1    # Early pandemic molecular diversity of SARS-CoV-2 in children

2

3    Ahmed M. Moustafa PhD[1]; William Otto MD[1], Xiaowu Gai PhD[2,3], Utsav Pandey PhD[2],

4    Alex Ryutov PhD[2], Moiz Bootwalla MS[2], Dennis T Maglinte[2], Lishuang Shen PhD[2],

5    David Ruble[2], Dejerianne Ostrow PhD[2], Jeffrey S. Gerber MDPhD[1,4], Jennifer Dien Bard

6    PhD[2,3], Rebecca M. Harris[1,5] PhD, Paul J. Planet[1,4,6*] MDPhD

7

8    1. Division of Pediatric Infectious Diseases, Children's Hospital of Philadelphia,

9    Philadelphia, PA 19104, USA.

10    2. Department of Pathology and Laboratory Medicine, Children's Hospital Los Angeles,

11    Los Angeles, CA,

12    3. Keck School of Medicine, University of Southern California, Los Angeles, CA

13    4. Department of Pediatrics, Perelman College of Medicine, University of Pennsylvania,

14    Philadelphia, PA 19104, USA.

15    5. Department of Pathology and Laboratory Medicine, Perelman School of Medicine,

16    University of Pennsylvania, Philadelphia, PA, USA

17    6. Sackler Institute for Comparative Genomics, American Museum of Natural History,

18    New York, NY 10024, USA.

19    **Emails**

20    **PJP: planetp@email.chop.edu**

21    **AMM: moustafaam@chop.edu**

22    ***Corresponding Author (AMM is alternate corresponding author)**

23

24 **Keywords**

25 *COVID-19, lineages, GNUVID, wgMLST, clonal complex*

26

27 **Running title**

28 SARS-CoV-2 diversity in children

29

30 **Summary**

31 Using sequencing and a novel technique for quantifying SARS-CoV-2 diversity, we

32 investigated 169 SARS-CoV-2 genomes (83 <21 years old). This analysis revealed

33 unexpected diversity especially in children. No clear differences in clinical presentation

34 were associated with the different virus lineages.

35 **Abstract**

36 **Background**

37 In the US, community circulation of the SARS-CoV-2 virus likely began in February

38 2020 after mostly travel-related cases. Children's Hospital of Philadelphia began testing

39 on 3/9/2020 for pediatric and adult patients, and for all admitted patients on 4/1/2020,

40 allowing an early glimpse into the local molecular epidemiology of the virus.

41 **Methods**

42 We obtained 169 SARS-CoV-2 samples (83 from patients <21 years old) from March

43 through May and produced whole genome sequences. We used genotyping tools to

44 track variants over time and to test for possible genotype associated clinical

45 presentations and outcomes in children.

46 **Results**

47 Our analysis uncovered 13 major lineages that changed in relative abundance as cases

48 peaked in mid-April in Philadelphia. We detected at least 6 introductions of distinct viral

49 variants into the population. As a group, children had more diverse virus genotypes than

50 the adults tested. No strong differences in clinical variables were associated with

51 genotypes.

52 **Conclusions**

53 Whole genome analysis revealed unexpected diversity, and distinct circulating viral

54 variants within the initial peak of cases in Philadelphia. Most introductions appeared to

55 be local from nearby states. Although limited by sample size, we found no evidence that

56 different genotypes had different clinical impacts in children in this study.

57

58 **Background**

59       After an initial period in January 2020 when most severe acute respiratory

60 coronavirus 2 (SARS-CoV-2) infections in the US were travel-related, the virus quickly

61 established itself during February with sustained, community spread[1]. Studies tracking

62 the spread of the virus using whole genome phylogenetics suggested multiple

63 introductions during this time period from Europe and Asia [2-7], as well as multiple

64 waves of transmission of distinct variants that differ locally[8].

65       Understanding genotypic diversity in local molecular epidemiology is critical for

66 tracking spread and new introductions, identifying hotspots, and enhancing contact

67 tracing[2, 4, 5, 7]. However, the biological significance of viral diversity is not known. For

68 instance, it is unclear if lineages differ in virulence or transmissibility[4, 9]. It is also

69 unclear if the immune response will be equally protective against all variants of the

70 virus, highlighting the need to understand SARS-CoV-2 diversity and evolution for

71 vaccine development[2, 10]. Moreover, there is little known about viral diversity across

72 the lifespan, with limited data on SARS-CoV-2 genomic diversity in pediatric

73 populations[11].

74       The first case of coronavirus disease 2019 (COVID-19) in Philadelphia was

75 reported on March 10, 2020 (https://www.media.pa.gov/pages/health-

76 details.aspx?newsid=734), 14 days after the first non-travel related case was confirmed

77 in California[1] and less than a week after the first cases of community spread in New

78 York State (https://www.governor.ny.gov/news/during-coronavirus-briefing-governor-

79 cuomo-signs-40-million-emergency-management-authorization). On March 9[th] the

80 infectious disease diagnostic laboratory (IDDL) at Children's Hospital of Philadelphia

81  (CHOP) became one of the first locations in the region to offer PCR-based testing for

82  SARS-CoV-2, and worked with local authorities to provide testing for both children and

83  adults in the community. On April 1$^{st}$, CHOP instituted universal screening for all

84  admitted children.

85      To track the molecular epidemiology of the virus locally in Philadelphia, and

86  especially in a pediatric population, we obtained 169 samples from the initial period of

87  testing between 3/19/2020 to 5/4/2020 and performed whole genome sequencing

88  (WGS).  Eighty-three samples were from patients less than 21 years old. We used our

89  genotyping tool GNUVID[8] to classify and compare these strains to the growing global

90  database of SARS-CoV-2 sequences at GISAID[12] (Supplementary Table 1)[13]. Here

91  we show that the early pandemic and peak in Philadelphia were characterized by

92  multiple, diverse, circulating viral variants, especially amongst children. We also

93  observed multiple introductions from distinct geographical origins. We report statistics

94  for clinical presentation and outcomes associated with each viral genotype in children.

95

96  **Methods**

97      All nasopharyngeal swab samples that had residual volume after initial laboratory

98  processing, from individuals that had positive PCR testing for SARS-CoV-2, were

99  obtained for this study. RNA was extracted from nasopharyngeal swab samples using

100  either the Roche MagNA Pure LC (Roche) or EZ1 virus mini kit (Qiagen) using magnetic

101  bead technology. Whole genome sequencing was done by the Children's Hospital Los

102  Angeles (CHLA) Center for Personalized Medicine and the Virology Laboratory. Briefly,

103  WGS of extracted viral RNA was performed as previously described using Paragon

104    Genomics CleanPlex SARS-CoV-2 Research and Surveillance NGS Panel[11, 14].

105    Libraries were quantified using the Agilent High Sensitivity D1000 ScreenTape assay

106    then normalized and pooled on the Biomek i7 liquid handler (Beckman Coulter Life

107    Sciences) to approximately 1nM. The resulting pool was quantified again using the

108    TapeStation High Sensitivity D1000 assay and diluted to a final concentration of 500pM;

109    libraries were denatured and diluted according to Illumina protocols and loaded on the

110    NextSeq 500 at 0.6pM. Paired-end and dual-indexed 2x150bp sequencing was done

111    using NextSeq 500 High Output Kit (300 Cycles).

112        All SARS-CoV-2 genomes (n=169)[13] were queried against the GNUVID

113    database (version August 17th 2020) that has 32,719 high coverage complete

114    genomes[8, 12]. Each genome was assigned an ST profile and CC. A minimum

115    spanning tree (MST) was then constructed using the goeBURST algorithm[15, 16] to

116    group STs into larger taxonomic units, clonal complexes (CCs), which we define as

117    clusters of >20 STs that are single or double allele variants away from a "founder"[8,

118    17]. Temporal plots were extracted using a custom script and plotted in GraphPad

119    Prism v7.0a. The genomes were also assigned to a lineage[2] using pangolin

120    (https://github.com/hCoV-2019/pangolin). A custom script was used to check the

121    specific combinations of 9 GISAID genetic markers, and genomes were assigned to the

122    GISAID clades. The genomes were grouped by different age groups and the relative

123    abundance of the STs and the 13 CCs were calculated. To compare the Shannon

124    diversity index between the different groups[18], a t-test was used to determine whether

125    the indices were significantly different[19].

126    To show the relationship amongst the genomes of the 169 isolates and the global

127    diversity of SARS-CoV-2, a maximum likelihood tree was constructed. Briefly,

128    consensus SARS-CoV-2 sequences for the 169 CHOP isolates were combined with full-

129    length SARS-CoV-2 sequences of 25,807 additional isolates from GISAID[12] that are

130    part of the GNUVID August database release[17] and have an assigned CC and date of

131    isolation (Supplementary Table 1)[13] to generate a multiple sequence alignment using

132    MAFFT's FFT-NS-2 algorithm[20] (reference MN908947.3[21], options: --add --

133    keeplength). The 5' and 3' untranslated regions were masked in the alignment file using

134    a custom script. A maximum likelihood tree using IQ-TREE 2[22] was then estimated

135    using the HKY model of nucleotide substitution[23], default heuristic search options, and

136    ultrafast bootstrapping with 1000 replicates[24]. The tree was rooted to MN908947.3.

137    The tree and the six GISAID clades data were visualized in iTOL[25]. The tree and the

138    tip dates were then used in TempEst[26] to estimate the evolutionary rate. Similar

139    procedures were used to construct two trees for both CC4 and CC258 and then

140    estimate the evolutionary rates. Commands used for producing the figures are available

141    in Supplementary Material.

142    Manual review of the electronic health record was performed for all patients who

143    tested positive for SARS-CoV-2 to obtain data on test characteristics, demographic

144    data, exposures, comorbidities, symptomatology, clinical severity, and treatment

145    information and deidentified. Samples were obtained under CHOP IRB protocol 17-

146    014648 as part of routine clinical care, solely for non-research purposes, carrying

147    minimal risk, and were therefore granted a waiver of informed consent. Summary

148    statistics were used to describe demographic and outcome data. Non-parametric

149    methods were used due to our small samples size, and to minimize the effect of outliers

150    on statistical associations. Multivariable logistic regression was used to evaluate the

151    association between viral sequence types and clinical outcomes. All statistics were

152    performed with STATA version 15.0, (Stata Corp., College Station, TX).

153

154    **Results**

155        Over the time period of this study, CHOP IDDL performed 4486 tests for SARS-

156    CoV-2 of which 246 (5.48%) were positive. Of the 246 positives in patients <21 years of

157    age, we were able to obtain samples from 71 patients. Of the 71 patients, 15 were

158    admitted, 3 to the intensive care unit (ICU), and 2 needed respiratory support. We also

159    obtained samples from 12 other children and 86 adults tested by the CHOP IDDL for a

160    total of 169 sequences in this study.

161        Using the GNUVID classifier[8, 12], we genotyped all 169 genomes and assigned

162    a sequence type (ST), which we define as the group of sequences that have exactly the

163    same allelic haplotype. When possible, each ST was then classified into a clonal

164    complex (CC), defined as a group of STs that differ by only one or two alleles from a

165    central "founder" sequence determined by minimum-spanning clustering[8]. Overall, we

166    identified 112 distinct STs in our data, 108 (165 genomes) of which could be assigned

167    to one of 13 CCs when compared to the most recent global GISAID genome

168    database[8, 12, 17]. While 13 STs (56 genomes) had an exact genotype match in the

169    global database, 99 STs (113 genomes) were novel, with previously unobserved alleles

170    that were not due to sequencing ambiguity based on sequence quality. The genomes

171    were widely distributed across the global SARS-CoV-2 phylogeny suggesting multiple

8

172    introductions (Figure 1A, Supp Fig. 1A). Temporal mapping of the viral CCs by week of

173    isolation showed the persistent predominance of CC258, but also persistence of

174    multiple, diverse haplotypes in the population (Figure 1B).

175         We estimated the number of putative introductions into our population by

176    comparing our data to high quality sequences from the global GISAID dataset[8, 12,

177    17], and requiring an identical ST to have been isolated in another geographic location

178    at least 10 days prior to the isolation date in our sample. Using this criterion, we

179    identified 6 independent STs that were likely introductions into our population (Table 1).

180    One of these putatively introduced genotypes, ST6228, had only ever been observed in

181    New York State before, and thus likely represents an introduction from this neighboring

182    state. ST338 and ST258 were also observed in New York State in the 10 days prior to

183    appearing in our population, but they were also widespread internationally during this

184    time period, and therefore could have been introduced from other sources. For ST258,

185    isolates were observed during this time window in 24 countries and 22 States including

186    Pennsylvania and other nearby states such as New Jersey. ST4 and ST1531 were

187    observed closest to Philadelphia in Washington DC and Virginia in the 10 days prior to

188    appearing in our population. The most likely international introduction was ST6134,

189    which was seen previously only in Australia. If we shortened the criterion to isolation 5

190    days prior, we detected 3 more putative introductions. All 3 of these STs were first

191    observed in New York.

192         To detect any exportations of viral genotypes, we looked for STs that were seen

193    in our dataset 10 days prior to isolation in another geographic location. Only one

194    possible exportation event was detected of ST13162 to Wisconsin.

195     It should be noted that our method of detecting introductions relies on robust

196     sampling both in our population and in other locations. The detected number of

197     importations and exportations is likely much higher than the numbers we were able to

198     find here, and estimates may grow as more genome sequences are added from

199     retrospective sampling.

200     The relative abundance of the 13 CCs found in our dataset was distributed

201     differently between children and adults, with the pediatric population showing

202     considerably more diversity (Shannon Entropy=1.815 vs 1.412, P = 0.0132). CC4, an

203     early lineage originally seen in Wuhan, was more prevalent in pediatric cases (20%)

204     compared to adults (14%). CC258, a lineage that predominated in Europe and New

205     York, was more prevalent in adults (55%) compared to children (40%). A more granular

206     analysis of STs recapitulated the higher diversity of viral types in the pediatric

207     population, but did not achieve statistical significance (Shannon Entropy= 2.624 vs

208     2.456, P= 0.3557).

209     One clear difference between our dataset and data from neighboring states over

210     the same time period is the increased diversity of CCs and the presence of the early

211     genotype CC4 (e.g., for NY v. our sample Shannon Entropy=1.69 vs 1.15, P = 4.23E-7).

212     It is unclear whether this reflects specific epidemiology of Philadelphia, our focus on

213     pediatric samples, or other biases in this convenience sample. Interestingly, while there

214     were only 6 STs observed in CC4 (5 STs in children and 2 in adults), there were 57 STs

215     from CC258 (25 STs in children and 38 in adults) demonstrating the much higher

216     diversity of genotypes associated with the CC258 lineage, and potentially the large

217     amount of diversification of this lineage as it peaked to very high numbers in nearby

10

218    New York City. To address the cause of this diversity, we calculated mutation rates for

219    CC4 and CC258 genomes using our genomes as well as genomes from the GISAID

220    database using TempEst[26] (Supp Fig 1B). The mutation rate for CC4 was $2.2 \times 10^{-4}$

221    sites/year while the mutation rate for CC258 was $5.9 \times 10^{-4}$ sites/year. The rate across all

222    GISAID sequences was $7.1 \times 10^{-4}$ in line with previous estimates. It is possible that both

223    had a higher mutation rate and a large effective population size through increased

224    transmission contributed to the higher diversity seen in CC258.

225        To assess the possibility that different genotypes were associated with distinct

226    clinical outcomes and presentations, we collected demographic and clinical information

227    for 71 pediatric viral genomes from patients in the CHOP Care Network. Although

228    limited by the sample size, we were unable to detect any significant differences in

229    specific clinical variables associated with the different genotypes (Tables 2 and 3 and

230    Figure 2). However, exploratory analysis of the data suggested that pediatric patients

231    infected with CC4 lineage virus and early pandemic genotypes (e.g., GISAID lineage L

232    and Pangolin lineage B) may have had increased rates of admission to the hospital

233    (odds ratio, OR 17.2, 95% confidence interval 2.23 to 132.13, P = 0.006) compared to

234    those infected with the CC258 lineage (Supplementary Tables 2, 3 and 4) and lineages

235    considered to be more derived (eg., GISAID lineage GH and Pangolin lineage B.1). In

236    addition, two of the single nucleotide polymorphisms (SNPs) (Table 4 and

237    Supplementary Tables 5 and 6) from more ancestral haplotypes (e.g., C241T, C3037T)

238    were also significantly associated with admission (Supplementary Tables 7, 8, 9 and

239    10). The D614G (SNP; A23403G) spike protein mutation was associated with less

240    hospital admission, albeit not statistically significant (OR 0.23, 95% CI 0.05-1.13)

241  (Supplementary Table 11), but it was the only SNP tested that was significantly

242  associated with decreased odds of being asymptomatic (OR 0.11, 95% CI 0.01-0.92)

243  (Supplementary Table 12).

244

245  **Discussion**

246       We have shown that the early pandemic in Philadelphia was diverse and

247  dynamic, with multiple likely introductions, most probably from local spread of the virus

248  from neighboring states. Although CC258, the clonal complex thought to have been

249  introduced from Europe that dominated in New York[4, 8], also predominated in our

250  sample across the early pandemic, other CCs were robustly present. For instance, CC4,

251  one of the earliest genotypes seen in Wuhan, persisted throughout the study period

252  demonstrating sustained spread in the community. Other CCs (e.g., CC3530, CC300,

253  CC1508) were also seen persistently in this sample implying sustained community

254  spread. This finding suggests that there was enough viral diversity early in the

255  pandemic that contact tracing may have been significantly enhanced by whole genome

256  (or targeted SNP detection) comparisons.

257       It is important to note that most of the putative introductions into our population

258  could be traced to nearby states surrounding the Philadelphia area, and only one

259  putative international introduction was detected. This may reflect international travel

260  restrictions in place at this time, but it also suggests that most spread was local, and

261  that there were missed opportunities to limit these events particularly in travel to and

262  from New York. It is important to note that as the database of SARS-CoV-2 genomes

263  grows and more genome sequences are available from the Philadelphia area, we may

12

264    find new evidence for introductions or importations, which likely far outnumber those

265    detected in our analysis.

266        Although the viral genotypes in our sample differed at several putatively key

267    amino acid locations, we did not detect any stark differences in clinical presentation or

268    outcome in children (Tables 2 and 3). Previous studies have shown that different

269    nucleotide variants or deletions may be associated with higher or lower severity [27, 28].

270    However, the small sample size and higher than expected viral diversity might have led

271    to an inability to discriminate smaller effect sizes. It should also be noted that the

272    retrospective nature of this study, incomplete sampling, and inconsistent capture of

273    symptoms and severity, could have biased these data. Nonetheless, it is still possible

274    that genetic differences between viral lineages may have an impact on virulence or

275    clinical outcome, and our observed differences in admission rates raises the possibility

276    that larger studies may uncover differences in the future. Notably, another recent

277    pediatric study of 141 SARS-CoV-2 in California, which assessed clinical characteristics

278    of 88 patients, demonstrated a possible association between a specific genotype and

279    disease severity[11].

280        It is also possible that genetic variants may have differential transmission

281    abilities, which could not have been detected directly using our data. However, it is

282    worth noting that the genotypes (CC4, CC750 and CC1508) that have the ancestral

283    alanine residue at position 614 in the spike protein persisted and spread throughout the

284    study period, suggesting that the derived allelic form (A23403G; D614G) that has been

285    proposed to be more transmissible[29] and is predominantly represented by CC258 in

286    our analysis, did not completely dominate the ancestral form over this amount of time.

287   Here we also showed much higher diversity in the CC258 lineage and a higher

288   estimated mutation for this CC in general. It is possible that this diversity is driven by

289   higher transmissibility and a large effective population size.

290        Overall, our findings suggest that whole genome sequencing and genotyping of

291   circulating clones could be used to track viral spread and identify opportunities for

292   intervention to stop spread from specific hotspots. The relationship between viral

293   genotype, rate of transmission, and clinical presentation and outcomes deserves further

294   exploration with increased sample size.

295

296   **List of abbreviations**

297   GNUVID      Gene Novelty Unit-based Virus Identification

298   ST      Sequence Type

299   CC      Clonal Complex

300   SARS-CoV-2  Severe Acute Respiratory Syndrome Corona Virus 2

301   COVID-19    Corona Virus Disease 2019

302   wgMLST  whole genome Multilocus Sequence Typing

303
304   **Data Sharing Statement**

305   The 169 genomes from our dataset will be available from the corresponding author and

306   available online for download through a permanent Zenodo DOI[13]. Other de-identified

307   clinical data used in the manuscript are available upon request from the corresponding

308   author. The GNUVID compressed database and GNUVID source code can be found in

309   its most up-to-date version here, https://github.com/ahmedmagds/GNUVID, under the

310   GNU General Public License.

311 **Authors' contributions**

312 AMM & PJP designed and conceptualized the study. AMM contributed to the data

313 collection, data analysis, coding, data interpretation, figures, literature review and

314 tables. WO contributed to the data collection, data analysis, data interpretation, tables,

315 and writing. XG, UP, AR, MB, DTM, LS, DR and DO contributed to the data collection,

316 data analysis, and data interpretation. JDB, JSG, RMH and PJP supervised the study

317 and contributed to the data collection, data analysis, data interpretation, and literature

318 review. AMM and PJP wrote the first draft of the manuscript. All authors reviewed and

319 approved the final manuscript.

320

334

**Conflict of Interest Statement**

336    The authors declare that they have no competing interests and they do not have a

337    commercial or other association that might pose a conflict of interest.

338

**References**

340    1. Jorden MA, Rudman SL, Villarino E, et al. Evidence for Limited Early Spread of COVID-19
341    Within the United States, January–February 2020. *Morbidity and Mortality Weekly Report*.
342    2020;69:680-684.
343    2. Rambaut A, Holmes EC, O'Toole A, et al. A dynamic nomenclature proposal for SARS-CoV-2
344    lineages to assist genomic epidemiology. *Nat Microbiol*. 2020.
345    3. Deng X, Gu W, Federman S, et al. Genomic surveillance reveals multiple introductions of
346    SARS-CoV-2 into Northern California. *Science*. 2020;369:582-587.
347    4. Worobey M, Pekar J, Larsen BB, et al. The emergence of SARS-CoV-2 in Europe and North
348    America. *Science*. 2020.
349    5. Bedford T, Greninger AL, Roychoudhury P, et al. Cryptic transmission of SARS-CoV-2 in
350    Washington state. *Science*. 2020.
351    6. Shen L, Dien Bard J, Biegel JA, Judkins AR, Gai X. Comprehensive Genome Analysis of 6,000
352    USA SARS-CoV-2 Isolates Reveals Haplotype Signatures and Localized Transmission Patterns by
353    State and by Country. *Frontiers in Microbiology*. 2020;11.
354    7. Gonzalez-Reiche AS, Hernandez MM, Sullivan MJ, et al. Introductions and early spread of
355    SARS-CoV-2 in the New York City area. *Science*. 2020;369:297-301.
356    8. Moustafa AM, Planet PJ. Rapid whole genome sequence typing reveals multiple waves of
357    SARS-CoV-2 spread. *bioRxiv*. 2020:2020.2006.2008.139055.
358    9. Korber B, Fischer WM, Gnanakaran S, et al. Spike mutation pipeline reveals the emergence of
359    a more transmissible form of SARS-CoV-2. *bioRxiv*. 2020:2020.2004.2029.069054.
360    10. Forster P, Forster L, Renfrew C, Forster M. Phylogenetic network analysis of SARS-CoV-2
361    genomes. *Proc Natl Acad Sci U S A*. 2020;117:9241-9243.
362    11. Pandey U, Yee R, Shen L, et al. High Prevalence of SARS-CoV-2 Genetic Variation and D614G
363    Mutation in Pediatric Patients with COVID-19. *Open Forum Infectious Diseases*. 2020 (In Press).
364    12. Shu Y, McCauley J. GISAID: Global initiative on sharing all influenza data - from vision to
365    reality. *Euro Surveill*. 2017;22.
366    13. Moustafa AM, Otto W, Gai X, et al. Dataset for Early pandemic molecular diversity of SARS-
367    CoV-2 in children 2020. Available at: published online Nov 20. DOI:10.5281/zenodo.4282048.
368    14. Li C, Debruyne DN, Spencer J, et al. Highly sensitive and full-genome interrogation of SARS-
369    CoV-2 using multiplexed PCR enrichment followed by next-generation sequencing. *bioRxiv*.
370    2020:2020.2003.2012.988246.

371    15. Francisco AP, Vaz C, Monteiro PT, et al. PHYLOViZ: phylogenetic inference and data
372    visualization for sequence based typing methods. *BMC Bioinformatics*. 2012;13:87.
373    16. Francisco AP, Bugalho M, Ramirez M, Carriço JA. Global optimal eBURST analysis of
374    multilocus typing data using a graphic matroid approach. *BMC Bioinformatics*. 2009;10:152.
375    17. Moustafa AM, Planet PJ. Moustafa AM, Planet PJ. ahmedmagds/GNUVID: GNUVID v2.0:
376    Globally circulating clonal complexes as of 2020-10-20 2020. Available at:
377    DOI:10.5281/zenodo.4313855. Published online Dec 9. .
378    18. Shannon CE. The mathematical theory of communication. 1963. *MD Comput*. 1997;14:306-
379    317.
380    19. Hutcheson K. A test for comparing diversities based on the Shannon formula. *J Theor Biol*.
381    1970;29:151-154.
382    20. Katoh K, Misawa K, Kuma K, Miyata T. MAFFT: a novel method for rapid multiple sequence
383    alignment based on fast Fourier transform. *Nucleic Acids Res*. 2002;30:3059-3066.
384    21. Wu F, Zhao S, Yu B, et al. A new coronavirus associated with human respiratory disease in
385    China. *Nature*. 2020;579:265-269.
386    22. Minh BQ, Schmidt HA, Chernomor O, et al. IQ-TREE 2: New Models and Efficient Methods
387    for Phylogenetic Inference in the Genomic Era. *Mol Biol Evol*. 2020;37:1530-1534.
388    23. Hasegawa M, Kishino H, Yano T. Dating of the human-ape splitting by a molecular clock of
389    mitochondrial DNA. *J Mol Evol*. 1985;22:160-174.
390    24. Hoang DT, Chernomor O, von Haeseler A, Minh BQ, Vinh LS. UFBoot2: Improving the
391    Ultrafast Bootstrap Approximation. *Mol Biol Evol*. 2018;35:518-522.
392    25. Letunic I, Bork P. Interactive Tree Of Life (iTOL) v4: recent updates and new developments.
393    *Nucleic Acids Res*. 2019;47:W256-W259.
394    26. Rambaut A, Lam TT, Max Carvalho L, Pybus OG. Exploring the temporal structure of
395    heterochronous sequences using TempEst (formerly Path-O-Gen). *Virus Evol*. 2016;2:vew007.
396    27. Young BE, Fong SW, Chan YH, et al. Effects of a major deletion in the SARS-CoV-2 genome
397    on the severity of infection and the inflammatory response: an observational cohort study.
398    *Lancet*. 2020;396:603-611.
399    28. Voss JD, Skarzynski M, McAuley EM, et al. Variants in SARS-CoV-2 Associated with Mild or
400    Severe Outcome. *medRxiv*. 2020:2020.2012.2001.20242149.
401    29. Korber B, Fischer WM, Gnanakaran S, et al. Tracking Changes in SARS-CoV-2 Spike: Evidence
402    that D614G Increases Infectivity of the COVID-19 Virus. *Cell*. 2020;182:812-827 e819.
403
404

405        **Table 1: Introductions to Philadelphia.**

| Specimen Date | ST | CC | Days | Countries in last 10 days before appearance | First time Seen\|Date |
|---|---|---|---|---|---|
| 3/19/20 | 4 | 4 | 10 | China Iceland Malaysia Singapore United Kingdom USA (CA, MI, WI) | China/Wuhan\|2019-12-30 |
| 3/24/20 | 258 | 258 | 10 | Australia Austria Canada Chile Colombia Costa Rica Czech Republic Denmark France Germany Greece Iceland Israel Luxembourg Netherlands Portugal Russia Singapore South Korea Sweden Taiwan United Kingdom USA (AZ, CA, CO, CT, FL, GA, IL, IN, ME, MI, MN, NJ, NM, NY, PA, TX, UN, VA, VI, VT, WA, WI) | Singapore\|2020-02-16 |
| 3/30/20 | 1531 | 258 | 10 | Denmark USA (DC, VA, CA) | USA/NY\|2020-03-14 |
| 3/31/20 | 6134 | 258 | 10 | Australia | Australia\|2020-03-19 |
| 3/31/20 | 6228 | 258 | 10 | USA (NY) | USA/NY\|2020-03-21 |
| 4/6/20 | 338 | 338 | 10 | Australia, Colombia, USA (NY, WI, MA, CA, CT, MD, FL) | USA/CA\|2020-02-29 |
| 3/20/20 | 1623 | 258 | 5 | USA (NY) | USA/NY\|2020-03-12 |
| 3/24/20 | 2261 | 258 | 5 | USA (NY) | USA/NY\|2020-03-19 |
| 3/27/20 | 1841 | 3530 | 5 | New Zealand, USA (FL) | USA/NY\|2020-03-18 |

406

407

18

408 **Table 2: Overall characteristics, grouped by clonal complex (excluding those with single isolate or no clonal**

409 **complex identified).**

| | | Clonal Complex | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | Total | CC258 | CC4 | CC3530 | CC300 | CC255 | CC844 | CC1508 | CC750 |
| | 71 | 32 | 10 | 7 | 6 | 4 | 3 | 2 | 2 |
| **Age (years), median (IQR)** | 10.91 (5.6, 17.0) | 11.18 (7.1, 17.2) | 7.32 (2.55, 14.75) | 9.96 (4.73, 18.11) | 5.37 (.45, 13.19) | 10.8 (6.6, 12.8) | 8.84 (8.84, 8.84) | 16.5 (15.8, 19) | 8.5 (7.6, 9.3) |
| **Age Group** | | | | | | | | | |
|   **0-12 months** | 6 (8%) | 1 (3%) | 1 (10%) | 1 (14%) | 3 (50%) | 0 (0%) | 0 (0%) | 0 (0%) | 0 (0%) |
|   **1-5 years** | 12 (17%) | 6 (19%) | 4 (40%) | 1 (14%) | 0 (0%) | 1 (25%) | 0 (0%) | 0 (0%) | 0 (0%) |
|   **6-11 years** | 20 (28%) | 10 (31%) | 1 (10%) | 2 (29%) | 1 (17%) | 2 (50%) | 1 (100%) | 0 (0%) | 2 (100%) |
|   **12-18 years** | 24 (34%) | 12 (38%) | 4 (40%) | 1 (14%) | 2 (33%) | 1 (25%) | 0 (0%) | 2 (67%) | 0 (0%) |
|   **18-21 years** | 9 (13%) | 3 (9%) | 0 (0%) | 2 (29%) | 0 (0%) | 0 (0%) | 0 (0%) | 1 (33%) | 0 (0%) |
| **Male sex** | 32 (45%) | 13 (41%) | 4 (40%) | 5 (71%) | 3 (50%) | 3 (75%) | 2 (67%) | 0 (0%) | 1 (50%) |
| **Race/Ethnicity** | | | | | | | | | |
|   **Non-Hispanic White** | 19 (27%) | 8 (25%) | 5 (50%) | 0 (0%) | 2 (33%) | 0 (0%) | 1 (100%) | 0 (0%) | 1 (50%) |
|   **Non-Hispanic Black** | 38 (54%) | 18 (56%) | 3 (30%) | 7 (100%) | 2 (33%) | 3 (75%) | 0 (0%) | 2 (67%) | 1 (50%) |
|   **Hispanic or Latino** | 7 (10%) | 3 (9%) | 1 (10%) | 0 (0%) | 1 (17%) | 1 (25%) | 0 (0%) | 0 (0%) | 0 (0%) |
|   **Multi-racial** | 2 (3%) | 1 (3%) | 1 (10%) | 0 (0%) | 0 (0%) | 0 (0%) | 0 (0%) | 0 (0%) | 0 (0%) |

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| Hawaiian or Pacific Islander | 1 (1%) | 0 (0%) | 0 (0%) | 0 (0%) | 0 (0%) | 0 (0%) | 0 (0%) | 1 (33%) | 0 (0%) |
| Other Race or Unknown | 4 (6%) | 2 (6%) | 0 (0%) | 0 (0%) | 1 (17%) | 0 (0%) | 0 (0%) | 0 (0%) | 0 (0%) |
| **Insurance status** | | | | | | | | | |
| Commercial Insurance | 26 (37%) | 13 (41%) | 4 (40%) | 3 (43%) | 2 (33%) | 0 (0%) | 0 (0%) | 0 (0%) | 0 (0%) |
| Government or Public Insurance | 40 (56%) | 18 (56%) | 4 (40%) | 4 (57%) | 4 (67%) | 4 (100%) | 1 (100%) | 2 (67%) | 1 (50%) |
| Self-pay | 1 (1%) | 0 (0%) | 0 (0%) | 0 (0%) | 0 (0%) | 0 (0%) | 0 (0%) | 0 (0%) | 1 (50%) |
| Other or Unknown | 4 (6%) | 1 (3%) | 2 (20%) | 0 (0%) | 0 (0%) | 0 (0%) | 0 (0%) | 1 (33%) | 0 (0%) |
| **Previously Healthy** | 23 (32%) | 12 (38%) | 3 (30%) | 2 (29%) | 2 (33%) | 1 (25%) | 0 (0%) | 1 (33%) | 1 (50%) |
| **Admitted** | 15 (21%) | 3 (9%) | 5 (50%) | 3 (43%) | 1 (17%) | 1 (25%) | 1 (33%) | 0 (0%) | 1 (50%) |
| **ICU admission** | 3 (4%) | 0 (0%) | 0 (0%) | 0 (0%) | 1 (17%) | 1 (25%) | 0 (0%) | 0 (0%) | 1 (50%) |
| **Need for respiratory support** | 2 (3%) | 0 (0%) | 1 (10%) | 0 (0%) | 1 (17%) | 0 (0%) | 0 (0%) | 0 (0%) | 0 (0%) |
| **Clinical Severity** | | | | | | | | | |
| Asymptomatic | 7 (10%) | 2 (6%) | 3 (30%) | 1 (14%) | 1 (17%) | 0 (0%) | 0 (0%) | 0 (0%) | 0 (0%) |
| Mild | 60 (86%) | 29 (94%) | 6 (60%) | 6 (86%) | 4 (67%) | 4 (100%) | 3 (100%) | 2 (100%) | 1 (50%) |
| Severe | 3 (4%) | 0 (0%) | 1 (10%) | 0 (0%) | 1 (17%) | 0 (0%) | 0 (0%) | 0 (0%) | 1 (50%) |

410

411 **Table 3: Symptoms, grouped by clonal complex (excluding those with single isolate or no clonal complex**

412 **identified).**

| Factor | Total | CC258 | CC4 | CC3530 | CC300 | CC255 | CC844 | CC750 | CC1508 |
|---|---|---|---|---|---|---|---|---|---|
| | 71 | 32 | 10 | 7 | 6 | 4 | 3 | 2 | 2 |
| No Symptoms | 8 (11%) | 2 (6%) | 4 (40%) | 1 (14%) | 1 (17%) | 0 (0%) | 0 (0%) | 0 (0%) | 0 (0%) |
| Fever or cough or shortness of breath | 57 (80%) | 28 (88%) | 5 (50%) | 5 (71%) | 5 (83%) | 2 (50%) | 3 (100%) | 2 (100%) | 2 (100%) |
| Fever | 38 (54%) | 18 (56%) | 2 (20%) | 5 (71%) | 2 (33%) | 2 (50%) | 2 (67%) | 2 (100%) | 1 (50%) |
| Cough | 41 (58%) | 20 (62%) | 4 (40%) | 4 (57%) | 4 (67%) | 1 (25%) | 2 (67%) | 1 (50%) | 2 (100%) |
| Shortness of Breath | 13 (18%) | 8 (25%) | 1 (10%) | 1 (14%) | 1 (17%) | 0 (0%) | 0 (0%) | 0 (0%) | 1 (50%) |
| Anosmia | 5 (7%) | 3 (9%) | 0 (0%) | 1 (14%) | 0 (0%) | 0 (0%) | 1 (33%) | 0 (0%) | 0 (0%) |
| Aguesia | 4 (6%) | 3 (9%) | 0 (0%) | 0 (0%) | 0 (0%) | 0 (0%) | 0 (0%) | 0 (0%) | 1 (50%) |
| Sore Throat | 13 (18%) | 8 (25%) | 0 (0%) | 1 (14%) | 1 (17%) | 0 (0%) | 2 (67%) | 0 (0%) | 0 (0%) |
| Chest Pain | 4 (6%) | 2 (6%) | 0 (0%) | 0 (0%) | 0 (0%) | 0 (0%) | 0 (0%) | 0 (0%) | 1 (50%) |
| Myalgias | 12 (17%) | 5 (16%) | 0 (0%) | 1 (14%) | 0 (0%) | 0 (0%) | 2 (67%) | 1 (50%) | 0 (0%) |
| Chills | 5 (7%) | 2 (6%) | 0 (0%) | 1 (14%) | 0 (0%) | 0 (0%) | 0 (0%) | 0 (0%) | 1 (50%) |
| Headache | 23 (32%) | 11 (34%) | 1 (10%) | 3 (43%) | 0 (0%) | 2 (50%) | 2 (67%) | 1 (50%) | 1 (50%) |
| Fatigue | 7 (10%) | 5 (16%) | 0 (0%) | 1 (14%) | 0 (0%) | 1 (25%) | 0 (0%) | 0 (0%) | 0 (0%) |
| Gastrointestinal Symptoms | 12 (17%) | 8 (25%) | 0 (0%) | 1 (14%) | 1 (17%) | 0 (0%) | 0 (0%) | 1 (50%) | 1 (50%) |

413

414

415 **Table 4: Outcomes, grouped by SNP (excluding those with single type)**

| | | C2411 | | C3037 | | A23403 | | C8782 | |
|---|---|---|---|---|---|---|---|---|---|
| | Value | C | T | C | T | A | G | C | T |
| **N** | 71 | 22 | 49 | 15 | 56 | 12 | 59 | 69 | 2 |
| **Admitted** | 15 (21%) | 10 (45%) | 5 (10%) | 7 (47%) | 8 (14%) | 5 (42%) | 10 (17%) | 15 (22%) | 0 (0%) |
| **ICU admission** | 3 (4%) | 3 (14%) | 0 (0%) | 1 (7%) | 2 (4%) | 0 (0%) | 3 (5%) | 3 (4%) | 0 (0%) |
| **Need for respiratory support** | 2 (3%) | 2 (9%) | 0 (0%) | 1 (7%) | 1 (2%) | 1 (8%) | 1 (2%) | 2 (3%) | 0 (0%) |
| **Clinical Severity** | | | | | | | | | |
| Asymptomatic | 7 (10%) | 4 (18%) | 3 (6%) | 3 (20%) | 4 (7%) | 3 (25%) | 4 (7%) | 7 (10%) | 0 (0%) |
| Mild | 60 (86%) | 15 (68%) | 45 (94%) | 10 (67%) | 50 (91%) | 8 (67%) | 52 (90%) | 58 (85%) | 2 (100%) |
| Severe | 3 (4%) | 3 (14%) | 0 (0%) | 2 (13%) | 1 (2%) | 1 (8%) | 2 (3%) | 3 (4%) | 0 (0%) |

416

417 **Table 4: Outcomes, grouped by SNP (excluding those with single type)**

| | | G25563 | | T28144 | | G28882 | |
|---|---|---|---|---|---|---|---|
| | Value | G | T | C | T | A | G |
| **N** | 71 | 23 | 48 | 2 | 69 | 6 | 65 |
| **Admitted** | 15 (21%) | 8 (35%) | 7 (15%) | 0 (0%) | 15 (22%) | 1 (17%) | 14 (22%) |
| **ICU admission** | 3 (4%) | 3 (13%) | 0 (0%) | 0 (0%) | 3 (4%) | 1 (17%) | 2 (3%) |
| **Need for respiratory support** | 2 (3%) | 2 (9%) | 0 (0%) | 0 (0%) | 2 (3%) | 1 (17%) | 1 (2%) |
| **Clinical Severity** | | | | | | | |
| Asymptomatic | 7 (10%) | 4 (17%) | 3 (6%) | 0 (0%) | 7 (10%) | 1 (17%) | 6 (9%) |
| Mild | 60 (86%) | 16 (70%) | 44 (94%) | 2 (100%) | 58 (85%) | 4 (67%) | 56 (88%) |
| Severe | 3 (4%) | 3 (13%) | 0 (0%) | 0 (0%) | 3 (4%) | 1 (17%) | 2 (3%) |

418

419

420

421
422
423

424 **Figure Legends**

425 **Figure 1.** SARS-CoV-2 diversity from testing at our center. **A.** Minimum spanning tree

426 (MST) of 32,719 SARS-CoV-2 genomes showing 17,615 Sequence Types (STs) and 70

427 clonal complexes (CCs). The MST represents the most recent dataset used in GNUVID

428    as of August 17[th]. The reported 13 CCs at CHOP are in black. The pie charts show the

429    percentage distribution of genomes from the different geographic regions in each CC.

430    **B.** Temporal Plot of 13 circulating CCs representing the 169 genomes in this study and

431    their relative abundance in Pennsylvania (PA) and the neighboring states; New York

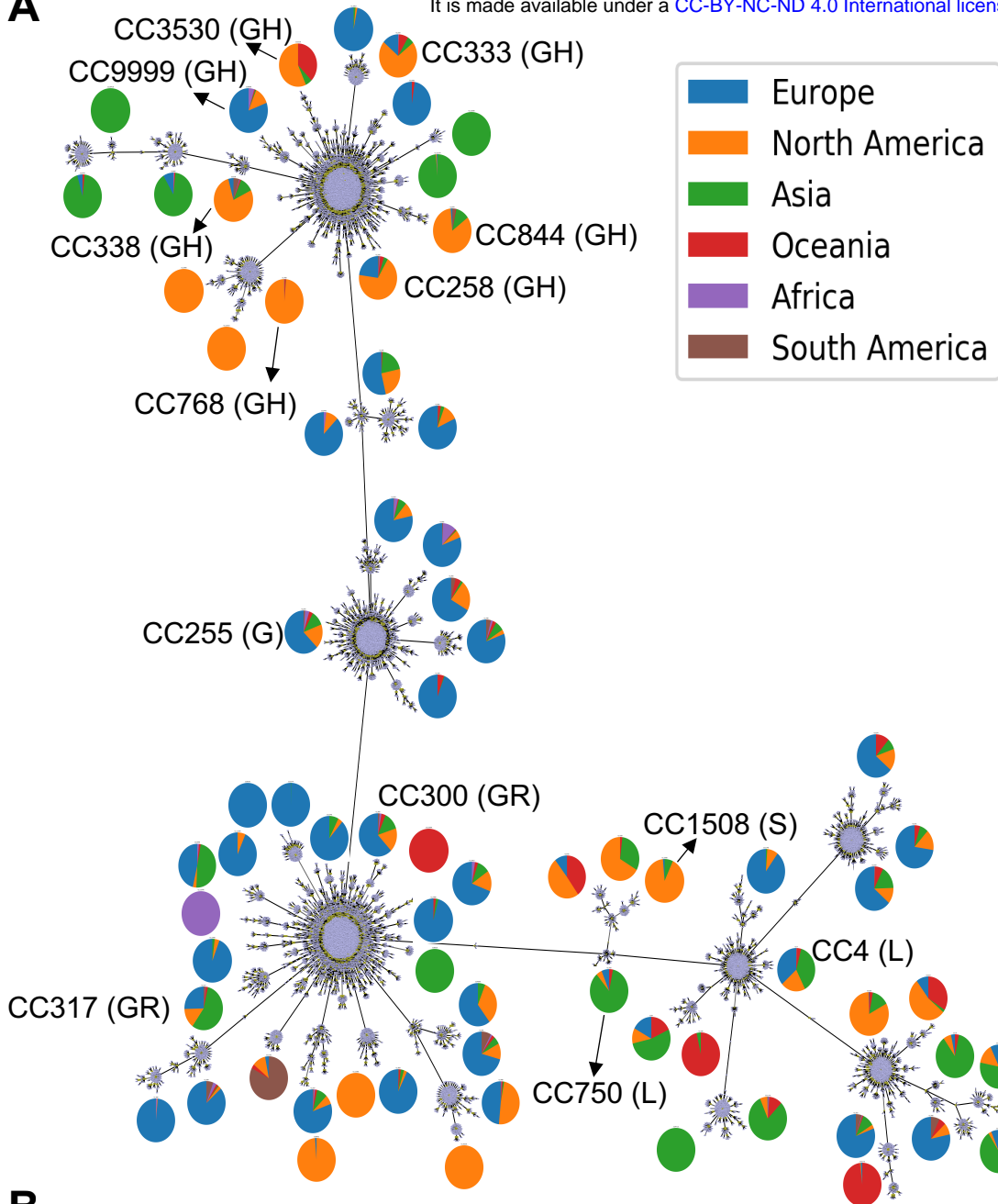432    (NY), New Jersey (NJ), Virginia (VA), Maryland (MD) and District of Columbia (DC).

433    Weeks 1, 2, 3, 4and 5 are from 03/19-03/25, 03/26-04/1, 04/02-04/08, 04/23-04/29 and

434    04/30-05/04, respectively. The GISAID clades corresponding to the CCs are reported in

435    parentheses.

436

437

438    **Figure 2.** SARS-CoV-2 diversity across different age groups in our sample. **A.** Relative

439    abundance of circulating CCs between pediatrics ($\leq$ 21 years old) and adults. **B.**

440    Relative abundance of circulating STs between children ($\leq$ 21 years old) and adults. **C.**

441    Relative abundance of circulating CCs in 5-year age groups. **D.** Relative abundance of

442    circulating CCs in childhood age ranges ($\leq$ 21 years old). Relative abundance is the

443    ratio of the number of genomes belonging a certain CC (lineage) divided by the total

444    number of genomes in a certain time window. The numbers on the bars represent the

445    total number of genomes in each group.