidence of sustained viral spread among sporadic clusters acting as "hidden

52    reservoirs" during summer 2020. Mathematical modelling shows that increased mobility

53    among residents eventually catalyzed the coalescence of such clusters, thus driving up the

54    number of infections and initiating a new epidemic wave. Our results suggest that the

55    efficacy of public health interventions is, ultimately, limited by the size and structure of

56    epidemic reservoirs, which may warrant prioritization during vaccine deployment.

57

58

59

60    **Keywords:** SARS-CoV-2; Italy; pandemic; genomic epidemiology.

61

62

63

64

65

66

67

**Main text**

69

70 On December 31st 2019, the World Health Organization (WHO) China Country Office was informed of pneumonia cases of unknown aetiology detected in Wuhan City, Hubei Province[1,2]. By January 11th – 12th 2020, Chinese authorities identified a novel single stranded, positive-sense enveloped RNA *Betacoronavirus*, with genome of 30,000 nucleotides in length, belonging to the *Coronaviridae* family, related to the severe acute respiratory syndrome coronavirus (SARS-CoV) that caused a global outbreak in 2002–2004 [3]. Initially named nCoV-2019 (novel Coronavirus 2019), the virus likely emerged from several recombination events in bats and pangolins [4], and was subsequently introduced in the human population through zoonotic transmissions [1,5]; it was later renamed SARS-CoV-2, and recognized as the etiologic agent of Coronavirus Disease 2019 (COVID-19) [6]. Epidemiological investigations and phylogenetic analysis promptly confirmed airborne SARS-CoV-2 human-to-human transmission [3,7]. Following its worldwide spread, the WHO declared the outbreak as a Public Health Emergency of International Concern on January 30th, 2020, and a pandemic on March 11th, 2020. As of December 16th, 2020, SARS-CoV-2 has spread to 216 countries with nearly 74 million confirmed cases and over 1.6 million fatalities [8].

86

87 Italy was one of the first and most affected countries in the world. By October 31st 2020, the Italian Ministry of Health and the Civil Protection Department reported 1.38 million total SARS-CoV-2-related cases, and 49,261 deaths [9]. The first confirmed imported cases dated back to January 30th 2020 when two tourists from Wuhan, China, were tested positive for SARS-CoV-2 in Rome (**Figure 1A**). On February 17th 2020, the Italian government confirmed the first locally acquired case in a small city in Northern Italy (Codogno, Lombardy region) [10]. Three days later, the first COVID-19-related death in Italy, a 78-year old male, was reported in the city of Padova. As the epidemic quickly spread throughout the country, establishing Italy as one of the major SARS-CoV-2 hotspots [11], the Italian government declared a Public Health Emergency of National Importance, enabling the introduction of restriction measures to limit new infections [12]. In the effort to flatten the epidemic curve, Phase I lockdown measures were first introduced on March 7th – 8th 2020 in

99   11 municipalities of Northern Italy, where most cases had occurred, and extended by March

100  11th to the whole country (**Figure 1A**). Described as the largest lockdown in the history of

101  Europe [13], citizen mobility was restricted, except for "well grounded" work- or health-related

102  reasons. A universal mask mandate was required at all times outdoors. Schools, university

103  activities, public/cultural events, and sport competitions were also suspended nationwide, as

104  well as non-essential commercial activities. Borders with other states were closed, and within

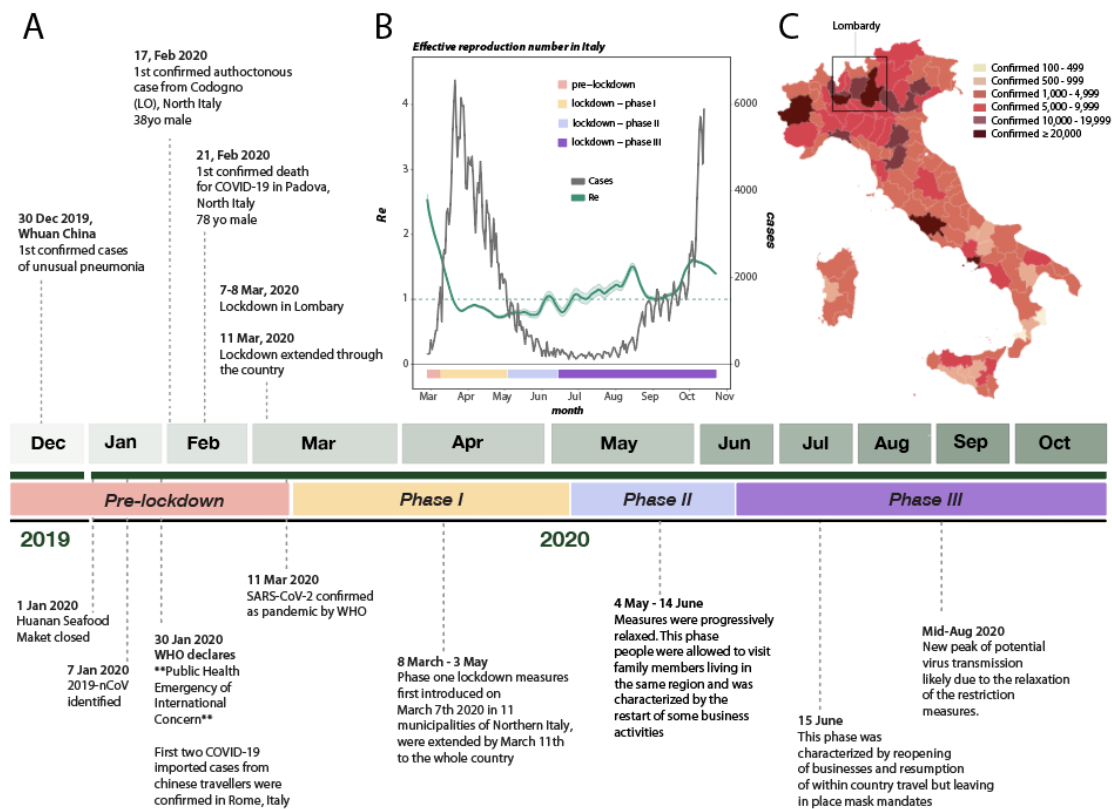105  the country public transport was limited or shut down.

106



107
108
109

110  **Figure 1. History of SARS-CoV-2 epidemic in Italy. A**) Timeline of key events following the first confirmed

111  cases of SARS-CoV-2 infection in Italy. **B**) Epidemic curve showing the progression of reported daily viral

112  infection numbers in Italy from the beginning of the epidemic in March (black) and changes in *Re* estimations in

113  the same period (green), with lockdown phases indicated along the bottom. **C**) Map of cumulative SARS-CoV-2

114  cases per 100,000 inhabitants in Italy up to Oct 2020.

115

116  As daily viral infection numbers decreased, public health measures were progressively

117  relaxed through a Phase II (May 4th), which allowed visits to family members living in the

118  same region and the restart of some business activities, and a Phase III (June 15th), which

119  allowed reopening of businesses and resumption of within country travel, but left in place

120  mask mandates and bans on large-scale meetings. A significant slowdown in the number of

121  infections since the beginning of May 2020 (**Figure 1B**) validated the effectiveness of Phase

122  I restrictions.

123

124  After a period of seemingly stable epidemic recession, with very few new cases detected

125  between June-August, a new epidemic wave hit the country, resulting in higher incidence

126  than before. Superimposition of the reported epidemic curve and dynamic estimates of the

127  effective reproduction number, $Re$, throughout the three major periods (first wave, recess,

128  second wave) of the Italian epidemic, revealed an interesting pattern (**Figure 1 B**). $Re$

129  provides a measure of the average number of secondary infections caused by a single infected

130  person: a growing epidemic is typically characterized by $Re > 1$, while $Re < 1$ indicates no

131  growth. As expected, $Re$ values were estimated to be $> 2$ at the beginning of SARS-CoV-2

132  exponential spread in Italy, and quickly fell to values $< 1$ after the start of Phase one

133  lockdown measures. Yet, between end of June and end of August, through Phase II and III

134  lockdowns, $Re$ values showed an oscillating behaviour, with progressively higher peaks ($>1$),

135  despite the consistently low number of newly detected infections. As infections and

136  hospitalizations began to climb in September, $Re$ temporarily decreased close to 1, to increase

137  again by mid-October, just before the beginning of the new exponential growth of infected

138  cases, currently ongoing. Indeed, by October 31[st], all Italian regions, albeit with different

139  rates were hit by the epidemic (**Figure 1C**). The rapid increase of COVID-19 patients

140  requiring hospitalization during the early months of 2020, as well as $Re$ oscillations during

141  the period of epidemic recession, suggest the virus was circulating cryptically among

142  undetected transmission clusters. During this time, there were possibly thousands of mild or

143  asymptomatic infections among undetected (hidden) reservoirs that preceeded each

144  exponential growth phase of each epidemic wave [14,15]. Indeed, dramatic resurgences in cases

145  after easing stringent public health interventions (i.e., stay-at-home orders) that temporarily

146  curtailed epidemic spread was also observed in several other European countries (e.g., UK,

147  France, and Germany, among others).

148

149       To investigate further, we coupled epidemiological data with phylodynamic analysis

150  of 714 viral sequences currently available from Italian patients, sampled between January

151  30[th] to October 1[st], 2020 (see Methods). Viral population dynamics were assessed using non-

152  parametric coalescent estimates of the effective population size ($Ne$) over time (a measure of

153    viral diversity representing the number of diverse genomes contributing to the next

154    generation), given a collection of plausible maximum likelihood (ML) evolutionary histories

155    inferred from viral sequence data. Although distinct patterns could be observed in $Ne$

156    estimates, all reconstructions agree on a rise in $Ne$ until the end of March 2020, matching the

157    rise in number of reported cases (**Supplementary Figure 1**). The best-fit model (i.e. the

158    collection of trees with the highest likelihood, $log$L > -49,120) also depicts a steady,

159    continuous decline in $Ne$ until October (**Supplementary Figure 1**, pattern A), possibly

160    reflecting the impact of lockdown measures on the viral population. As $Ne$ is related to viral

161    genetic diversity, this pattern may indicate that, despite the rapid rise of cases in late summer,

162    the viral population maintained lower diversity relative to the earlier months of the epidemic.

163    This is consistent with a reduction of viral importations, likely resulting from global public

164    health intereventions such as travel bans. Two alternative patterns inferred from trees similar

165    in likelihood value to the best-fit model, show either a similar downward trend followed by a

166    pronounced increase in $Ne$ between September and October (pattern B), or a slower but

167    steady increase in $Ne$ between April and October (pattern C). Both reconstructions are in

168    agreement with an increase of viral $Ne,$ corresponding with an exponential increase of SARS-

169    CoV-2 infections during the second epidemic wave. Together, with the inferred oscillations

170    of $Re$ values following the first epidemic wave, the analyses suggest the persistence of

171    complex transmission dynamics throughout the epidemic recession period, involving

172    undetected asymptomatic or mildly affected individuals. Even considering the $Ne$ values

173    inferred from ML trees with lower likelihood, we arrive at an analogous conclusion – overall

174    reduction in $Ne$ after April but repeated fluctuations throughout recession and second

175    epidemic wave (**Supplementary Figure 1**, black curves).

176

177        Longitudinal comparison of SARS-CoV-2 dissemination patterns over time among

178    different Italian regions shows that the pre-lockdown phase was characterized by an

179    exponential growth of the number of daily-confirmed COVID-19 cases and deaths, with

180    highest incidence in the Northwest, followed by a significant decrease across all regions in

181    the aftermath of lockdown measures (**Supplementary Figure 2**). By the end of August 2020,

182    epidemiological data also show increased and sustained transmission in the South and Insular

183    regions, possibly driven by interregional spreading through small family/social network

184    clusters. These regions are the main touristic destination for Italians, and most of the

185    restrictions on international travel were still in place during Phase III [16]. Lineages proportion

186    and regional-specific distribution in different parts of the country are indicative of several

187    independent founder events (**Figure 2B**). For example, lineage A, predominant in Sicily, has

188    been detected in epidemiologically linked transmission chains that appear to be related to

189    immigrants arrived from North Africa during the late Phase III [17]. Interestingly, the number

190    of circulating lineages have changed over time (**Supplementary Figure 3**). Sub lineage B.2

191    was the first one identified in January, marking the primary introduction of imported cases

192    from China (no shown in Figure 2). Between February and April, additional sub lineages,

193    such as

194



195

196    **Figure 2.** Frequency and distribution of SARS-CoV-2 lineages and sub lineages in Italy. **A**) Frequency of the

197    lineages and sub lineages of SARS-CoV-2 among Italian macro regions. **B**) Distribution of the most prevalent

198    lineage and sub lineage across the country.

199

200    B.1, B.1.1, and B.1.5 emerged in Northern and Central Italy, the epicenter of the first

201    epidemic wave, likely reflecting subsequent importations. At the beginning of Phase II

202    lockdown in May, which followed a dramatic decrease in cases, only B.1. and B.1.1 sub

203    lineages were detected. During the period of epidemic recession between June and July,

204    multiple sub lineages co-circulated again. However, the subsequent second wave was

205    dominated by B.1.1 (September) and B.1. (October) (**Supplementary Figure 3**). Since Phase

206    II and III measures permitted intra- and then inter-regional travel, respectively, while

207    country borders remained mostly closed (except with European countries part of the Shengen

208    agreement), it is plausible that in the first epidemic wave lineages' heterogeneity resulted

209    from intial founder events associated with international travel, and then propagated through

210    within state mobility during epidemic recession. Such sequence-based inferences, however,

211    should be interpreted with caution because of the inherent sampling bias in SARS-CoV-2

212    full-length genomes currently available from Italian patients, which could affect results and

213    limit their generalizability [18].

214

215       In our sequence dataset, only Lombardy (Northwest, most affected region so far), has

216 provided a robust number of viral genomes (n=405), which in turn corresponds

217 approximatively to one genome available every 450 positive cases. Abruzzo in Central Italy

218 is the second most represented region in terms of available genomes (n=87), while many

219 other regions, including Liguria (Northwest), Umbria (Central) and Calabria (South) are not

220 comprehensively represented (**Figure 3A**), thus affecting our ability to characterize in-depth

221 SARS-CoV-2 molecular epidemiology at a regional level. Nevertheless, phylogeny-inferred

222



223

224 **Figure 3. Phylogenetic characterization of Italian SARS-CoV-2 sequences.** A) Map of Italy showing the

225 number of SARS-CoV-2 genome sequences by region. The size of the circles indicates the number of new

226 genomes available since the beginning of the epidemic in Italy; B) Time-resolved maximum likelihood tree of

227 1421 SARS-CoV-2 sequences including 714 from Italy (red circles). C) Chord diagram of estimated numbers of

228 migration flows between the geographic areas. D) Frequency of estimated geographical origins for identified

229 transmission clusters involving Italy and originating in the months of January through October of 2020.

230 E) Frequency of Italian sequences (sampled from January through October) classified as unclustered (grey) or

231 belonging to clusters with Italian (white) or non-Italian origins (black).

232

233 virus evolutionary patterns are useful to corroborate epidemiological data, test hypotheses

234 regarding factors driving epidemic dynamics, and assess public health interventions such as

235 stay-at-home orders. To this end, we time-scaled the best 100 ML trees of all available

236 SARS-CoV-2 full genomes from Italian patients, and inferred the most likely location of each

237 internal node (ancestral sequence) in the trees (see Methods for details). The overall

238 topologies of the inferred trees were highly similar, and linear regression of root-to-tip

239 genetic distances against sampling dates indicated sufficient temporal signal in the sequence

240     data (**Supplementary Figure 4**). Although SARS-CoV-2 evolutionary rate in Italy was

241     somewhat lower (1.44 $10^{-04}$ nucleotide substitutions/site/year) than values obtained for the

242     worldwide epidemic [19,20], the most recent common ancestor (TMRCA) of the available Italian

243     sequences, ranged between January 2nd and January 26th (mean Jan 14th) 2020, consistent

244     with the date of the first confirmed case (Jan 30th). Similarly, the root node (origin) of a time-

245     scaled ML tree including both Italian (n=714) and worldwide reference sequences (n=1,421)

246     was placed in China (99.8% probability), with a TMRCA dating back to early December

247     2019, in agreement with available epidemiology data [21,22], further validating our phylogeny

248     inference. The tree (**Figure 3B**) consistently shows most of the Italian sequences interspersed

249     with virus strains collected in other countries. This pattern, alike the one observed elsewhere

250     [23], confirms that emergence of SARS-CoV-2 strains during the first epidemic wave was

251     primarily fostered by travel exposure during the pre-lockdown phase, rather than

252     interregional spreading. According to the estimation of migration flows, we further examined

253     the potential Italian role as an exporter of SARS-CoV-2. The number of state transitions into

254     and from Italy (**Figure 3 panel C**) heavily relies on the number and nature of the sequences

255     that are included from other locations. Independently of the dataset, and in line with the

256     epidemiological information, most of the geographical sources of the introductions are

257     attributed to Europe  (**Figure 3C**). Well supported (bootstrap values > 90%) putative

258     transmission clusters within the phylogeny were identified based on a pre-defined genetic

259     distance threshold likely to detect epidemiologically linked sequences (see Methods).

260     Clusters containing at least one Italian sequence were considered of interest for the estimation

261     of temporal and spatial origins of the transmission. Temporal origins of each cluster were

262     derived from the clock-estimated age of the MRCA of all sequences belonging to the cluster.

263     Spatial origins were inferred using joint likelihood ancestral state reconstruction, given

264     known country of sampling of tip nodes (sampled sequences) within the tree. As expected,

265     the number of (well supported) clusters formed over the course of the epidemic was largely

266     influenced by the number of contemporaneous samples (**Figure 3D**), limiting conclusions

267     regarding the rate of cluster formation over time. The estimated geographic origins of each

268     cluster reflected the distribution of samples among the reference sequences, largely limited to

269     Europe and North America. However, after April, transmission clusters could only be traced

270     back to Italy, suggesting highly localized transmission following the implementation of Phase

271     one lockdown measures. Each Italian sequence was then classified either as unclustered (i.e.

272     no cluster with any other sequence with bootstrap >90%), or belonging to a local (all Italian)

273     cluster. Italian sequences within well supported clusters including and originating from non-

274  Italian strains were classified as belonging to "outside" clusters. Finally, each well supported
275  cluster for which a single country could not be assigned with >90% probability as the one at
276  the origin of that cluster, was also considered to be an outside (albeit unknown in origin)
277  cluster. This revealed distinct patterns between January, February-July, and August (**Figure**
278  **3E**). All Italian sequences obtained in January belonged to clusters of foreign origin,
279  demonstrating the influence of outside introductions before lockdowns were put into place.
280  The predominant fraction was quickly replaced by sequences belonging to clusters of local
281  origin and unclustered sequences, which suggests potential undersampling. The month of
282  September, when the second epidemic wave was increasing, sequences of local origin, with
283  no sequences of foreign origin, largely dominated. The fraction of sequences sampled in
284  August (75%) was outside the 95% confidence interval (~50%) for the fraction in remaining
285  months, emphasizing the significant contribution of local transmission on sequences sampled
286  in September. However, the specific mutational profile of the Italian sequences (**Figure 4A**),
287  relative to the Wuhan reference (NC_045512), also provided some evidence of recently
288  imported strains during Phase III lockdown, when travel bans began to be eased. In
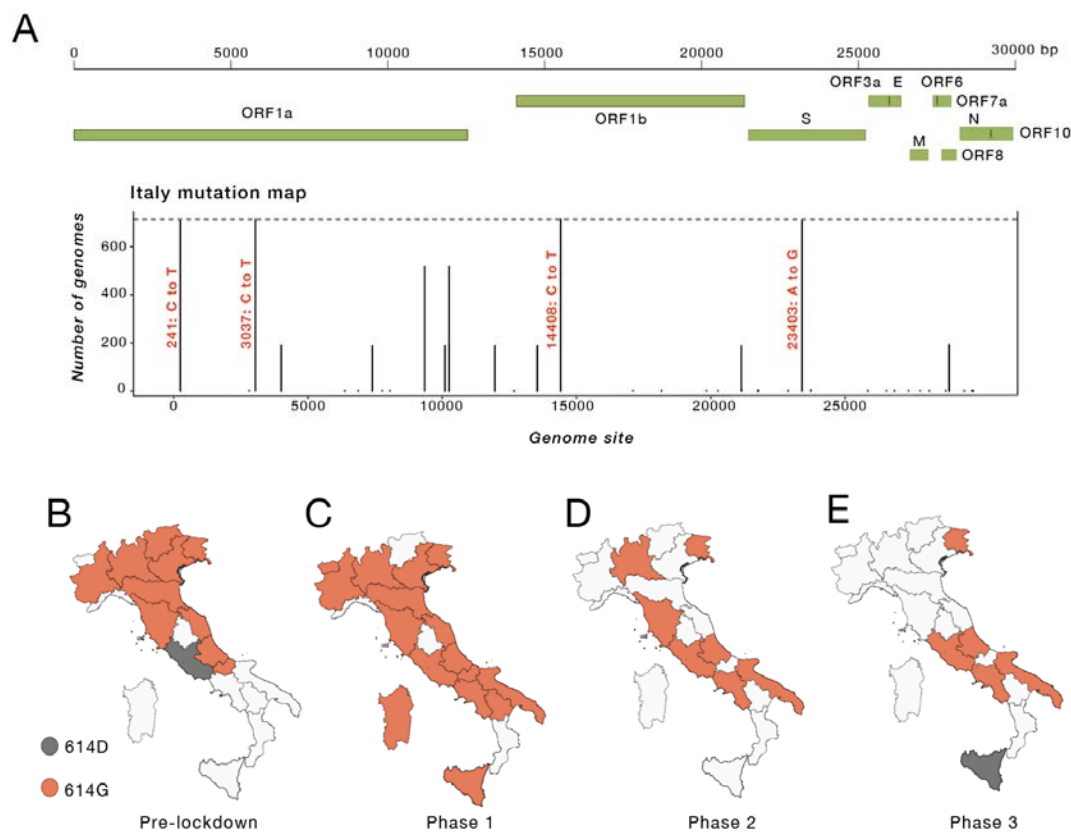289  particular, 97.34% (n=695) of the available sequences carried the



290

291

292 **Figure 4. Italian strains mutations pattern.** A) Variant maps of the most common mutations mapped against
293 the SARS-CoV-2 genomes. Most common mutations defined as mutations present in >90% of the genomes in
294 that group (black lines). B-E) Change in frequency of the D614G mutation in the Spike Protein across Italian
295 regions during epidemic phases.
296

297

298 mutation encoding for the amino acid change D614G (genomic coordinate: *23403A>G*) in

299 the Spike protein of SARS-CoV-2, while the remaining 2.66% (n=19) sequences displayed

300 the nucleotide sequence encoding for the D614 wild type. The D614G mutation has been

301 associated with higher infectivity and greater transmissibility with no effects on disease

302 severity outcomes [24–26], although some of these findings have recently been questioned [27].

303 The frequency of the D614G polymorphism among Italian regions over time (**Figure 4B**)

304 shows that while G614 quickly became rapidly dominant during the first epidemic wave, and

305 remained the only one detected in the available sequences through the first two lockdown

306 phases. Thereafter, the D614 variants re-emerged following the relaxation of Phase III

307 measures in Sicily (Insular Italy), possibly due to the epidemiologically-linked transmission

308 chains related to immigration flow from North Africa [17], a scenario reinforced by SARS-

309 CoV-2 lineage A prevalence in that region (**Figure 2B**).

310

311 Results of cluster analysis indicate maintained local transmissions fostered by

312 relatively small transmission chains during the months of low case reports and through the

313 beginning of the second wave. This observation suggests that epidemic resurgence was

314 associated with a relaxation of lockdown measures that led to increased local transmission,

315 rather than a large number of virus re-introductions into the country. Such a scenario is also

316 supported by surveys showing a significant reduction in the number of foreign tourists (about

317 -65.9%), but an increase, albeit small (1.1%), of domestic tourism during the summer season

318 after restrictions on interregional travel were relaxed [16]. In order to explore whether

319 increased mobility could explain the second surge of cases in Italy, we carried out stochastic

320 agent-based epidemic simulations. Mobility data across three different modes of

321 transportation (walking, public, and personal vehicle), derived from Apple Mobility trends

322 reports, were used as a proxy for the number of individuals with whom an infected individual

323 comes into contact, which was allowed to vary over time (see Methods). As the number of

324 hospitalizations also dropped drastically (and stayed low) following the first surge in cases,

325 the role of removal of infected individuals from the population *via* hospitalization was also

326 tested, by allowing the probability of an infected individual exiting the simulation to be

327    proportional to the standardized hospitalization rates (also varying in time). The simulated

328    number of active infections over time using the mobility data alone, hospitalization data

329    alone, and combined were then compared to the empirical case data. Whereas all three

330    models produced a similar rate of new infections during the first epidemic wave

331    (**Supplementary Figure 5**), time-varying rate of removal based on hospitalization rates

332    (without mobility data) produced a continuing exponential growth of infections (**Figure 5A**).

333    As in empirical data, time-varying number of contacts based on mobility produced two

334    distinct waves, which were the most similar to the epidemic curve (**Figure 5B**). The model

335    incorporating both mobility and hospitalization rates produced a first wave that was of too

336    large a magnitude and a second wave too early in origin than the previous model (**Figure**

337    **5C**). The model incorporating mobility data alone resulted in the lowest mean absolute error

338    (**Figure 5D**), producing a first wave of similar timing and magnitude and a delayed second

339    wave, closer to the empirical epidemiology data (**Figure 1B**).
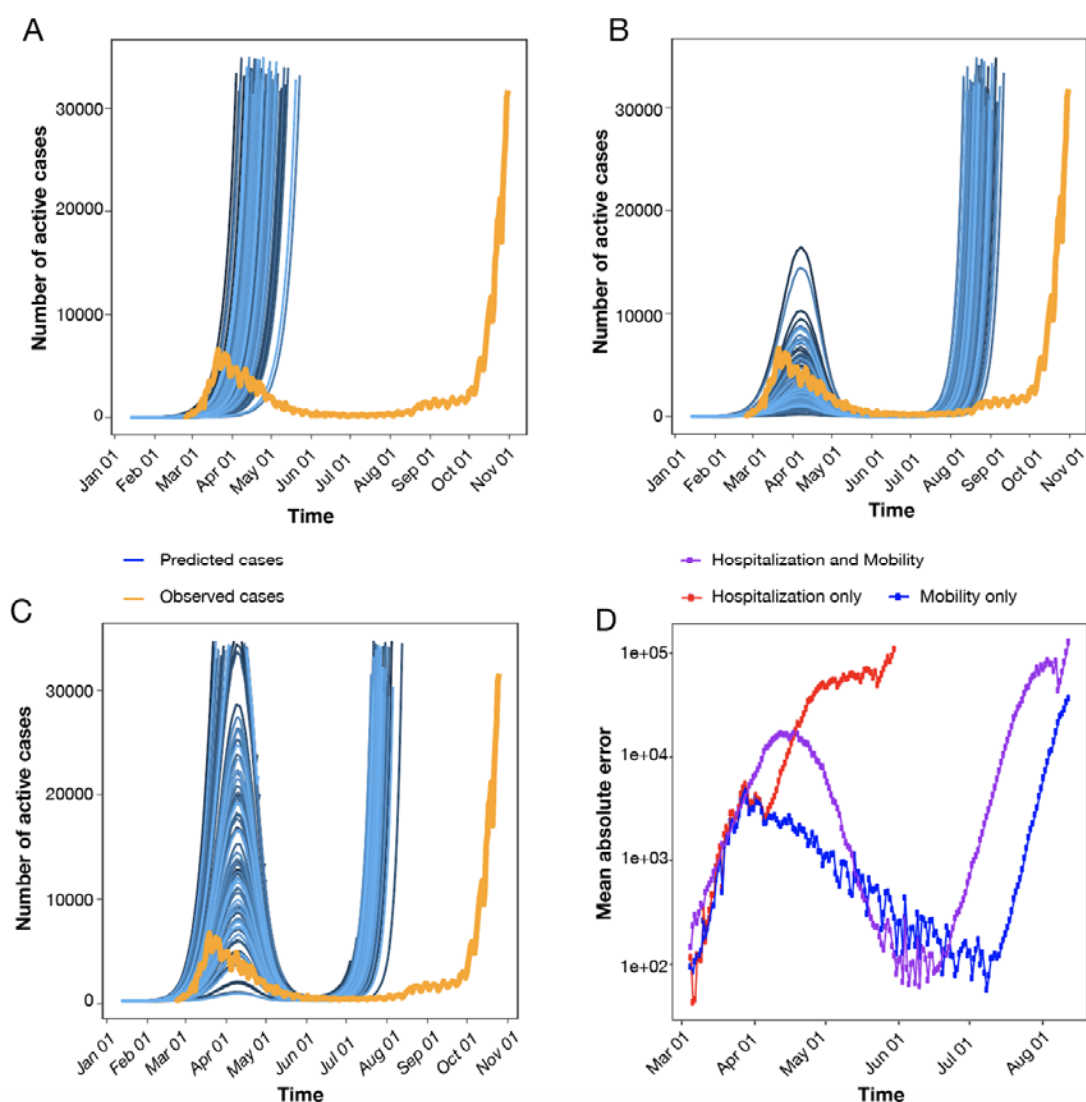
340

12

341

**Figure 5. Simulated epidemics under scenarios involving time-varying mobility and hospitalization rates.** A) The probability of removal of an infected individual was proportional to the empirical rate of hospitalization in the simulation of active infected cases over one year (blue). B) The number of individuals with whom an infected individual comes in contact was proportional to the empirically determined number of individuals utilizing walking, as well as public and personal modes of transportation, as primary means of mobility in the simulation of active infected cases over one year (blue). C) Hospitalization rates (as in A) and mobility data (as in B) were combined in the simulation of active infected cases over one year (blue). In A-C, orange represents the number of empirically observed infections. D) The absolute error was calculated for each time point of collected observations between simulated infections using hospitalization rates only (red), mobility data only (blue), and the combination (purple). Mean absolute error was calculated as the averaged error across 1,000 simulations. Simulation assumed a single outside introduction.

353

354 While results indicate that mobility data could reproduce epidemic wave patterns in Italy, we

355 cannot exclude additional factors, that might have played a role delaying the second epidemic

356 wave, not fully captured by our simulations, such as retention of restriction measures,

357 different "types" of mobility between first and second wave, or higher temperatures in the

358 summer [28].

359   By coupling phylodynamic analysis of viral genetic and epidemiology data, we show
360 how the interplay between public health intervention and shifting SARS-CoV-2 transmission
361 dynamics in Italy may explain the oscillation between times of relatively stable epidemic
362 recession and dramatic resurgences, as it is currently being observed. This pattern of
363 "rubberbanding" or "snapping back" after public health restrictions are lifted has,
364 unfortunately, been followed by several other European countries. Overall, we show the
365 critical role played by small transmission clusters, acting as "hidden reservoirs" during
366 epidemic recession following aggressive lockdown measures, in maintaining SARS-CoV-2
367 low-level circulation in Italy, which eventually seeded a new epidemic wave. Despite the
368 consistent agreement between different viral phylogeny-based and epidemiology data
369 analyses, however, limitations of our work need to be acknowledged. Availability of a large
370 number of viral sequences, collected over an extended period of time and sufficiently
371 representative of the ongoing epidemic, is crucial for prompt genomic surveillance, and the
372 evaluation and planning of effective and opportune control strategies. The number of Italian
373 SARS-CoV-2 full genomes currently deposited in public databases represents a very small
374 fraction (0.05%) of the documented number of confirmed cases in Italy, and sampling bias
375 across regions differently affected by the epidemic further limits generalizability of the
376 results. Moreover, our definition of putative transmission clusters (see Methods) does not
377 require the sampling and inclusion of all the strains involved in a transmission chain,
378 although it does allow for detection of monophyletic clades that likely comprise sequences
379 epidemiologically linked through a transmission chain, whilebeit not directly. Nevertheless,
380 epidemiology observations, corroborated by phylodynamic analyses based on available
381 sequences, depicted a coherent picture. The first epidemic wave in Italy appears to have
382 largely been linked to outside introductions leading to large transmission clusters,
383 concomitant with high number of infections. Subsequent implementation of a three-phase
384 nationwide lockdown strategy greatly mitigated numbers of infection and
385 hospitalization during summer 2020. Yet, once mobility increased and social distancing
386 decreased due to the progressive easing of lockdown measures, a sudden spike of infectious
387 cases was observed, promptly followed by new hospitalizations. Our agent-based
388 mathematical model recapitulates this phenomenon, further supporting the hypothesis that the
389 small clusters observed during the summertime were acting, essentially, as "hidden
390 reservoirs" that likely merged following the increased in mobility and reduction of social
391 distancing measures. This in turn provided the "spark" for the sudden increase of infections
392 observed at the end of summer, which led to the subsequent second wave of exponential

393　grow. In other words, the drivers of SARS-CoV-2 transmission dynamics shifted from high

394　levels of community transmission, likely involving mass super spreader events, in the early

395　Italian epidemic, to sustainment by smaller family/social network clusters later in the

396　epidemic. Unfortunately, this also suggests that no amount of community level interventions

397　may be sufficient to curb the epidemic as long as people do not adhere to individual level

398　measures such as mask use, hand hygiene, and social distancing. New lockdown measures

399　are likely to provide only temporary relief, as has already happened in the first months of the

400　epidemic in Italy and many other countries. Indeed, an important debate is currently ongoing

401　about vaccine deployment, given financial and logistic restrictions mandating a very long

402　phased deployment, based on prioritization policies. In this context, our results suggest that

403　hidden transmission reservoirs may continue to sustain local outbreaks into late 2021, as

404　vaccine rollout will likely take months before reaching the necessary herd-immunity

405　threshold. Ultimately, our ability to curb successfully the current pandemic, may be linked to

406　our ability to determine number and structure of such reservoirs within the social and

407　behavioral context of specific locales.

408

409　**Acknowledgments**

417

418　**Authors' contributions**

419　Conception and design: MG, EC, BRM, MS and MC; Data collection: MG, EC;

420　Investigations: MG, EC, FB, BRM, VF, SF, and JL; Data Analysis: MG, EC, BRM, EW, VF,

421　NED, CB and JL; Writing – Original: MG, EC, FB, BRM, MS, DZ and MC; Draft

422　Preparation: MG, EC, FB, BRM, JL, VT, MS, DZ and MC; Revision: GC, SP, AR, AB, VT,

423　CS, AM, TA, EW, TdO, LCJA, GC, RC, AC, JL, MS, DZ, and MC.

424

425　**Competing Interests Statement**

426　The authors declare no competing interests.

427

428 **REFERENCES**

429 1.  Andersen, K. G., Rambaut, A., Lipkin, W. I., Holmes, E. C. & Garry, R. F. The proximal

430      origin of SARS-CoV-2. *Nature Medicine* vol. 26 450–452 (2020).

431 2.  Yang, X.-L. *et al.* Isolation and Characterization of a Novel Bat Coronavirus Closely

432      Related to the Direct Progenitor of Severe Acute Respiratory Syndrome Coronavirus. *J.*

433      *Virol.* **90**, 3253–3256 (2016).

434 3.  Wu, F. *et al.* A new coronavirus associated with human respiratory disease in China.

435      *Nature* **579**, 265–269 (2020).

436 4.  Massimiliano S. Tagliamonte *et al.* Recombination and purifying selection preserves

437      covariant movements of mosaic SARS-CoV-2 protein S | bioRxiv. (2020)

438      doi:http://doi.org/10.1101/2020.03.30.015685.

439 5.  Zhou, P. *et al.* A pneumonia outbreak associated with a new coronavirus of probable bat

440      origin. *Nature* **579**, 270–273 (2020).

441 6.  Giovanetti, M., Angeletti, S., Benvenuto, D. & Ciccozzi, M. A doubt of multiple

442      introduction of SARS-CoV-2 in Italy: A preliminary overview. *J. Med. Virol.* **92**, 1634–

443      1636 (2020).

444 7.  The 2019-new coronavirus epidemic: Evidence for virus evolution - Benvenuto - 2020 -

445      Journal of Medical Virology - Wiley Online Library.

446      https://onlinelibrary.wiley.com/doi/full/10.1002/jmv.25688.

447 8.  World Health Organization (WHO). Coronavirus disease (COVID-19) Situation Report –

448      103, 02 May 2020 [Internet]. 2020. Available from: https://www.who.int/docs/default-

449      source/coronaviruse/situation-reports/20200502-covid-19-sitrep-

450      103.pdf?sfvrsn=d95e76d8_4.

451     9.    Civil Protection Department. COVID-19 Italia.

452           http://opendatadpc.maps.arcgis.com/apps/opsdashboard/index.html#/b0c68bce2cce4

453           78eaac82fe38d4138b1.

454     10.   La Repubblica. Coronavirus, i contagi nel Lodigiano sono 15: i primi sono un 38enne di

455           Codogno e sua moglie. In isolamento 250 persone.

456           https://milano.repubblica.it/cronaca/2020/02/21/news/coronavirus_a_milano_contagg

457           iato_38enne_e_un_italiano_ricoverato_a_codogno-249121707/.

458     11.   Percivalle, E. *et al.* Prevalence of SARS-CoV-2 specific neutralising antibodies in blood

459           donors from the Lodi Red Zone in Lombardy, Italy, as at 06 April 2020. *Eurosurveillance*

460           **25**, (2020).

461     12.   Civil Protection Department. State of epidemiological emergency.

462           http://www.protezionecivile.gov.it/media-comunicazione/news/dettaglio/-

463           /asset_publisher/default/content/coronavirus-dichiarato-lo-stato-di-emergenza.

464     13.   Jason Horowitz & Emma Bubola. On Day 1 of Lockdown, Italian Officials Urge Citizens to

465           Abide by Rules. *NYTimes.com* (2020).

466     14.   The first two cases of 2019-nCoV in Italy: Where they come from? - Giovanetti - 2020 -

467           Journal of Medical Virology - Wiley Online Library.

468           https://onlinelibrary.wiley.com/doi/full/10.1002/jmv.25699.

469     15.   Stefanelli, P. *et al.* Whole genome and phylogenetic analysis of two SARS-CoV-2 strains

470           isolated in Italy in January and February 2020: additional clues on multiple introductions

471           and further circulation in Europe. *Eurosurveillance* **25**, 2000305 (2020).

472     16.   CST Firenze for Assoturismo Confesercenti. *Turismo Estate 2020 Italia: mancano gli*

473           *stranieri, calo della domanda del -30,4%.*

474     http://centrostudituristicifirenze.it/blog/turismo-estate-2020-italia-mancano-stranieri-

475     calo-della-domanda/.

476   17. Irene Schöfberger & Marzia Rango. Migration in West and North Africa and across the

477     Mediterranean - COVID-19 and migration in West and North Africa and across the

478     Mediterranean - | IOM Online Bookstore. (2020).

479   18. Mavian, C., Marini, S., Prosperi, M. & Salemi, M. A Snapshot of SARS-CoV-2 Genome

480     Availability up to April 2020 and its Implications: Data Analysis. *JMIR Public Health*

481     *Surveill.* **6**, (2020).

482   19. Andrew Rambaut. Phylodynamic Analysis | 176 genomes | 6 Mar 2020. *Virological*

483     https://virological.org/t/phylodynamic-analysis-176-genomes-6-mar-2020/356 (2020).

484   20. Su, Y. C. *et al.* Discovery of a 382-nt deletion during the early evolution of SARS-CoV-2.

485     *bioRxiv* 2020.03.11.987222 (2020) doi:10.1101/2020.03.11.987222.

486   21. Lu, J. *et al.* Genomic Epidemiology of SARS-CoV-2 in Guangdong Province, China. *Cell*

487     **181**, 997-1003.e9 (2020).

488   22. Duchene, S. *et al.* Temporal signal and the phylodynamic threshold of SARS-CoV-2. *Virus*

489     *Evol.* **6**, (2020).

490   23. Deng, X. *et al.* Genomic surveillance reveals multiple introductions of SARS-CoV-2 into

491     Northern California. *Science* **369**, 582 (2020).

492   24. Korber, B. *et al.* Tracking Changes in SARS-CoV-2 Spike: Evidence that D614G Increases

493     Infectivity of the COVID-19 Virus. *Cell* **182**, 812-827.e19 (2020).

494   25. Plante, J. A. *et al.* Spike mutation D614G alters SARS-CoV-2 fitness. *Nature* 1–6 (2020)

495     doi:10.1038/s41586-020-2895-3.

496   26. Benvenuto, D. *et al.* Evidence for mutations in SARS-CoV-2 Italian isolates potentially

497     affecting virus transmission. *J. Med. Virol.* **92**, 2232–2237 (2020).

498    27. van Dorp, L. *et al.* Emergence of genomic diversity and recurrent mutations in SARS-

499        CoV-2. *Infect. Genet. Evol.* **83**, 104351 (2020).

500    28. Benedetti, F. *et al.* SARS-CoV-2: March toward adaptation. *J. Med. Virol.* (2020)

501        doi:10.1002/jmv.26233.

502    29. Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local alignment

503        search tool. *J. Mol. Biol.* **215**, 403–410 (1990).

504    30. Müller, H. *et al.* VCF.Filter: interactive prioritization of disease-linked genetic variants

505        from sequencing data. *Nucleic Acids Res.* **45**, W567–W572 (2017).

506    31. Nakamura, T., Yamada, K. D., Tomii, K. & Katoh, K. Parallelization of MAFFT for large-

507        scale multiple sequence alignments. *Bioinformatics* **34**, 2490–2492 (2018).

508    32. Larsson, A. AliView: a fast and lightweight alignment viewer and editor for large

509        datasets. *Bioinformatics* **30**, 3276–3278 (2014).

510    33. O'Toole Á & McCrone J. Phylogenetic Assignment of Named Global Outbreak LINeages.

511        https://github.com/hCoV-2019/pangolin (2020).

512    34. Nguyen, L.-T., Schmidt, H. A., von Haeseler, A. & Minh, B. Q. IQ-TREE: A Fast and

513        Effective Stochastic Algorithm for Estimating Maximum-Likelihood Phylogenies. *Mol.*

514        *Biol. Evol.* **32**, 268–274 (2015).

515    35. Volz, E. M. & Frost, S. D. W. Scalable relaxed clock phylogenetic dating. *Virus Evol.* **3**,

516        (2017).

517    36. R Core Team. R: A language and environment for statistical computing. in (2019).

518    37. Volz, E. M. & Didelot, X. Modeling the Growth and Decline of Pathogen Effective

519        Population Size Provides Insight into Epidemic Dynamics and Drivers of Antimicrobial

520        Resistance. *Syst. Biol.* **67**, 719–728 (2018).

19

521    38. Prosperi, M. C. F. *et al.* A novel methodology for large-scale phylogeny partition. *Nat.*

522        *Commun.* **2**, 321 (2011).

523    39. Liam Revell. phytools: an R package for phylogenetic comparative biology (and other

524        things). *Methods Ecol. Evol.* (2012).

525    40. Pupko, T., Pe'er, I., Shamir, R. & Graur, D. A fast algorithm for joint reconstruction of

526        ancestral amino acid sequences. *Mol. Biol. Evol.* **17**, 890–896 (2000).

527    41. Paradis, E. & Schliep, K. ape 5.0: an environment for modern phylogenetics and

528        evolutionary analyses in R. *Bioinformatics* **35**, 526–528 (2019).

529    42. Wickham, H., François, R., Henry, L., Müller, K. & RStudio. *dplyr: A Grammar of Data*

530        *Manipulation.* (2020).

531    43. Henry, L., Wickham, H. & RStudio. *purrr: Functional Programming Tools.* (2020).

532    44. Kun Ren. rlist: A Toolbox for Non-Tabular Data Manipulation version 0.4.6.1 from CRAN.

533        https://rdrr.io/cran/rlist/.

534    45. Yu, G., Jones, B. & Arendsee, Z. *tidytree: A Tidy Tool for Phylogenetic Tree Data*

535        *Manipulation.* (2020).

536    46. Wickham H. *ggplot2: Elegant Graphics for Data Analysis.* (2016).

537    47. Dowle, M. *et al. data.table: Extension of 'data.frame'.* (2020).

538    48. Wickham, H. Reshaping Data with the reshape Package. *J. Stat. Softw.* **21**, 1–20 (2007).

539    49. Grolemund, G. & Wickham, H. Dates and Times Made Easy with lubridate. *J. Stat. Softw.*

540        **40**, 1–25 (2011).

541    50. Yu, G., Smith, D. K., Zhu, H., Guan, Y. & Lam, T. T.-Y. ggtree: an r package for visualization

542        and annotation of phylogenetic trees with their covariates and other associated data.

543        *Methods Ecol. Evol.* **8**, 28–36 (2017).

544    51. Wickham, H. & RStudio. *tidyr: Tidy Messy Data.* (2020).

545   52. Cori, A., Ferguson, N. M., Fraser, C. & Cauchemez, S. A New Framework and Software to

546        Estimate Time-Varying Reproduction Numbers During Epidemics. *Am. J. Epidemiol.* **178**,

547        1505–1512 (2013).

548   53. Jérémie Scire *et al. A method to monitor the effective reproductive number of SARS-CoV-*

549        *2.*

550   54. Gostic, K. M. *et al.* Practical considerations for measuring the effective reproductive

551        number, Rt. *medRxiv* (2020) doi:10.1101/2020.06.18.20134858.

552   55. Goldstein, E. *et al.* Reconstructing influenza incidence by deconvolution of daily

553        mortality time series. *Proc. Natl. Acad. Sci. U. S. A.* **106**, 21825–21829 (2009).

554   56. Lequime, S., Bastide, P., Dellicour, S., Lemey, P. & Baele, G. nosoi: A stochastic agent-

555        based transmission chain simulation framework in r. *Methods Ecol. Evol.* **11**, 1002–1007

556        (2020).

557
558
559
560
561
562
563
564
565   **Methods**
566
567   *Sequence data collection*
568   To perform a comprehensive analysis of the genomic epidemiology of SARS-CoV-2 in Italy,
569   after excluding low-quality genomes (> 10% of ambiguous positions), we dowloaded all
570   Italian full-length viral genomes available on GISAID (https://www.gisaid.org/) (n=714) up
571   to October 31[h] 2020. Appropriate acknowledgement was given to the sequencing laboratories
572   (**Supplementary Data S1**). Sampling locations of available genomes in this dataset included
573   17 of 20 regions in Italy, and collection dates spanned from January 30[th] (the first two
574   imported cases in Italy) to October 1[th] 2020. Each Italian sequence was used in a local

21

575    alignment (BLAST) [29] search for the most (genetically) similar non-Italian sequence in the

576    GISAID database as of Oct 31st, 2020, and linked to two reference sequences including the

577    best match (highest E-value) with a date occurring within one month following, as well as

578    one month prior to the sampling date of the Italian sequence (although, in some cases, only a

579    single non-Italian reference sequence fulfilling one of the inclusion criteria could be found

580    for multiple Italian query sequences). After removing duplicate sequences and masking

581    mutations potentially associated with common sequencing errors, using a vcf filter [30], a final

582    dataset of 1,421 reference sequences was assembled (**Supplementary Table 2**).

583

584    *Sequence alignments and phylogenetic analysis*

585    Sequences (Italian + reference strains) were aligned using MAFFT (FF-NS-2 algorithm)

586    employing default parameters [31]. The alignment was manually curated to optimize number

587    and location of gaps using Aliview [32]. A site-specific mutational comparison of the 714

588    Italian genomic sequences obtained from the GISAID database was made with the MAFFT-

589    aligned SARS-CoV-2 reference genome (RefSeq: NC_045512.2), obtained from the

590    GenBank database. Lineage assessment was conducted using the Phylogenetic Assignment of

591    Named Global Outbreak LINeages tool available at https://github.com/hCoV-2019/pangolin

592    [33]. Phylogenetic analysis of was performed using the maximum likelihood (ML) method

593    implemented in IQ-TREE (version 1.6.10), employing the best-fit model of nucleotide

594    substitution according to the Bayesian Information Criterion (BIC), as indicated by the Model

595    Finder application implemented in IQ-TREE [34]. The statistical robustness of individual nodes

596    was determined using 1000 bootstrap replicates.

597

598    *Molecular clock calibration and estimation of virus effective population size*

599    ML trees were inspected in TempEst v1.5.3 for the presence of temporal signal (i.e., linear

600    relationship between genetic distance and sampling time in the available sequences) [21]. The

601    treedater package in R v3.6.0 [35,36] was used for molecular clock calibration of the Italy-only

602    data, as well as the combined Italy and reference data. The top 100 maximum likelihood

603    (ML) trees (i.e. the trees with the 100 lowest *-log*[likelihood] values), were chosen for

604    calibration according to a strict clock (no branch specificity) among the Italy-only data,

605    whereas a single-best ML tree was chosen for the combined dataset. Individual *taxa* sampling

606    times were used to rescale branch lengths to time in each tree using a starting value of $8 \times 10^{-4}$

607    substitutions/site/year. The skygrowth non-parametric demographic model [37] was then used

608    in R with time-scaled trees to estimate of median virus effective population size (*Ne*) and

609 95% high posterior density intervals for each week during the epidemic in Italy (Italy-only

610 dataset) using the default smoothing parameter value (tau) of 0.1.

611

612 *SARS-CoV-2 transmission cluster identification and characterization*

613 Transmission clusters were identified using Phylopart v2 [38] applied to the ML tree of

614 combined sequence data (scaled in substitutions/site). A range of percentile thresholds

615 spanning $10^{-6}$ % – 15% of the whole-tree patristic distance distribution was used to choose an

616 optimal threshold point and to verify robustness of cluster composition. The minimum

617 percentile threshold that maximized the number of clusters was chosen as the optimal

618 threshold by performing multiple clustering runs on randomly sampled patristic distance

619 distributions (1 million for each run). Well-supported sub-trees (bootstrap values > 90%)

620 with mean pairwise patristic distances among *taxa* within the chosen threshold were

621 considered putative transmission clusters (i.e. clusters comprising sequences

622 epidemiologically linked through a transmission chain, although some of the direct links may

623 be missing). Only clusters containing at least 1 Italian sequence were considered in

624 downstream analyses. The phytools package [39] in R was used for joint likelihood

625 reconstruction of discrete ancestral origins [40] according to country (and associated

626 uncertainty) for the most recent common ancestor (MRCA) of each transmission cluster

627 within the ML tree (scaled in substitutions/site) for the combined dataset. Transition rates

628 among discrete states (countries) along tree nodes were considered to be equal. The tree

629 scaled in time was used to attribute temporal origins to each cluster, or time of MRCA

630 (TMRCA). The following R packages were used in the manipulation of data for cluster

631 characterization and visualization: ape [41], dplyr [42], purr [43], rlist [44], tidytree [45], ggplot2 [46],

632 data.table [47], reshape2 [48], lubridate [49], ggtree [50], tidyr [51], parallel [36].

633

634 *Estimation of basic reproduction number*

635 Estimates for daily basic reproduction number, *Re*, of SARS-CoV-2 in Italy were obtained

636 from the COVID-19-re data repository (https://github.com/covid-19-Re/dailyRe-Data) as at

637 20th September 2020. The effective reproductive number describes the average number of

638 secondary infections caused by an infected individual. The relevant method of calculation of

639 Re builds upon another method developed by Cori et al.[52], accessible through EpiEstim R

640 package. Instead of using a time series of infection incidence, which cannot be observed

641 directly, the relevant method infers the infection incidence time series based on secondary

642 sources of information such as COVID-19 confirmed case data, hospital admissions, and

23

643      deaths. This was considered in combination with two other sets of time variables: i) the
644      duration of SARS-CoV-2 incubation period and ii) the time delays between the onset of the
645      symptoms and a positive test, a hospital admission or the death of a patient. The relevant
646      method infers infection time series from the stated observed incidence data by deconvolution
647      [52–55].

648

649      *Epidemiology data assembly*
650      We analysed COVID-19 cases counts in Italy from publicly released data up to October 31[st]
651      2020 from the Italian Civil Protection Department repository (https://github.com/pcm-
652      dpc/COVID-19) that releases daily updates on the number of new confirmed cases, deaths
653      and recoveries, with a breakdown by region. To illustrate the epidemic progression, the daily
654      number of confirmed cases of people infected with SARS-Cov-2 in Italy was plotted
655      alongside a timeline of lockdown phases and variation in estimated virus reproduction
656      number until October 31[st] 2020. For convenience the geographical locations were aggregated
657      by Italian macro regions: Northeast, Northwest, Central, South, and Insular, which are basic
658      regions for the application of regional policies (Italian regions). Mobility data over time,
659      combining data on three different forms of transportation - personal vehicle, public, and
660      walking - were obtained from Apple Mobility Trends Reports
661      (https://covid19.apple.com/mobility).

662

663      *Agent-based stochastic model simulation of the Italian epidemic*
664      The Italian epidemic was simulated using the forward-time, agent-based stochastic
665      transmission chain simulator, nosoi [56], which allows for time-varying parameterization. The
666      simulation was initiated with a single infected individual with a probability of transmission of
667      0.02 per day following an incubation period, which was set to a mean of 5 days and standard
668      deviation of 2 days. The rate of transmission was fixed throughout the simulation at this
669      value. Following the incubation period, each individual was considered infectious for
670      approximately 9 days, exiting the simulation at a mean time of 14 days (standard deviation
671      equal to 2) after infection. Three different simulation scenarios were tested, assuming a direct
672      relationship between 1) rate of mobility and the number of other individuals with which each
673      infected individual comes into contact, 2) hospitalization rate and the rate at which an
674      infected individual was removed from the simulation during the infectious period
675      (approximately 5 to 14 days), or both. Hospitalization and mobility data were standardized by
676      (observed-minimum)/(maximum-minimum) to a range of 0-1 and modeled using a Fourier

677 series periodic function, with mobility data comprised of 2 sine/cosine terms (linear model

678 regression R2=0.8484) and hospitalization data of 4 sine/cosine terms (linear model

679 regression $R^2$=0.984). In scenarios 1 and 3, the probability of exiting the simulation was

680 allowed to vary over time proportionally to the standardized rate of hospitalization, resulting

681 in a maximum of ~35% of infected individuals hospitalized during peak hospitalization of the

682 epidemic. In scenarios 2 and 3, the number of contacts for each infected individual was

683 allowed to vary over time proportionally to the mobility rate, resulting in a mean of

684 approximately 15 individuals in contact with each infected individual. For scenario 1, a static

685 mean removal rate (during the infectious period described above) over time was set to 0.04

686 (standard deviation of 0.01). For scenario 2, the mean number of contacts per individual was

687 set to 15 (standard deviation of 8). The mean absolute error for each time point was

688 calculated to assess the deviation of the simulated number of actively infected individuals for

689 each of the three scenarios from the true number of cases, provided by Italian Ministry of

690 Health and the Civil Protection Department.

691

692 *Data Availability*

693 Sequence and epidemiology raw data utilized, generated or analyzed during these studies are

694 available from the authors upon request (including sequence alignment and R scripts for the

695 phylodynamic analyses).

696

697 **Supplementary Figures Legend**

698

699 **Supplementary Figure 1.** Estimates of the viral effective population size (Ne) in the Italian

700 epidemic. Estimates of the viral effective population size (Ne) using a sample of 100

701 phylogenetic trees with the highest log-likelihood values. Encircled letters (A, B and C) label

702 the three major patterns inferred from collection of trees with the highest likelihood values.

703

704 **Supplementary Figure 2.** Spatial and temporal distribution of SARS-CoV-2 trough Italian

705 regions. Italian regions were aggregated into five macro regions (NUTS, Nomenclature of

706 territorial units for statistics): Northeast, Northwest, Central, South and Insular. The left-hand

707 Y-axis (left) represents incidence (cases per 100K population, black curve) of COVID-19,

708 while the secondary Y-axis (right) represents the number of deaths (red curve) related to

709 COVID-19 through the country. Y-axis numbers are represented as log10 for visual purposes.

710 Colours at the bottom represent the epidemic phases in Italy: pre-lockdown (early
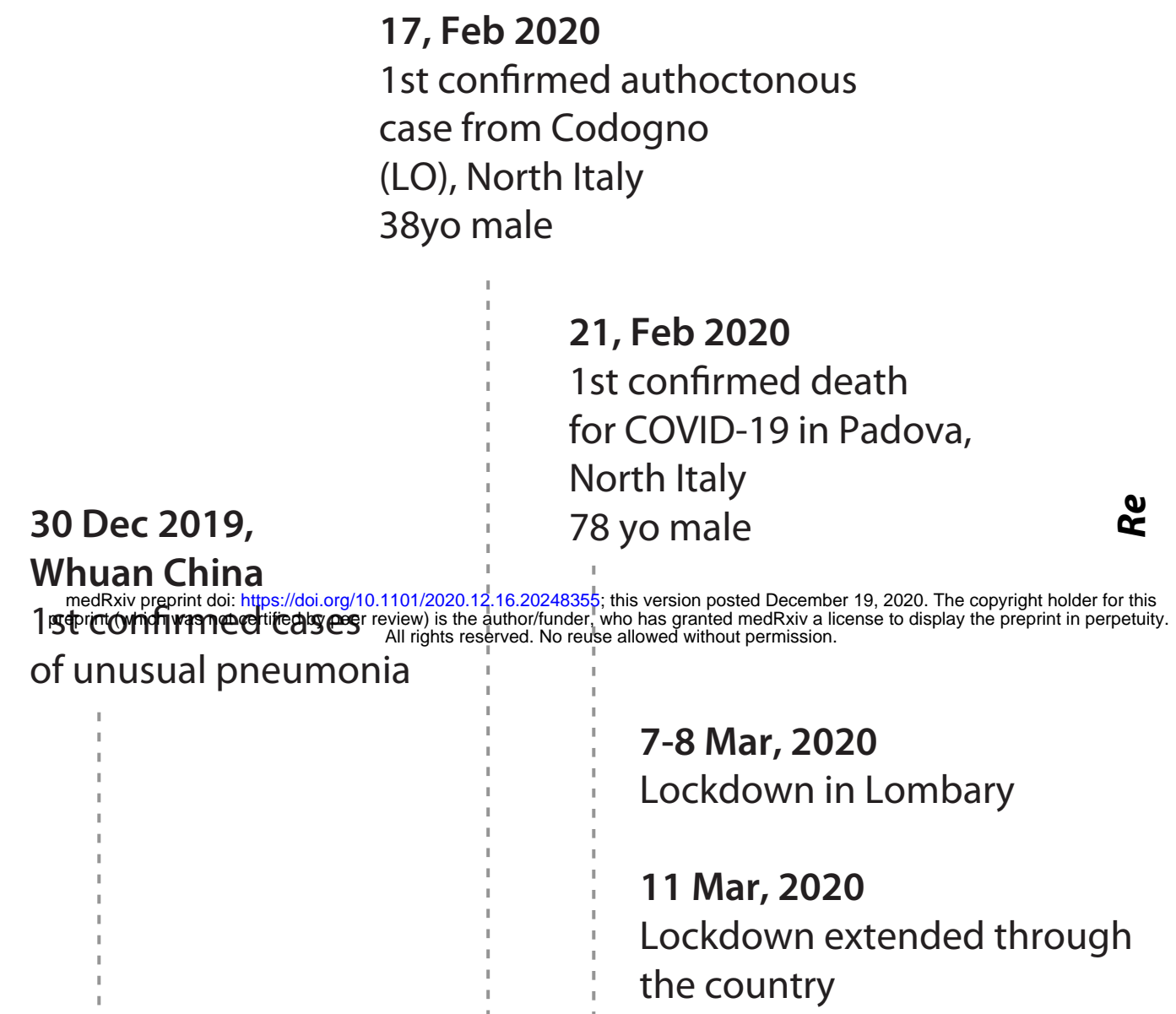
711    transmission) in light red, phase 1 in yellow, phase 2 in light violet and phase 3 in purple.

712    Maps of Italy were superimposed to showing to exact location of each region: Northwest;

713    Northeast; Central; South and Insular.

714

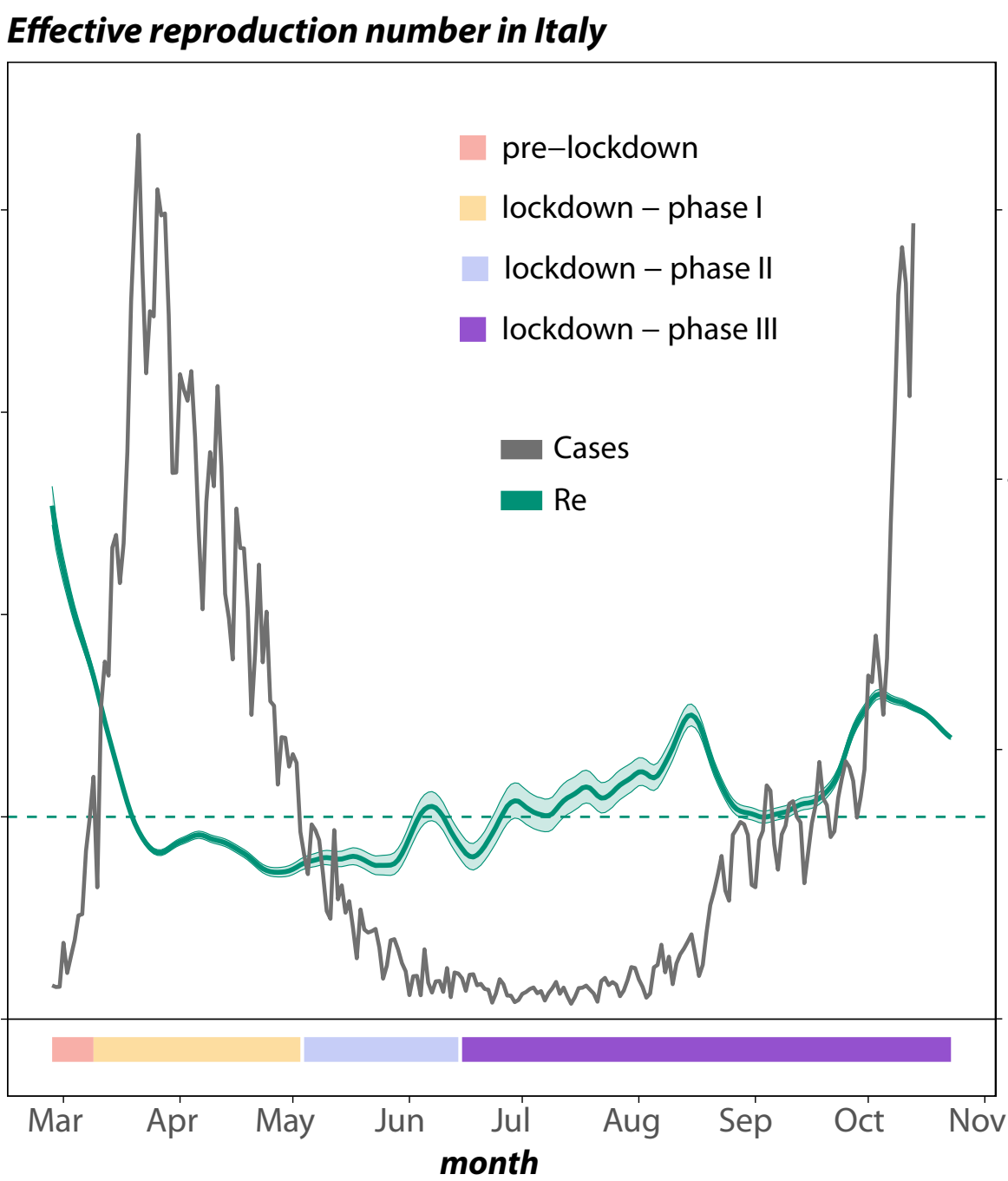715    **Supplementary Figure 3.** Frequency and distribution of SARS-CoV-2 lineages in Italy over

716    time.

717

718    **Supplementary Figure 4.** Analysis of temporal structure. **A.** Root-to-tip genetic divergence

719    of Italian sequences against time of sampling. **B.** Root-to-tip genetic divergence for the whole

720    dataset (Italian strains + reference sequences) against time of sampling.

721

722    **Supplementary Figure 5.** Empirical observations (points) and chosen model (lines) for

723    hospitalization rates (red) and number of mobile individuals (blue) over time. Hospitalization

724    rates represent the number hospitalized per positive case, whereas mobility represent the

725    number of individuals utilizing walking, as well as public and personal modes of

726    transportation, as primary means of mobility (blue). Values were standardized for

727    comparison. Empirically observed number of new infections (standardized) are also shown
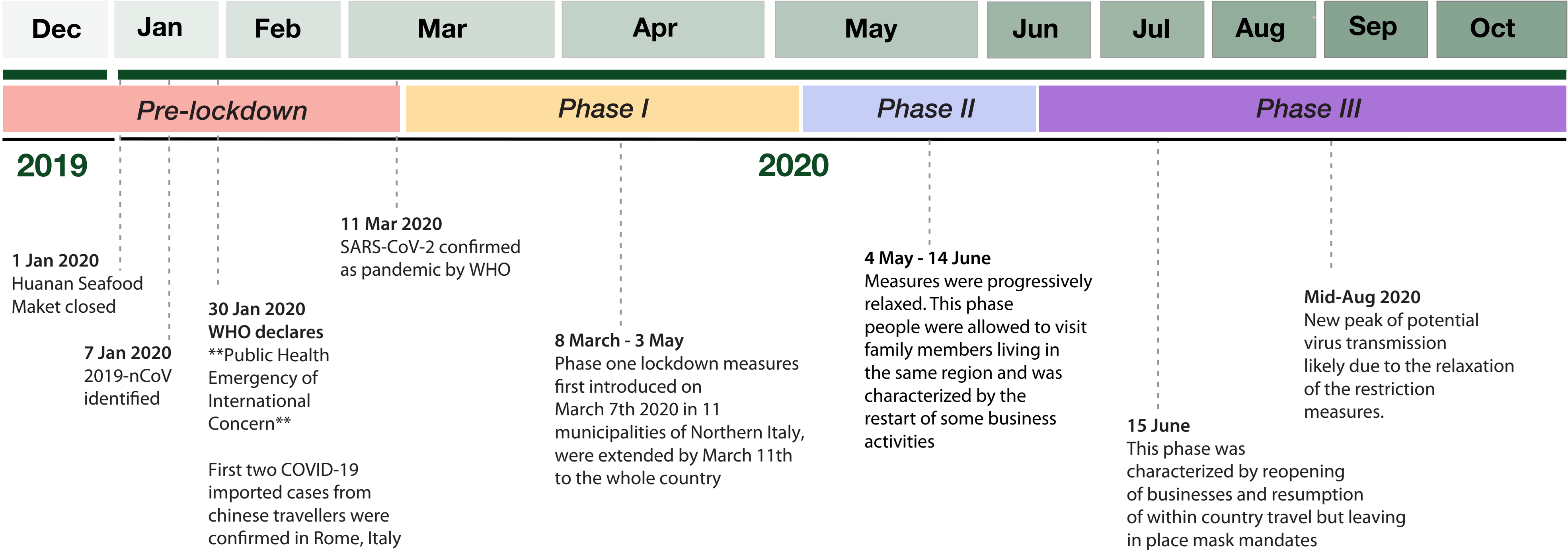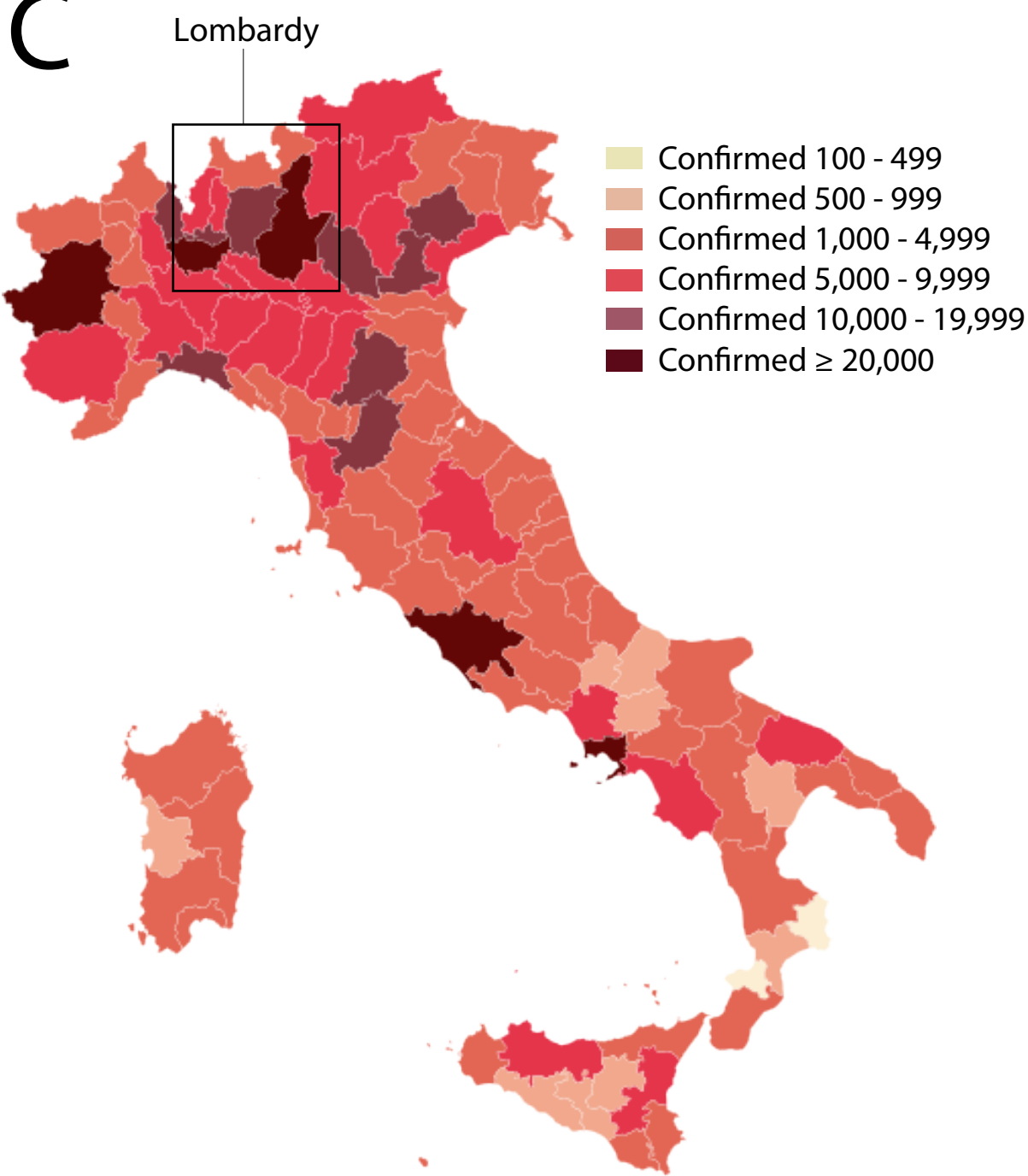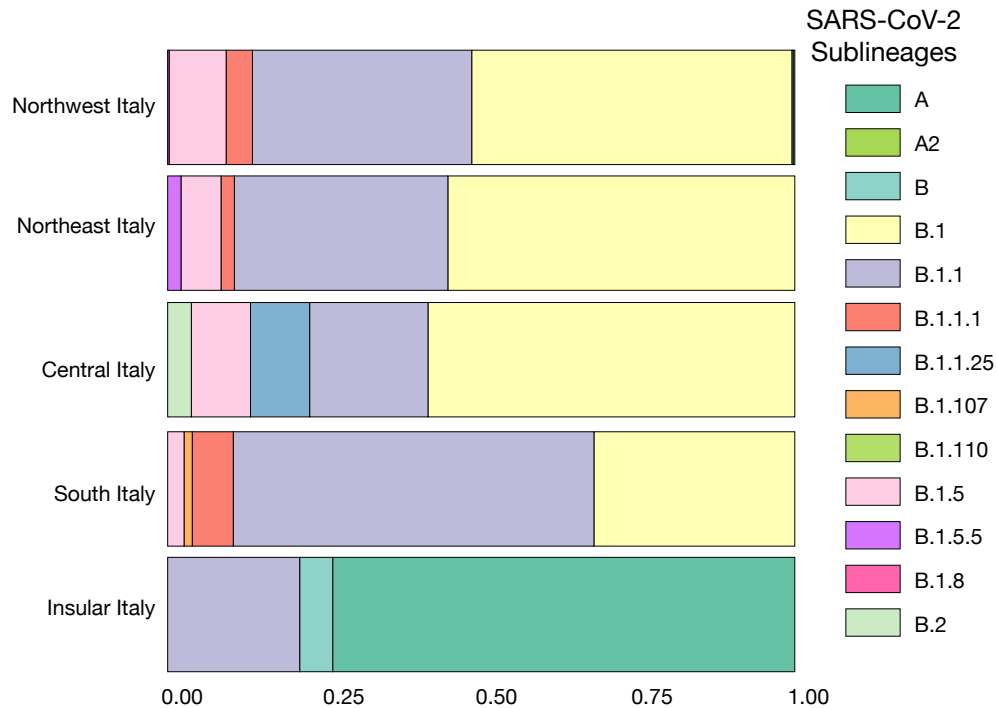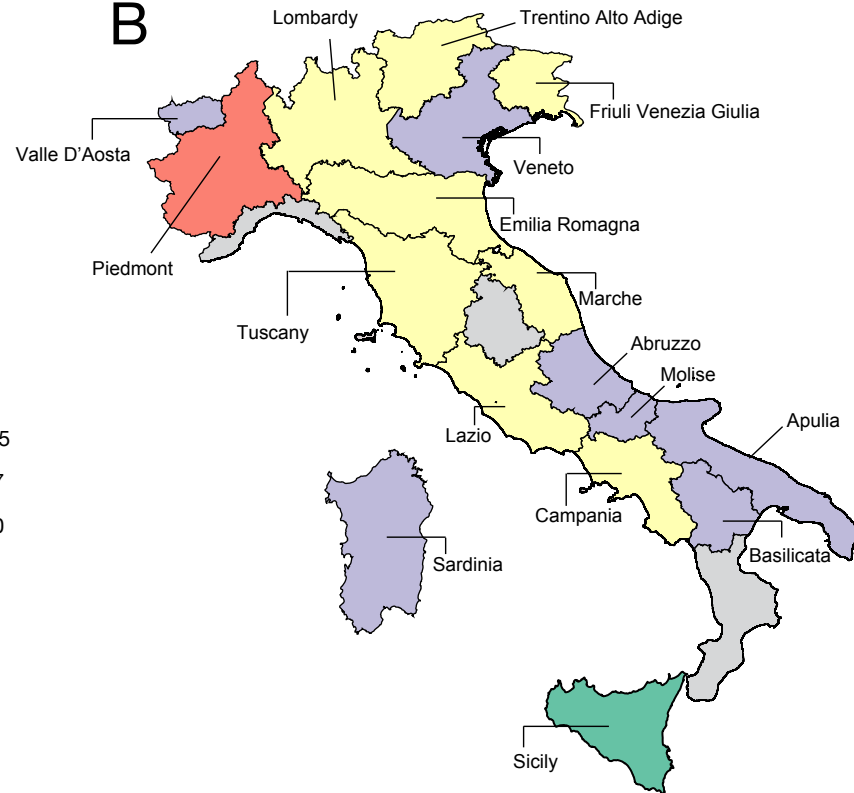
728    for comparison in orange.

729

730

**A**

**30 Dec 2019, Whuan China**
1st confirmed cases of unusual pneumonia

**17, Feb 2020**
1st confirmed authoctonous case from Codogno (LO), North Italy
38yo male

**21, Feb 2020**
1st confirmed death for COVID-19 in Padova, North Italy
78 yo male

**7-8 Mar, 2020**
Lockdown in Lombary

**11 Mar, 2020**
Lockdown extended through the country

**B**

*Effective reproduction number in Italy*

- pre-lockdown
- lockdown – phase I
- lockdown – phase II
- lockdown – phase III

— Cases
— Re

*Re*

*cases*

*month*

Mar Apr May Jun Jul Aug Sep Oct Nov

**C**

Lombardy

- Confirmed 100 - 499
- Confirmed 500 - 999
- Confirmed 1,000 - 4,999
- Confirmed 5,000 - 9,999
- Confirmed 10,000 - 19,999
- Confirmed ≥ 20,000

| Dec | Jan | Feb | Mar | Apr | May | Jun | Jul | Aug | Sep | Oct |

*Pre-lockdown* | *Phase I* | *Phase II* | *Phase III*

**2019**

**2020**

**1 Jan 2020**
Huanan Seafood Maket closed

**7 Jan 2020**
2019-nCoV identified

**30 Jan 2020**
WHO declares
**Public Health Emergency of International Concern**

First two COVID-19 imported cases from chinese travellers were confirmed in Rome, Italy

**11 Mar 2020**
SARS-CoV-2 confirmed as pandemic by WHO

**8 March - 3 May**
Phase one lockdown measures first introduced on March 7th 2020 in 11 municipalities of Northern Italy, were extended by March 11th to the whole country

**4 May - 14 June**
Measures were progressively relaxed. This phase people were allowed to visit family members living in the same region and was characterized by the restart of some business activities

**15 June**
This phase was characterized by reopening of businesses and resumption of within country travel but leaving in place mask mandates

**Mid-Aug 2020**
New peak of potential virus transmission likely due to the relaxation of the restriction measures.

**A**

Northwest Italy

Northeast Italy

Central Italy

South Italy

Insular Italy

0.00    0.25    0.50    0.75    1.00

SARS-CoV-2
Sublineages

- A
- A2
- B
- B.1
- B.1.1
- B.1.1.1
- B.1.1.25
- B.1.107
- B.1.110
- B.1.5
- B.1.5.5
- B.1.8
- B.2

**B**

Lombardy
Trentino Alto Adige
Friuli Venezia Giulia
Valle D'Aosta
Veneto
Piedmont
Emilia Romagna
Marche
Tuscany
Abruzzo
Molise
Lazio
Apulia
Campania
Basilicata
Sardinia
Sicily

A

Italy mutation map

241: C to T

3037: C to T

14408: C to T

23403: A to G

614D

614G

Pre-lockdown

Phase 1

Phase 2

Phase 3

A

B

C

D

— Predicted cases
— Observed cases

— ■ — Hospitalization and Mobility
— ■ — Hospitalization only       — ■ — Mobility only

**A**

Cofficient correlation = 0.68

R squared = 0.47

**B**

**Location**  ● Italy  ● Other

Cofficient correlation = 0.48

R squared = 0.23