

# Real Time Viral Sub-Strains Discovery in Emerging Infectious Disease Situation – The African Perspective

Moses Effiong Ekpenyong (✉ [mosesekpenyong@uniuyo.edu.ng](mailto:mosesekpenyong@uniuyo.edu.ng))

University of Uyo <https://orcid.org/0000-0001-6774-5259>

**Faith-Michael Uzoka**

Mount Royal University

**Mercy Edoho**

University of Uyo

**Udoinyang G. Inyang**

University of Uyo

**Ifioek J Udo**

University of Uyo

**Nseobong P. Uto**

Saint Andrews University

**Itemobong S. Ekaidem**

University of Uyo

**Anietie E. Moses**

University of Uyo

**Enoabasi D. Anwana**

University of Uyo

**Youchou M. Tاتفeng**

Niger Delta University

**Geoffery Joseph**

University of Uyo

**Emmanuel A. Dan**

University of Uyo

---

## Research article

**Keywords:** cognitive mining, genome expression pattern, open-source framework, COVID-19, SARS-CoV-2 sub-strains diversity, infectious disease, intelligent prediction, self-organizing map, transmission pathway

**Posted Date:** January 19th, 2021

**DOI:** <https://doi.org/10.21203/rs.3.rs-147634/v1>

**License:**  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

---

1 **Real Time Viral Sub-Strains Discovery in Emerging Infectious Disease Situation – The**  
2 **African Perspective**

3  
4  
5  
6  
7  
8

Moses Ekpenyong<sup>1\*</sup>, Faith-Michael Uzoka<sup>2</sup>, Mercy Edoho<sup>1</sup>, Udoinyang Inyang<sup>1</sup>, Ifiok Udo<sup>1</sup>, Nseobong Uto<sup>3</sup>,  
Itemobong Ekaidem<sup>4</sup>, Anietie Moses<sup>4</sup>, Enobasi Anwana<sup>5</sup>, Youchou Tاتفeng<sup>6</sup>, Geoffery Joseph<sup>1</sup>, Emmanuel  
Dan<sup>1</sup>, Juliana Ndunagu<sup>7</sup>

9 <sup>1</sup>Department of Computer Science, University of Uyo, Nigeria  
10 <sup>2</sup>Department of Mathematics and Computing, Mount Royal University, Canada  
11 <sup>3</sup>School of Mathematics and Statistics, University of Saint Andrews, United Kingdom  
12 <sup>4</sup>College of Health Sciences, University of Uyo, Nigeria  
13 <sup>5</sup>Department of Botany and Ecological Studies, University of Uyo, Nigeria  
14 <sup>6</sup>College of Health Sciences, Niger Delta University, Nigeria  
15 <sup>7</sup>Department of Computer Science, National Open University, Nigeria

16  
17 \* – Corresponding Author  
18

19 **Abstract**

20 **Background:** The increased number of accessible genomes has prompted large-scale comparative studies for  
21 discerning evolutionary knowledge of infectious diseases, but challenges such as non-availability of close  
22 reference sequence(s), incompletely assembled or large number of genomes, preclude real time multiple  
23 sequence alignment and sub-strain(s) discovery. This paper introduces a cooperatively inspired open-source  
24 framework, for intelligent mining of severe acute respiratory syndrome coronavirus-2 (SARS-CoV-2) genomes.  
25 We situate this study within the African context, to drive advancement on state-of-the-art, towards intelligent  
26 infectious disease characterization and prediction. The outcome is an enriched Knowledge Base, sufficient to  
27 provide deep understanding of the viral sub-strains’ identification problem. We also open investigation by  
28 gender, which to the best of our knowledge has been ignored in related research. Data for the study came from  
29 the Global Initiative on Sharing All Influenza Data database (<https://gisaid.org>) and processed for precise  
30 discovery of viral sub-strains transmission between and within African countries. To localize the transmission

31 route(s) of each isolate excavated and provide appropriate links to similar isolate strain(s), a cognitive solution  
32 was imposed on the genome expression patterns discovered by unsupervised self-organizing map (SOM)  
33 component planes visualization. The Freidman-Nemenyi's test was finally performed to validate our claim.

34

35 **Results:** Evidence of inter- and intra-genome diversity was noticed. While some isolates (or genomes) clustered  
36 differently, implying different evolutionary source (or high-diversity), others clustered closely together,  
37 indicating similar evolutionary source (or less-diversity). SOM component planes analysis revealed multiple  
38 sub-strains patterns, strongly suggesting local- or intra-community and country to country transmissions.  
39 Cognitive maps of both male and female isolates revealed multiple transmission routes. Freidman's test results  
40 showed highly significant difference ( $p < 0.01$ ) among the various isolate groups. Nemenyi's test revealed groups  
41 that differed in their isolates.

42

43 **Conclusion:** The proposed framework offers explanations to SARS-CoV-2 diversity and provides real time  
44 identification to disease transmission routes, as well as rapid decision support for facilitating inter- and intra-  
45 country contact tracing of infected case(s). Intermediate data produced in this paper are helpful to enrich the  
46 genome datasets for intelligent characterization and prediction of COVID-19 and related pandemics, as well as  
47 the construction of intelligent device for accurate infectious disease monitoring.

48

49 **Key words:** cognitive mining, genome expression pattern, open-source framework, COVID-19, SARS-CoV-2  
50 sub-strains diversity, infectious disease, intelligent prediction, self-organizing map, transmission pathway.

51

## 52 **Background**

53 The coronavirus disease 2019 (COVID-19)—the disease caused by the severe acute respiratory syndrome  
54 coronavirus-2 (SARS-CoV-2), is a data-driven pandemic, because massive data and information are constantly  
55 being released and shared at very unprecedented scale. Artificial Intelligence (AI) has recently aided this  
56 process, as a growing amount of data and tools are constantly being explored to repurpose traditional approaches  
57 toward intelligent techniques, for early disease detection, real time contact tracing, new diagnostic tools  
58 development, informed policy formulation and implementation, and swift public health response, recovery and  
59 resilience. Taxonomically, SARS-CoV is one of the 36 coronaviruses in the family Coronaviridae within the  
60 order Nidovirales. Members of this family are known to cause respiratory or intestinal infections in humans and

61 other animals. However, despite a marked degree of phylogenetic divergence from other known coronaviruses,  
62 SARS-CoV together with bat SARS-CoV are now considered group 2b beta-coronavirus<sup>[1]</sup>. It has since been  
63 confirmed that two strains (the L- and S-strains) of the new coronavirus are spreading around the world today<sup>[2]</sup>,  
64 and the fact that the L-type is more prevalent suggests that it is “more aggressive” than the S-strain. To date, the  
65 most widely spread is the G-strain, while the L- and S-strains are fast disappearing<sup>[3]</sup>. Although a lot is yet to be  
66 understood about SARS-CoV-2, we know that it is among the six known species that can infect humans, with  
67 one of the species assuming two different genetic variants, making a total of seven species. Out of these species,  
68 four produce generally mild symptoms, while the remaining three produce potentially severe symptoms<sup>[4]</sup>.

69 A greater proportion of research progress on SARS-CoV-2 utilize the biotechnology dimension<sup>[5,6]</sup>, strictly  
70 focusing on species characterization and variants analysis through features extraction. However, AI and  
71 Machine Learning (ML) methods are expanding biotechnology capacity into the bioinformatics realm, through  
72 further features probing for precise classification and prediction. Nonetheless, most of these methods lack  
73 cooperatively inspired flavors to drive intelligent solutions to the problem. In general, AI/ML research on  
74 SARS-CoV-2 has permeated four key areas of healthcare services, discussed within the African context as  
75 follows:

76 *Screening and treatment*—Real time reverse transcriptase polymerase chain reaction (RT-PCR) appears the “gold  
77 standard” for detecting SARS-COV-2 in Africa<sup>[7]</sup>, but prolonged latency is associated with this test from when  
78 samples are collected to when they are processed, for results to be made available. Although serological tests  
79 that detect specific antibodies to the virus appear a quicker alternative to the RT-PCR test, they are not  
80 completely free of limitations<sup>[8]</sup>. Recently, several rapid diagnostic tests (RDTs) have been developed and tested  
81 in single studies, but none has so far been validated or commercially available.

82 *Contact tracing*—Placing contact tracing in the hands of all is certain to enhance contact identification<sup>[9]</sup>. Model-  
83 based approaches<sup>[10]</sup> as well as smartphone-enabled approaches<sup>[11]</sup>, constitute some recent methods for managing  
84 the COVID-19 spread. These methods are cost-effective, given the pressure on healthcare systems in a global  
85 pandemic, but however, incipient in Africa.

86 *Prediction and forecasting*—Combination of clinical indicators are known to predict symptoms of infected  
87 persons. Hence, studies exploit predictive and forecasting systems<sup>[12-14]</sup>, to effectively manage the  
88 spread/transmission of SARS-CoV-2. Cabore et al.<sup>[15]</sup> combined virus transmission characteristics and country  
89 specific socioecological factors to predict the most likely outcome of widespread and sustained community  
90 transmission of SARS-CoV-2. Using Markov chain model, with the transition states and country specific

91 probabilities derived from available knowledge, they found that sparsely populated cities, including Algeria,  
92 South Africa and Cameroon showed high risks of exposure. Nigeria was predicted as having the largest number  
93 of infections, followed by Algeria and South Africa. Mauritania was predicted to have the least cases of  
94 infection, followed by Seychelles and Eritrea. Sun et al.<sup>[16]</sup> assuming conservative estimates of African cases,  
95 utilized the reported number of COVID-19 cases traced to 12 major countries in Europe and America to  
96 Singapore, as well as flight data, to estimate the number of cases imported into Africa. They propagated the  
97 uncertainty in imported case count estimates to simulate onward spread of the virus, until 10000 cases were  
98 reached. The sensitivity of their results was accomplished using 1000 simulation runs under 4 different  
99 parameter combinations. They found that Morocco, Algeria, South Africa, Egypt, Tunisia, and Nigeria had the  
100 largest number of COVID-19 cases.

101 *Drugs and vaccine discovery*–The use of machine learning and simulations<sup>[17-19]</sup> as well as nanoparticles<sup>[20-21]</sup>,  
102 are helping the repurposing and discovery of drugs, for efficient drug delivery. Furthermore, biotechnology- and  
103 bioinformatic-based approaches such as genome sequencing in identifying potential therapeutic options for  
104 vaccine discovery and for controlling COVID-19, have heightened. Initially, most genomes of African cases  
105 were sequenced abroad, but recently we have witnessed the setting up of molecular laboratories to handle this  
106 process within. Some African countries have also intensified efforts on alternative remedies to COVID-19, but  
107 there is need for established standards on their efficacy and safety.

108 This paper proposes an open source framework for advancing infectious disease research in Africa. The  
109 specific objectives of the paper are:

- 110 • To exhume relevant literature and assertions on SARS-CoV-2, for proper examination of its prevalence  
111 and transmission.
- 112 • To perform complete genome analysis by gender, for proper understanding and inference of SARS-  
113 CoV-2 diversity, and genome expression patterns across African countries.
- 114 • To mine cognitive knowledge from discovered patterns, for efficient (contact) tracing of the viral  
115 transmission routes.
- 116 • To report our findings and corroborate existing literature/assertions.

117

118 The contributions of this paper to knowledge include:

- 119 • *Open Source Framework*–Most of the biotechnology and bioinformatic tools are ‘black boxes’ and not  
120 open to contributions by the research community. This paper therefore encourages reproducible

121 research by introducing a set of rapid prototype modules capable of generating intermediate results that  
122 provide further insights into the prevalence and transmission of the pandemic.

123 • *Effectual Tracing of Undocumented Source of Infection*–Community transmission of viral and anti-  
124 viral treatments could engender novel mutations in the virus, leading to potentially evolving sub-strains  
125 with high mortality resistance. Consequently, tracing the routes of infection for efficient documentation  
126 of COVID-19 cases is very essential. Unsupervised genome pattern clustering and cognitive modeling  
127 are achieved in this study, to explain the genome diversity of the SARS-CoV-2 sub-stains as well as  
128 provide real time solution to the disease transmission pathway.

129 • *Intelligent Genome Surveillance*–It has been observed that when this virus transmits from one person to  
130 another over few months, it may acquire random sequence variations of its genetic material which  
131 serves as distinctive genomic “fingerprints”. This paper enables the accurate mining of newly infected  
132 patients, to know which sub-strain of the virus is spreading within a country or been acquired from a  
133 different country. By combining machine learning techniques with cognitive knowledge mining,  
134 hidden sub-strains are revealed and different expression patterns followed, for seamless navigation of  
135 specific disease prevalence and surveillance of global infectious disease threats.

136 • *Misinformation/Disinformation Management*–Establishing transmission pathways would help  
137 minimize the growing trend of misinformation/disinformation, as country specific/global transmissions  
138 and spread of the virus could easily be contained. This paper guarantees the identification of possible  
139 infection routes by comparing genome sequences from different locations, to discover genetic diversity  
140 among sub-strains, and future potentials for investigating its fatal nature and spread.

141 • *Inputs to Novel Vaccine Development*–Understanding infection and transmission pathways could  
142 provide meaningful contributions to vaccine development and the discovery of clinically active  
143 variants and prototype drugs/vaccines for curative purposes. This paper does not only discover the  
144 SARS-CoV-2 sub-strains but also computes dissimilarity/variability in emerging sub-strains—an  
145 essential variable for vaccine development. Probing underlying genetic variations of infected  
146 individuals by gender would certainly enhance comprehensibility of the viral strain patterns, impact on  
147 the affected cells and aid the development of both preventive and therapeutic vaccine prototypes for the  
148 disease.

149

150

## 151 **SARS-COV-2: Existing Assertions and Developments**

152 On March 11, 2020, the World Health Organization (WHO) declared the coronavirus disease a global pandemic.  
153 Although the disease is still spreading, the rate of spread has greatly declined. Hence, almost all the countries  
154 had reopened their economies after compulsory national lockdowns, and currently adapting to local  
155 circumstances, with reduced rate of contact tracing and follow-up. Accelerated research developments and  
156 competing demands to contain the virus however have opened several opportunities for clinicians and  
157 researchers to exploit available avenues for developing suitable treatment and vaccines. Consequently, a  
158 plethora of publications flooded the scientific and medical domains/journals, with majority of the contributions  
159 received from the Asian countries and China—the very source of the pandemic (<https://clinicaltrials.gov>). Several  
160 studies and investigations have resulted in the following assertions and developments:

- 161 1) WHO claimed that most transmissions of COVID-19 are attributed to symptomatic persons than  
162 asymptotically infected persons, with asymptomatic persons practically incapable of transmitting or  
163 spreading the virus; but recently, persistent replication of SARS-CoV-2 variant has been derived from  
164 an asymptomatic individual<sup>[22]</sup>. Furthermore, whole genome sequencing of the persistently replicating  
165 strain shows diversity in nucleotide positions leading to 6 non-synonymous ORF1ab protein  
166 substitutions.
- 167 2) Confirmed cases of COVID-19 have surpassed those of SARS<sup>[22-23]</sup>. Its genetic diversity in most  
168 countries is similar to what obtains globally, suggesting repeated inter- and intra-country spread by  
169 infected persons rather than by “patient zero.” While some studies claim that mutation of new strains of  
170 SARS-CoV-2 potentially escalate severity of the pandemic<sup>[24]</sup>, further analysis have confirmed  
171 premature conclusions—as there is currently not enough evidence to support the claim that mutation  
172 significantly impacts spread of the virus.
- 173 3) Non-pharmaceutical interventions including physical distancing, isolation, and the use of mask are the  
174 best approach to contain the outbreak and may assist flatten the peak in communities. However, the  
175 challenge of compliance resulting in the alleged fear of increased number of infection, especially in low  
176 and medium income countries or resource limited settings, such as Africa, remains an unresolved  
177 puzzle, as poor health facilities and confusable symptoms continue to becloud the true evidence of  
178 infected cases.
- 179 4) The question of how Africa has survived the COVID-19 surge thus far may lie in the herbal remedies  
180 that abound within the continent’s biodiversity-rich ecosystems<sup>[25]</sup>, widely used in most African

181 communities and typified by the recently announced Madagascan COVID-19 remedy<sup>[26]</sup>. While the  
182 unsubstantiated remedy still requires medical scrutiny to prove its efficacy by globally acclaimed  
183 standards stipulated by the WHO, it is also a pointer that Africa may be in a position to provide  
184 alternative solution to disease management in moments of distress such as the present pandemic.

185 5) Amid conspiracy theories, it has since been inferred that SARS-CoV-2 is not a laboratory engineered  
186 virus but a natural process, after a comparative analysis of the SARS-CoV-2 genomic data and related  
187 (reference) viruses was conducted—as the distinct features of mutation in the receptor-binding domain  
188 portion of the virus spike protein usually targets the outer cell (of humans) involved in regulating blood  
189 pressure; and the lack of evidence of the virus being engineered from previously known viruses,  
190 debunk the notion of SARS-CoV-2 from being biologically engineered.

191 6) All the three human CoVs (SARS, MERS and SARS-2) are the result of recombination among  
192 CoVs<sup>[27]</sup>, as recombination has been found to affect patterns of common variants as well as  
193 substitutions.

194 7) Like SARS-CoV and MERS-CoV, SARS-CoV-2 appears to be a zoonotic virus which is transmitted to  
195 humans through animals such as bats, because genomic sequences of SARS-CoV-2 isolates from  
196 patients share significant sequence identity with very high degree of certainty that suggests a host shift  
197 from bats into humans<sup>[14][28]</sup>.

198 8) Clinical specimens used for viral ribonucleic acid (RNA) detection of COVID-19 as reported in the  
199 literature include nasopharyngeal aspirates, throat and nose swabs, saliva, sputum, endotracheal  
200 aspirates, feces, and urine. Of these, saliva yields greater detection sensitivity and consistency with  
201 high viral load concentration<sup>[29-32]</sup>.

202 9) At present, the sensitivity of clinical nucleic acid detection appears limited without clear pointers to  
203 genetic variation. However, studies such as<sup>[29]</sup> specifically identified nucleotides at different sites to  
204 infer genotypic/genomic variants of SARS-CoV-2, hence, suggesting multiple outbreak and source of  
205 transmission. Further, the presence of more samples in certain sites may indicate increased  
206 transmissibility.

207 10) Emerging trend of the virus may impact human health outcomes, demanding close monitoring and  
208 characterization of the viral genetic patterns. However, this view has opened series of inconclusive  
209 debates, with many scientists arguing that the prevalence of genetic mutations could have increased as  
210 a result of random (stochastic) processes without increased fitness. A more formal analysis of the

211 frequency of mutation recently suggests decreased transmissibility and the fact that the position of the  
212 spike protein does not reside within the receptor-binding domain, nullifies existing notions that  
213 mutation confers greater transmissibility.

214 11) Although majority of the mutations arising from viral replication have shown very negligible effect on  
215 the virus, with no possibility of infection; analysis of mutations in the spike protein of SARS-CoV-2  
216 suggests increased mutation frequency<sup>[30]</sup>. However, mutation information is appropriate to track new  
217 variants of the virus with unique mutant genomes, improve understanding of transmission and quicken  
218 determination of whether new mutations are changing the virus properties.

219 12) The presence of near real time whole-genome sequence analysis has provided reliable assessments on  
220 the extent of SARS-CoV-2 transmission in communities, hence, facilitating early decision making to  
221 control the local spread of the virus.

222 13) The sudden appearance of various sub-strains of the virus may not be unconnected with the fact that the  
223 virus is influenced by the new physical or biochemical environment it finds itself and/or in its ability to  
224 adapt to such a new and changing environment. Consequently, studies have successfully traced the  
225 SARS-CoV-2 of infected patients using molecular and phylogenetic methods<sup>[33]</sup>—as most phylogenetic  
226 inferences substantially prove that the virus has evolved into several sub-strains or variants specific to  
227 regions of transmission. Some studies have also shown high similarity between strains in different  
228 countries—as genotyping analysis of SARS-CoV-2 isolates around the globe reveals that specific  
229 multiple mutations are predominant during similar pandemics. Hence, comparing genome sequences  
230 from different locations allows for the analysis of the genetic diversity among viral sub-strains, its fatal  
231 nature, pathogenicity, origin and spread.

232 14) Although people of all ages are prone to infection by this virus, elderly people with co-morbidities  
233 (underlying health conditions and compromised immune system) are more susceptible to severe  
234 infection and death. Presence of genetic variants among young men with severe COVID-19 have been  
235 confirmed in<sup>[34]</sup>—using whole-exome sequencing performed to identify potential monogenic cause. But  
236 uncertainty did set in among medical practitioners on whether COVID-19 is a viral disease or the  
237 response to a person’s immune system that invariably damages a patient’s organs. Also, confusion in  
238 treating diseases presenting COVID-19 symptoms, instigated difficulty for physicians to determine  
239 with confidence, the optimal means of caring for critically infected patients. Howbeit, available data

240 informs the role of immune system in either diminishing or aggravating the infection and optimal  
241 measures for resolving confusable symptoms.

242 15) Confidence in how to treat COVID-19 has tremendously grown, but uncertainty remains<sup>[35]</sup>. At the  
243 outset of the pandemic, there appeared to be no definite treatment and the fear as to whether physicians  
244 themselves would get sick griped almost all the health providers/centers, the world over; as some  
245 diagnosed with the virus were asymptomatic (showed no symptom). Currently, most COVID-19  
246 patients now have mild symptoms; but two important questions still linger: Will there be a next phase  
247 of the pandemic? Has most of the various communities suddenly reached “herd immunity”?

248 16) Development of high-throughput sequencing has contributed high quality datasets including whole  
249 genome sequences of viral isolates to the public domain. Analysis of genome sequence also provides  
250 insights into global spread patterns, genetic diversity, as well as the dynamics of sub-strains evolution.  
251 With continuous availability of new data, deeper investigation into new methods towards efficient  
252 candidate vaccines discoveries for emerging and re-emerging COVID-19 and related pandemics is  
253 ongoing.

254 17) Emergency lockdowns have returned to parts of Europe, as France and Germany struggle with a second  
255 wave while the COVID-19 outbreak resurges rapidly. The United States is experiencing a third wave,  
256 and the UK is set to impose another lockdown. As this uncertainty window reopens, deaths are flat,  
257 with cross-border countries becoming apprehensive.

258  
259

## 260 **SARS-COV-2 Genome Analysis of African Isolates**

261 In this section, we review existing works on SARS-CoV-2 analysis conducted on African genomes and present  
262 in Table 1, a summary of the viral isolates, their transmission history, intra-country sub-strains discovery and  
263 additional information about local transmission, mutation and spread.

264 *Egypt:* Sekizuka et. al.<sup>[36]</sup> characterized the possible origin of 10 SARS-CoV-2 positive travelers from  
265 Egypt together with their close contacts. The viral genome sequences of the 10 travelers were aligned with  
266 genome sequence retrieved from GISAID using MAFFT v7.222; two distinct genome lineages circulating  
267 mostly in Europe and South America were identified. They concluded that increased cases may complicate the  
268 identification of infection routes. The analysis and comparison of 2 Egyptian SARS-CoV-2 isolates using CLC  
269 Genomic Workbench version 20<sup>[37]</sup> yielded at least 99.9% similarity. However, variations occurred at 8658,

270 15907, 19906 and 18877 nucleotide sites. Comparable with the Wuhan reference genome, 5 mutations (C->T:  
271 C241T, C->T: C3037T, C->T: C14408T, A->G: A23260G and G->T: G25563T) were observed among the  
272 Egyptian sequences. Specifically, the genome sequence of patient 1 was discovered to shared similarities with  
273 Taiwanese isolate who traveled in February 2020 to Dubai and Egypt, few mutations were found at sites (C->G:  
274 A8658G, G->A: G15907A and C->T: C18877T). Patient 2 recorded similarities but with 3 specific variations  
275 (T->C: T4278C, G-T: G18963T, C26692) with the genome sequence of a Japanese on a Nile river ship cruise,  
276 who tested positive to SARS-CoV-2 on March 9, 2020. Amidst all the variations, the spike protein of the 2  
277 Egyptian SARS-CoV-2 isolates are identical with G614D but varies from the spike protein of the Wuhan  
278 reference genome. Further investigation involved the nucleotide alignment of the Egyptian sequence with other  
279 GISAID SARS-CoV-2 genome sequences using BioEdit version 7.0.5.3 and ClustalW. Phylogenetic tree  
280 constructed using MEGA7, showed the clustering of Egyptian sequences in clade A2a with Asian Europe,  
281 United States, Australian and African sequences.

282 *Kenya:* The first reported case of SARS-CoV-2 sequencing and analysis in Kenya consisted of positive  
283 samples from symptomatic and asymptomatic patients from Nairobi (20) and Coastal Kenya (102). Seventy-  
284 eight global sequences representing countries in Europe, Asia, America and Africa in GISAID were randomly  
285 sampled and retrieved. The alignment of the retrieved sequences together with the Kenyan sequences was  
286 realized using MAFFT v7.310. A maximum likelihood phylogenetic tree was established through RAXML-  
287 NGS v0.9.0 using a GTR+F0+G4m model in 1000 bootstraps run. The assignment of lineages to the Kenyan  
288 genome sequences was possible through PANGOLIN toolkit (v1.1.14). Evidently, the phylogenetic tree  
289 displayed 10 strains of SARS-CoV-2 circulating in Kenya, which evinces the multiple introduction of the virus  
290 into Kenya. Nonetheless, B.1 lineage discerned to be dominant and causing most of the infections in Coastal  
291 Kenya. However, all the viral strains were identifiable with the strains circulating globally and none of the  
292 strains was distinctive to Kenya.

293 *Morocco:* Laamarti et al.<sup>[38]</sup> studied the molecular distribution of 28 Moroccan SARS-CoV-2 strains  
294 isolated between March 3, 2020 and May 15, 2020, with 12 North African (Tunisia (7), Algeria (3) and Egypt  
295 (2)) viral sequences downloaded from GISAID and 6 Moroccan genome sequenced for the study. Specifically, 6  
296 sequences were mapped to the Wuhan reference genome using BWA-MEM v0.7.17-r1188, while Minimap  
297 v2.12-r847 was used in mapping the GISAID downloaded genomes. The analysis of all Moroccan genome  
298 sequences disclosed 61 mutations in comparison with the Wuhan genome: 27 synonymous, 5 intergenic, 27  
299 non-synonymous and 2 lost stops. These mutations were distributed among 5 genes (ORF1ab, S, M, N and

300 ORF3a); ORF1ab harbored the highest number (37.7%) of non-synonymous mutations. In like manner, the  
301 comparison and characterization of the 12 North African genome sequences together with the 28 Moroccan viral  
302 genome sequences, against the Wuhan genome, revealed a total of 118 mutations: 58 non-synonymous, 48  
303 synonymous and 12 inter-gene mutations. Regarding non-synonymous mutations, missense, lost stop and stop  
304 again were found to contribute 91.38%, 6.90% and 1.72%, respectively. Amongst the 58 non-synonymous  
305 mutations, 13 were repeatedly found in more than one genome. The most recurrent of the mutations observed  
306 within the four North African countries occurred in the S protein (D614G) and ORF3a (Q57H) with prevalence  
307 of 92.5% and 42.5%, respectively. The 11 (T265I, T5020I, K2798R, R203K, D1036E, V2047F, A2637V,  
308 T2648I, C4588F, S202N, L84S) other mutations were inconsistently observed among the four countries. Also,  
309 in addition to the 5 genes earlier discovered to harbor mutations, 2 more genes (E and ORF8) were discovered.  
310 However, ORF1ab remained the leading gene with 67.4% mutations, bearing two-third of the 118 mutations. In  
311 addition, with a focus on Morocco, the phylogenetic analysis conducted using 256 representative genomes  
312 constituting genome sequences from the 6 continents, the phylogenetic tree disclosed 5 major clades, 2 of which  
313 constituted main strains from Asia while the other clades consisted of strains belonging to various continents.  
314 This diversity in lineage infers the introduction of SARS-CoV-2 into Morocco from multiple routes. In<sup>[39]</sup>, the  
315 molecular analysis of SARS-CoV-2 genome sequences of 22 Moroccan isolates obtained from three laboratories  
316 in Morocco as at June 7, 2020 revealed 62 mutations. In comparison with the Wuhan reference genome and  
317 40366 viral genome sequences retrieved from GISAID, the Moroccan genome evinced similar mutations with  
318 the Wuhan reference genome and other strains circulating globally. Additional 6 mutations (NSP10\_R134S,  
319 NSP15\_D335N, NSP16\_1169L, NSP3\_L431H, NSP3\_P1292L and Spike\_V6F) particular to the Moroccan  
320 SARS-CoV-2 genome were also discovered. Their study was realized by performing MAFFT multi-alignment  
321 of all retrieved genomes from GISAID and phylogeny analysis of the aligned sequences created by maximum  
322 likelihood using IQTREE. The evolutionary analysis revealed 3 clades: 20A, 20B and 20C, and authenticates  
323 the findings of<sup>[40]</sup>, which used similar methodology to investigate the phylogenesis of 250 SARS-CoV-2 genome  
324 sequences from GISAID. Sixteen variants were detected in 6 Moroccan SARS-CoV-2 genome sequences.  
325 Among the variants were synonymous (F924F and L4715L), nonsynonymous (D->G: D614G) and intergenic  
326 (241C->T). Jouali et al.<sup>[41]</sup> corroborate the inference of<sup>[39]</sup> by comparatively studying the SARS-CoV-2 genome  
327 sequence of a mildly symptomatic Moroccan patient with other sequences from Morocco. The phylogenetic  
328 analysis of the genome was conducted using GISAID enabled Nextstrain tool. The genome under study was  
329 discovered to belong to clade B11 revealing high similarity with genome sequences from Florida, USA.

330 *Nigeria:* The comprehensive knowledge of phyloevolution and comparative discrimination of SARS-CoV-2  
331 molecular characterization can be useful in the critical investigation of the virus pathogenesis, disease control,  
332 treatment and vaccine development. The study of SARS-CoV-2 evolution in Nigeria by<sup>[42]</sup> exhibited a concerted  
333 similarity with the Wuhan reference genome, which introduction into Nigeria was inferred to arrive from  
334 Wuhan through an Italian traveler. The study involved a 2-step analysis of multiple sequence alignment and  
335 phylogenetic analysis of 39 complete genome of SARS-Cov-2 with their various travel history. The constructed  
336 phylogenetic tree separated the Nigerian strains into the cluster with a Wuhan subclade. Multiple sequence  
337 analysis using ClustalW revealed >70% similarity with the Wuhan reference sequence. In another study,  
338 representative whole genome sequences of each of the seven lineages of human SARS-CoV-2 circulating in  
339 Nigeria were obtained from GISAID and aligned with all full genomes from Nigeria using MAFFT v7.310<sup>[43]</sup>. It  
340 was found that 4 of the new sequences clustered closely together forming separate clade, strongly suggesting  
341 local community transmission. Similarly, other new sequences behaved in same way, revealing a follow-up  
342 from same patients. Inter-country analysis of lineages from Nigeria clustered with sequences from Asia, Europe,  
343 USA, Middle East, Australia, and other African countries, indicating multiple transmissions.

344 *South Africa:* From the consensus genomic sequence of a South African isolate, who travelled back South  
345 Africa from Italy, Allam et al.<sup>[44]</sup> identified 6 non-synonymous variants. This was attained using MAFFT v7.042  
346 from the multiple sequence alignment of 965 SARS-CoV-2 genome sequence, extracted from GISAID together  
347 with the isolate's sequence. As at time of the report, the mutations at location 13,620 bp and 21,595 bp were  
348 reported to be absent in every other SARS-CoV-2 genome. Furthermore, DUET web server prognosticated the  
349 destabilizing and stabilizing effect of D614G variant of the spike protein and P322L mutation of the nsp12  
350 respectively. As a supplemental study to the obstacles encountered during near-real time SARS-CoV-2  
351 genotyping during pandemic, Pillay et al.<sup>[45]</sup> established 3 lineages (B, B.1 and B.2) from a phylogenetic  
352 analysis of 54 near-full-length genome, using Phylogenetic Assignment of Named Global Outbreak LINEages  
353 (PANGOLIN) software suite. The 54 genome sequences and 10,959 GISAID reference genomes were aligned  
354 using MAFFT v7.313 and a maximum likelihood tree topology constructed in IQ-TREE. From the 54 genome  
355 sequences, lineage B.1 had the highest number of 50 samples clustering closely with 99.99% similarities.  
356 Moreover, the number of minimal mutations reported was quite minimal and not too divergent from those found  
357 in the public sequences.

358 *Uganda:* In Uganda, Bugembe et al.<sup>[46]</sup> reported on the genomic sequences of 14 travelers from SARS-CoV-  
359 2 dense region and 6 truck drivers returning to Uganda from Kenya, Tanzania and South Sudan. Phylogeny of

360 Ugandan genome sequences identified with 6 lineages (A, B, B.1, B.1.1, B1.1.1 and B4) was performed by  
361 comparing with globally detected genomes. The B.1 lineage accommodated greater number of sequences  
362 circulating in more than 20 countries in Europe, America and Australia and Asia. Although infection routes were  
363 mostly from the cargo truck drivers and air travelers, the viral genome of travelers from Dubai associated with 3  
364 lineages (A, B, B.1.1.1) while genomes from the cargo truck drivers entering Uganda from Tanzania belonged to  
365 lineage A and B.1. Whilst the sample size is small, the evidence suggests multiple sources of contact. Furthermore,  
366 the diversity of the Uganda genome sequence from the Wuhan reference genome occurred at 5-20 nucleotide  
367 positions across approximately 30kb genome, including the spike protein. In the spike protein, four viral sequences  
368 from lineage A encode D614, whereas sequences belonging the other clades encode G614.

369

370 [Table 1, here]

371

## 372 **Methods**

373 Studies on SARS-CoV-2 single nucleotide polymorphism and lineage discovery keep surging the literature;  
374 mainly exploiting a two-step algorithm (multi sequence alignment and phylogenetic analysis), with the use of  
375 common tools and techniques such as MAFFT and maximum likelihood tree topology that target specific  
376 nucleotide sites. Although aligning collected genomes against reference genome(s) has helped in the discovery  
377 of gene/genetic variability/ diversity, results of evolutionary analysis have consistently shown structured  
378 transmission with possible multiple introductions into the population<sup>[79]</sup>. Furthermore, most of the works on  
379 African genome isolates mine data from the Global Initiative on Sharing All Influenza Data: GISAID EpiFlu™  
380 (a database of SARS-CoV-2 partial and complete genome compilations distributed by clinicians and  
381 researchers, the world over).

382 To advance the current practice, an open source framework that combines biotechnology and bioinformatic  
383 approaches with AI methodologies, into a hybridized system, is proposed in this section, for in-depth  
384 understanding and further works development on SARS-CoV-2. Our framework generates interesting  
385 intermediate data including phylogenomic trees, pairwise nucleotide similarity matrices/scores, gene diversity  
386 plots, genome expression patterns analysis, essential for enriching the genome datasets, towards intelligent  
387 genome characterization and prediction. With this approach, community contribution is guaranteed, and  
388 reproducible research possible. Furthermore, intermediate data could be repurposed for building new concepts

389 and models. The general workflow describing the proposed system framework is shown in Fig. 1, and algorithm  
390 implementing the framework presented in Algorithm 1.

391

392

393

394 **Fig. 1. Workflow describing the proposed system framework.** The workflow begins with the excavation of FASTA files of human SARS-  
395 CoV-2 genome sequences from GISAID. These files are stripped and processed into a genome database (DB) as multiple columns of  
396 nucleotide sequence. A series of AI/ML techniques are applied to extract knowledge from the genome datasets as follows: Compute  
397 dis(similarities) scores between the various pairs of genome sequences and obtain a genomic tree of highly dis(similar) isolates grouped in the  
398 form of a dendrogram/phylogenomic tree. Determine the optimal number of natural clusters—to provide additional knowledge. Separate the  
399 viral sub-strains using self-organizing map (SOM) component planes—for transmission pathways visualization. Perform direct pairwise  
400 nucleotide alignment of the entire genome sequences—to yield a nucleotide similarity matrix. Generate cognitive map—for intelligent sub-  
401 strains contact tracing and prediction.

402

403 **Algorithm 1. Steps implementing the workflow in Fig. 1.**

---

```
1. import necessary libraries
2. set path to current directory
3. #Genome nucleotide fragments processing
4. create a list of FASTA files (fasta_list) to process
5. for file_name in fasta_list:
6.     open FASTA_file for read
7.     store a line of genome sequence
8.     for line in file_name:
9.         strip line into a list of nucleotide fragments (nucleotide_fragments)
10.    for line in nucleotide_fragments:
11.        write nucleotide code into complete genome file (complete_genome)
12.    close FASTA_file
13. #Direct nucleotide alignment and similarity scores generation
14. open complete_genome for read
15. store a line of nucleotide code
16. for line in complete_genome:
17.    align nucleotide pair and compare nucleotide code
18. build pairwise dis(similarity) matrix using a suitable distance metric (e.g., Euclidean distance)
19. #AGNES/hierarchical clustering: generate phylogenomic tree and cluster plots
20. treat observations (nucleotides) as cluster points and compute AGNES distance coefficients between clusters
21. compute scores between genome isolates clusters
```

---

- 
22. build and visualize genomic tree
  23. discover and validate optimal natural clusters (k) using any k-means based N approaches (N>2) (elbow, silhouettes, gap-statistics etc.).
  24. partition the tree into k clusters
  25. #Genome expression patterns discovery
  26. perform SOM clustering on complete\_genome
  27. obtain SOM component planes of learned genome expression patterns
  28. obtain pairwise correlation coefficients
  29. label target (output) classes using dis(similarity) and genome expression clusters—indicating mutant sub-strains and viral expression patterns, to form enriched genome datasets.
  30. generate cognitive maps with embedded links of genome isolates.
  31. close complete\_genome.
- 

404

405

406

#### 407 **Data Source and Genome Sequences Selection**

408 Publicly available datasets of coronavirus deposited between December 2019 and October 12, 2020 were  
409 excavated from GISAID for the purpose of this study, and complete genome sequences of human SARS-CoV-2  
410 isolates from selected countries within Africa, collected. Useful metadata on the extracted genome sequences  
411 (Country, State, Abbreviation, Accession No., Length, Gender, Age, Travel history, Specimen source, Status,  
412 Submitting Lab, Authors), with Gender and Age removed, are documented (see SupplData1\_1.xlsx). The  
413 preprocessed FASTA files of genome isolates excavated from GISAID, striped and dumped as column sequences  
414 for male and female patients are found in (SupplData2\_1.xlsx). Although some of the datasets had incomplete  
415 information (e.g. age, specimen source and status), gender served as a compulsory criterion for profiling the  
416 excavated genomes. Male isolates excavated according to *Country (number of sequence, state(s) extracted)*,  
417 include: Algeria (2, 1); Benin (5, 2); Cameroon (1, 1); DRC(15, 4); Egypt (5, 2); Gambia (6, 2); Ghana (4, 1);  
418 Kenya (4, 1); Mali (5, 2); Morocco (5, 3); Nigeria (13, 5); Senegal (14, 6); South Africa (65, 20); Tunisia (1, 1).  
419 Female isolates excavated according to *country (number of sequence, state(s) extracted)*, include: Algeria (1, 1);  
420 Benin (4, 1); DRC(13, 4); Egypt (4, 2); Gambia (4, 2); Ghana (4, 1); Kenya (4, 1); Madagascar (2, 2); Mali (4, 1);  
421 Morocco (2, 1); Nigeria (11, 5); Senegal (11, 3); South Africa (80, 24); Tunisia (1, 1). Hence, a total of 290  
422 genome sequences (145 males, 145 females) with genome lengths of over 29000 nucleotides, were excavated.  
423 Specimen sources include swabs (nasal, oral, throat, nasal and oral); fluids (bronchoalveolar lavage, saliva,  
424 sputum) and unknown. Status of patients include hospitalized, not hospitalized, acute bronchitis, symptomatic,

425 asymptomatic, alive and unknown. Age range of 2 months and 99 years were considered. In the Male datasets the  
426 ages of 9 patients (Ghana (1), Kenya (4), Morocco (1), Mali (1) and Nigeria (2)), were unknown; while in the  
427 female datasets, the ages of 10 patients (Kenya (4), Morocco (2), and Nigeria (4)), were unknown. Finally, about  
428 1.75% and 2.11% of errors in sequencing (noise) were noticed in the male and female genome datasets,  
429 respectively.

430

431

### 432 **Unsupervised Genome Clustering**

433 Self-organizing map (SOM) has been used extensively in the field of bioinformatics, for visual inspection of  
434 biological processes, genes pattern expressions—as maps of (input) component planes analysis. SOM is an  
435 unsupervised artificial neural network (ANN), learned to produce a low-dimensional (typically two-  
436 dimensional), discretized representation of the training samples input space, known as a map. Patterns exhibited  
437 by the different isolates confirm intra- and inter-country transmissions. The SOM algorithm locates a winning  
438 neuron, its adjusting weights and neighboring neurons. Using an unsupervised, competitive learning process,  
439 SOMs produce a low-dimensional, discretized representation of the input space of training samples, known as  
440 the feature map (see Fig. 2). During training, weights of the winning neuron and neurons in a predefined  
441 neighborhood are adjusted towards the input vector using equation (1),

$$442 \quad w_{id}^{t+1} = w_{id}^t + rf(i, q)(x_d - w_{id}^t); 1 \leq d \leq D. \quad (1)$$

443 where  $r$  is the learning rate and  $f(i, q)$  is the neighborhood function, with value 1 at the winning neuron  $q$ ; and  
444 decreases as the distance between  $i$  and  $q$  increases. At the end, the principal features of the input data are  
445 retained, hence, making SOM a dimension reduction technique. The batch unsupervised weight/bias algorithm  
446 of MATLAB (*trainbu*) with mean squared error (MSE) performance evaluation, was adopted to drive the  
447 proposed SOM. This algorithm trains a network with weight and bias learning rules using batch updates. The  
448 training was carried out in two phases: a rough training with large (initial) neighborhood radius and large  
449 (initial) learning rate, followed by a finetuned training phase with smaller radius and learning rate. The rough  
450 training phase can span any number of iterations depending on the capacity of the processing device. In this  
451 paper, we kept the number of iterations at 200 with initial and final neighborhood radius of 5 and 2,  
452 respectively, in addition to a learning rate in the range of 0.5 and 0.1. The fine training phase also had a  
453 maximum of 1000 epochs, and a fixed learning rate of 0.2. Selection of best centroids of the genome feature  
454 within each cluster was based on the Euclidean distance criterion. The algorithm configures output vectors into a

455 topological presentation of the original multi-dimensional data, producing a SOM in which individuals with  
456 similar features are mapped to the same map unit or nearby units, thereby creating smooth transition of related  
457 genome sequences to unrelated genome sequences over the entire map.

458

459

460

461 **Fig. 2. SOM showing the map topology and interactions between nodes.** Each neuron is assigned a vector of weights ( $w =$   
462  $w_{i1}, w_{i2}, \dots, w_{iN}$ ) with dimension similar to the input vector  $i$  ( $i = 1, 2, \dots, L$ ); where  $L$  is the total number of neurons in the network. The input  
463 nodes have  $p$  features, and the output nodes,  $q$  prototypes, with each prototype connected to all features. The weight vector of the  
464 connections consumes the prototype of each neuron and has same dimension as the input vector. SOMs differ from other artificial neural  
465 networks as they apply competitive learning, against error correction learning such as backpropagation, and the fact that they preserve the  
466 topological properties of the input space using a neighborhood function.

467

## 468 **Cognitive Knowledge Mining**

469 Knowledge mining has served huge benefits for quick learning from big data. We applied Natural Language

470 Processing of the genome datasets to extract knowledge of similar strains of the virus. A simple iterative

471 technique is imposed on the SOM isolates ( $i = 1, 2, 3, \dots, n$ ), where  $n$  is the maximum number of isolates, as

472 follows: For each isolate pattern, compile similar patterns with the rest of the isolates (i.e.,  $i + 1, i + 2, \dots, n$ ).

473 Concatenate compiled patterns into a list ( $j_1, j_2, \dots, j_m$ ) where  $j$  is an element of the list. Dump the compiled list

474 into  $CogMap(k_i \in j_1, j_2, \dots, j_m)$ .

475

476

## 477 **Configuration of Computing Device**

478 A HP laptop 15-bs1xx with up to 1TB storage running on Windows 10 Pro Version 10.018326 Build 18362 was

479 used for processing the excavated genome sequences, algorithms/programs and other ancillary data. The system

480 had an installed memory (RAM) of 16 GB with the following processor configuration: 1.60 GHz, 1801 MHz, 4

481 Core(s) and 8 logical processors. Although our system performed satisfactorily and produced the desired results,

482 higher system configurations would improve the computational speedup.

483

## 484 **Results**

485 **Low Cytosine to Guanine Transition and High Thymine Content in Human SARS-CoV-2**

486 The RNA sequence is composed of 4 nucleotides (adenine (A), cytosine (C), guanine (G) and thymine (T)), also  
487 considered as polymers of 16 (i.e.  $2^4$ ) dinucleotides. Yin<sup>[54]</sup> revealed the frequency of mutations in the spike-  
488 protein, RNA polymerase, RNA primase and nucleoprotein. The study algorithm involved the alignment of  
489 multiple genome subsets with SARS-CoV-2 reference genome using Clustal Omega and Jaccard distance in the  
490 variation among 558 complete genome sequence retrieved from GISAID on March 23, 2020. Comparably,  
491 Wang et al.'s<sup>[55]</sup> homology analysis and sequence alignment established a reference genome from 95 strains  
492 obtained from NCBI and GISAID on February 14, 2020. The study as well, found mutations at nt8782 of  
493 ORF1a, nt28144 of ORF8 and nt29095 of N region from the analysis at the nucleotide and amino acid levels.

494 The SARS-CoV-2 reference genome (29903 nucleotides<sup>[56]</sup>, sequence number NC\_045512) consists of  
495 29.94% of A, 18.37% of C, 19.61% of G and 32.08% of T nucleotides<sup>[57]</sup>. Hence, the expected frequency of CG  
496 dinucleotide in the viral genome is 3.60% (i.e.  $19.61\% \times 18.37\%$ ). Mercatelli and Giorgi<sup>[58]</sup> analyzed 48635  
497 complete genome sequences spread across geographic regions including, Africa (514), Asia (3340), Europe  
498 (31818), North America (10250), Oceania (2127) and South America (575). The obtained sequences were  
499 aligned over the Wuhan reference genome sequence (NC\_045512.2) using NUCMER version 3.1. The analyzed  
500 result exemplifies the nature of mutation across the world, per continent and per country. The number of  
501 mutations per sample was reported to be relatively low but with an average mutation rate of 7.23. Although  
502 number of mutations per continent did not differ significantly from the average mutation rate, the average  
503 number of mutations per country differed significantly. For two out of the three African countries included in  
504 the study; Congo had a high mutational burden of 8.30 while Kenya had a low mutation rate of 5.38. Single  
505 nucleotide polymorphism substitution accounted for 0.6% of all the observed mutations, making it more  
506 prevalent over insertion and deletion mutations. The transition from C->T makes up the 55.1% of all point  
507 mutation, A->G is the second leading transition (14.8%) globally and in Africa, America and Europe. The effect  
508 of A->G transition on the protein sequence of SARS-CoV-2 formulates the G-clade predominantly found in  
509 Africa, Europe, Oceania and South America. Sjaarda et al.<sup>[59]</sup> studied 25 SARS-CoV-2 genome samples from  
510 local cases of COVID-19 collected during the early days of the spread from eastern Ontario, Canada between  
511 March 18 and March 30, 2020, with 2 genomes belonging to the S-clade and the remaining 23 belonging to the  
512 G-clade of SARS-CoV-2; and contained 45 polymorphic sites with one shared missense and three unique  
513 synonymous variants in the gene encoding the spike protein. They found that most of the genomes had between  
514 6 to 8 variants when compared with the NC\_045512.2 reference genome. Also, the most common nucleotide  
515 substitution was from C->T (25/45 variants), followed by G->T (7/45 variants) and A->G (4/45).

516 From the genomes excavated for this study, the average frequency count of nucleotides for male and female  
517 isolates are roughly similar, as the proportion of each nucleotide has (A=29.8553%; C=18.3830; G=19.6366;  
518 T=32.1254) nucleotides for the male isolates and (A=29.3616; C=18.3765; G=19.6365; T=32.1254) nucleotides  
519 for the female isolates; with average frequency of 36.1% CG dinucleotide compared to 59.1% CT and 63.1 GT  
520 dinucleotides, in viral genome, for both genders. Our result corroborates the findings of<sup>[57]</sup> on CG reduction in  
521 SARS-CoV-2, achieved through C/G nucleotide mutating into A/T (a universally occurring process in all forms  
522 of life). Generally, the mutation spectrum of new genome mutants seems enriched in C->T and G->T mutations,  
523 as different studies also corroborate our findings of dominant transition and transversion mutants in human  
524 SARS-CoV-2 isolates<sup>[43][59][60]</sup>, but no strong evidence supports the claim that the virus rapidly or slowly  
525 mutates than expected, as most of the mutations are probably neutral or deleterious to the virus<sup>[61]</sup>. But while the  
526 T nucleotide is the most frequent nucleotide in the genome, its frequency seems to increase further across all  
527 samples, and the substitution process appears non-reversible and unbalanced.

528

### 529 **Genome Diversity Analysis**

530 Genes cover very large regions of chromosomes with most gene content having almost identical expressions  
531 with other chromosomes in the genome. We introduce the density plots, to examine whole genome sequences,  
532 for the discovery of variability in nucleotide distribution in male and female isolates, between and within  
533 countries.

534

535 **Inter-country Analysis:** Fig. 3 shows density plots for male and female isolates between African countries.  
536 Observe that the male isolates exhibit smoother distribution curve with most of the isolates having identical  
537 expression patterns; compared to female isolates, which distribution curve is partly influenced by dominant  
538 outliers of possible infested sequences from Gambia, Kenya, Mali, Morocco, Nigeria, Senegal and South Africa.  
539 The outliers may be as a result of observed sequencing errors, or extensive localized variation in DNA  
540 polymorphism and large regions of low gene density, diversity and recombination.

541

542

543

544 **Fig. 3. Density plots of entire genome datasets.** A Density plot visualizes the distribution of data over a continuous scale. It is a variation  
545 of the histogram that applies kernel smoothing to plot values, enabling smoother distributions by smoothening out the noise. The peaks of  
546 the density plot help display where values are concentrated over the interval.

547  
548  
549  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
560  
561  
562  
563  
564  
565  
566  
567  
568  
569  
570  
571  
572  
573  
574  
575  
576  
577  
578  
579

**Intra-country Analysis:** Density plots revealing genomic diversity between male and female isolates of selected countries with more genome samples are presented in Fig. 4, Fig. 5, Fig. 6 and Fig. 7, for DRC, Nigeria, Senegal and South Africa, respectively. We observe marked variabilities in some nucleotide sequences of Nigerian (Fig. 5), Senegalese (Fig. 6) and South African (Fig. 7) isolates—indicating gene expression patterns differences in some of the isolates.

**Fig. 4. Density plots of DRC’s isolates.** Both plots exhibit similar curve pattern with near aligned isolates, indicating identical genome expression between the isolates.

**Fig. 5. Density plots of Nigerian isolates.** Although both plots maintain same curve patterns, the male isolates exhibit diversity in most of the genomes. Two female isolates (Fig, 5b) present distinct variations from the rest of the isolates, indicating significant genome diversity between the male and female isolates.

**Fig. 6. Density plots of Senegalese isolates.** Aside the single isolate that fails to align with other male isolates (Fig. 6a), both plots maintain near-aligned isolates.

**Fig. 7. Density plots of South African isolates.** Both plots maintain same patterns with shades of similar isolates dominating the plot area.

### **Phylogenomic Analysis**

Genes content comparison has become commonplace, but associating its order is challenging. Phylogenomic trees appear not widely used because of computational difficulties (massive data, high processing cost and limited processing infrastructure). In this paper, we exploit complete genome sequences to construct hierarchical

580 cluster structures (dendrograms) that discriminate inter- and intra-country SARS-CoV-2 isolates. To achieve  
581 this, Hierarchical Clustering Analysis (HCA) also known as Agglomerative Nesting (AGNES)<sup>[62]</sup> was  
582 performed on the various isolates. While there are natural structural entities in some datasets that provide  
583 information on the number of clusters or classes, others including the dataset containing genome sequences are  
584 structured without boundaries. Cluster validation (an unsupervised methodology aimed at unravelling the actual  
585 count of clusters that best describes a dataset without any priori class knowledge) is therefore essential. This  
586 paper adopted the elbow criteria<sup>[63]</sup>, to validate the number of clusters available in the genome datasets. Hence,  
587 yielding 2 and 3 clusters for inter- and intra-country phylogenomic analysis, respectively.

588

589 **Inter-Country Analysis:** Fig. 8 shows phylogenomic trees for male and female patients. Both trees suggest  
590 inevitable sub-strains (independent) mutant accumulation in different countries, resulting in highly dense  
591 clusters (encircled in red), while few mild divergent strains with specific mutations are geographically different,  
592 hence, occupying smaller disparate clusters.

593

594

595

596 **Fig. 8. Phylogenomic trees for African male and female isolates.** For full names of country codes, see Additional file 1:

597 SupplData1\_1.xlsx.

598

599

600

601 **Intra-country Analysis:** In Fig. 9, phylogenomic trees of male and female isolates from DRC are shown. For  
602 male patients (Fig. 9a), Haut-Katanga and Sud Kivu isolates cluster differently with isolates from other states.  
603 However, Kinshasha isolate clusters closely with Nord Kivu and Kongo Central isolates, while other Kinshasha,  
604 Sud Kivu and Haut-Katanga isolates cluster together showing less genome diversity. For female patients (Fig.  
605 9b), Kinshasha isolate clusters differently with isolates from Kongo Central, but clusters together with Sud Kivu  
606 isolate. However, the Sud Kivu and Kongo Central isolates (independently) cluster together indicating intra-  
607 specific genome similarity.

608

609

610 **Fig. 9. Phylogenomic trees for DRC's male and female isolates.** For full names of country codes, see Additional file 1:  
611 SupplData1\_1.xlsx.

612  
613

614 In Fig. 10, phylogenomic trees of male and female isolates from Nigeria are shown. For male patients (Fig. 10a),  
615 Osun isolate clusters differently with isolates from Lagos, Kwara, Oyo and other Osun isolates, including the  
616 unknow (infected) Nigerian isolate (NGA) who travelled to Greece. Oyo isolate clusters differently with isolates  
617 from Kwara, Osun and other Oyo isolate, indicating high genome diversity, even between same state. Also,  
618 Kwara isolate clusters differently with isolates from Osun and Oyo. However, Osun isolate clusters closely with  
619 isolates from Oyo and Kwara, indicating less genome diversity. For female patients (Fig. 10b), the  
620 phylogenomic tree present a near-flat structure with Ekiti isolate closely clustering with Osun, Ogun and the  
621 unknown Nigerian isolates (NGN)—infected through a local infection, and from Ogun. Furthermore, Ekiti and  
622 Ondo isolates cluster differently with isolates from other states.

623  
624

625 **Fig. 10. Phylogenomic trees for Nigerian male and female isolates.** For full names of country codes, see Additional file 1:  
626 SupplData1\_1.xlsx.

627  
628

629 In Fig. 11, phylogenomic trees of male and female isolates from Senegal are presented. For male patients (Fig.  
630 11a), Pikine isolate clusters differently with isolates from the rest of the states, while Diourbel and Dakar  
631 isolates cluster differently with other isolates, save Pikine. However, other Dakar isolates cluster closely with  
632 remaining isolates from Diourbel, St. Louis, Kolda and Thies, indicating less genome diversity between these  
633 isolates. For Female patients (Fig. 11b), Thies isolates cluster differently with isolates from Dakar and Diourbel  
634 as well as another isolate from Thies, indicating high genome diversity. However, Diourbel isolates cluster  
635 closely with remaining Diourbel isolate, Dakar and Thies isolates, indicating less genome diversity between  
636 these isolates.

637  
638

639 **Fig. 11. Phylogenomic trees for Senegalese male and female isolates.** For full names of country codes, see Additional file 1:  
640 SupplData1\_1.xlsx.

641

642  
643  
644  
645  
646  
647  
648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672

In Fig. 12, phylogenomic trees of male and female isolates from South Africa are presented. In Fig. 12b, female isolates/patients show high diversity as more isolates cluster differently, compared to male isolates (Fig. 12a), which maintain a near-flat tree structure, with the eThekweni isolate clustering differently from all other isolates.

**Fig. 12. Phylogenomic trees for South African male and female isolates.** For full names of country codes, see Additional file 1: SupplData1\_1.xlsx.

### **Nucleotide Similarity Analysis**

Several techniques for biological sequence alignment (multiple or pairwise) have flourished the literature<sup>[64]</sup>, but most of these techniques suffer from the lack of accuracy and partial interpretations. A direct pairwise genome sequence alignments (embedded in Algorithm 1) was carried out to match each nucleotide pair at the exact nucleotide positions of the SARS-CoV-2 genome, extending the alignments across other genome pairs. The output is a matrix of similarity scores. Fig. 13 shows inter- and intra-nucleotide similarities with strong evolutionary relationships highlighted in green color (for more clearer view of the nucleotide similarity matrices, see SupplData2\_1.xlsx). For male isolates (Fig. 13a), inter-nucleotide similarities cut across the following countries with strong evolutionary relationships observed for the countries listed beside them: Algeria (Senegal). Benin (DRC; Gambia; Greater Ghana; Kenya; Mali; Tunisia). Cameroon (DRC; Nigeria; Senegal; South Africa). Kenya (Mali; Nigeria; Tunisia). Morocco (South Africa). Nigeria (Senegal; South Africa; Tunisia). Senegal (South Africa). Intra-nucleotide similarities exist for the following isolates, with strong evolutionary relationships between the (states) listed beside them: Algeria (Bilda-Bilda). Benin (Cotonou-Cotonou, Cotonou-Oueme). DRC (Haut-Katanga-Haut-Katanga, Haut-Katanga-Kinshasha, Haut-Katanga-Kongo Central). Egypt (Cairo-Cairo). Gambia (Kombo-West Coast Region). Ghana (Greater Ghana-Greater Ghana). Kenya. Mali (Bamako-Bamako, Bamako-Mopti). Nigeria (Lagos-Osun; Lagos-Oyo, Osun-Oyo). Senegal (Dakar-Kolda; Dakar-St. Louis; Dakar-Thies; Kolda- St. Louis; Kolda-Thies; St. Louis-St. Louis; Thies-Thies). South Africa has the highest number of isolates with intra-nucleotide similarities spread across the various cities. The following states have strong evolutionary relationships with those listed beside them: Amajuba (eThekweni; Iembe; King Cetshwayo; Kwazulu Natal; Umgungundlovu; Umkhanyakude;

673 Umzinyathi; Uthukela; Uthungulu). Berea (Berea; Kwazulu Natal). Cape Town Metro (eThekwini; Harry  
674 Gwala; King Cetshwayo; Uthukela). eThekwini however, exhibits similar nucleotide relationship with Ilembe,  
675 King Cetshwayo, Kwazulu Natal, Umgungundlovu, Umkhanyakude, Umzinyathi; and share evolutionary  
676 relationships with some of these states (Ilembe; King Cetshwayo; Kwazulu Natal; Umgungundlovu;  
677 Umkhanyakude; Umzinyathi; Uthukela; Uthungulu).

678 For female isolates (Fig. 13b), inter-nucleotide similarities cut across the following countries with strong  
679 evolutionary relationships observed for the countries listed beside them: Algeria (DRC; Nigeria). Benin  
680 (Gambia; Ghana; Kenya; Mali; Morocco; Nigeria; Tunisia). DRC (Nigeria; Senegal; South Africa). Gambia,  
681 Ghana and Kenya exhibit similar nucleotide relationship with some of these countries (Kenya; Mali; Nigeria;  
682 South Africa and Tunisia). Morocco (Nigeria). Nigeria (South Africa; Tunisia). Senegal (South Africa). Intra-  
683 nucleotide similarities exist for the following isolates with strong evolutionary relationships in the states listed  
684 beside them: Benin (Cotonou-Cotonou). DRC (Haut-Katanga-Kongo Central, Kinshasha-Kinshasha, Haut-  
685 Katanga-Sud Kivu). Gambia (Kombo-West Coast Region). Ghana (Greater Ghana-Greater Ghana). Kenya.  
686 Madagascar (Fenoarivo-Antananarivo). Mali (Bamako-Bamako). Morocco (Rabat-Rabat). Nigeria (Ekiti-Ondo;  
687 Ekiti-Osun). Senegal (Dakar-Diourbel; Dakar-Thies; Diourbel-Diourbel). South Africa has the highest number  
688 of isolates with intra-nucleotide similarities spread across the various cities. The following states have strong  
689 evolutionary relationships with those listed beside them: Amajuba (Cape Town; Ilembe; North West; Overport;  
690 Umgungundlovu; Uthukela; Uthungulu; Zululand). Berea (North West; Umbilo; Umgungundlovu;  
691 Umkhanyakude; Uthungulu). Cape Town Metro (Cape Town; Ilembe; North West; Umgungundlovu; Uthukela;  
692 Zululand). EC (eThekwini). eThekwini (Harry Gwala; Ilembe; King Cetshwayo; Kwazulu Natal; Stanger;  
693 Uthukela; Umzinyathi). Free State (Free State; Harry Gwala; King Cetshwayo; Morningside; North West;  
694 Sisonke; Ugu; Umbilo). eThekwini however, exhibits similar nucleotide relationship with Harry Gwala, Ilembe,  
695 King Cetshwayo and Kwazulu Natal; and share evolutionary relationships with some of these states (Harry  
696 Gwala; Ilembe; King Cetshwayo; Kwazulu Natal; Stanger; Umkhanyakude; Umzinyathi).

697

698

699 **Fig. 13. Nucleotide similarity matrices.** Green colored cells are regions of high similarity that may indicate functional, structural and/or  
700 evolutionary relationships between nucleotide sequences.

701

702

703

704 **Sub-strain Pattern and Transmission Route Discoveries**

705 Comparing component planes of self-organizing maps (SOMs) can help detect genome expression patterns in  
706 identical positions (indicating correlation between the respective components), suitable for the discovery of sub-  
707 strains across the respective isolates. Component planes representation can enable the visualization of the  
708 relative component distributions of input data. Each component plane having the relative distribution of one data  
709 vector component. Local correlations can also occur if two parameter planes are similar in some regions.  
710 Furthermore, both linear and non-linear correlations including local or partial correlations between variables are  
711 possible<sup>[65]</sup>. We achieved the correlation hunting automatically, by decoupling the SOM correlations for  
712 correlation coefficients of at least 0.60, to explore patterns among pairwise genome samples for distinct  
713 identification of transmission pathways or routes. As can be seen in Fig 12, the SOM component planes reveal  
714 both inter- and intra-country sub-strains transmission for male and female patients. For male patients, the  
715 number of (sub-strains/isolate count) discovered by country include: Algeria (1/2), Benin (3/5), Cameroon (1/1),  
716 DRC (4/15), Egypt (3/5), Gambia (3/6), Ghana (2/4), Kenya (2/4), Mali (1/5), Morocco (3/5), Nigeria (5/13),  
717 Senegal (5/14), South Africa (7/65), and Tunisia (1). For female patients, the number of (sub-strains/isolate  
718 count) discovered by country include: Algeria (1/1), Benin (2/4), DRC (3/13), Egypt (2/4), Gambia (1/4), Ghana  
719 (2/4), Kenya (3/4), Madagascar (2/2), Mali (3/4), Morocco (2/2), Nigeria (5/11), Senegal (5/11), South Africa  
720 (9/80), and Tunisia (1).

721

722

723 **Fig. 14. SOM component planes visualization.** Maps are ordered by countries, with at least 1 isolate per country. The isolate numbers (1-  
724 145) represent the various states of the country excavated from GSAID (see SupplData1\_1.xlsx). The isolates are distributed by countries as  
725 follows. Male: Algeria (1-2), Benin (3-7), Cameroon (8), DRC (9-23), Egypt (24-28), Gambia (29-34), Ghana (35-38), Kenya (39-42), Mali  
726 (43-47), Morocco (48-52), Nigeria (53-65), Senegal (66-79), South Africa (80-144), Tunisia (145). Female: Algeria (1), Benin (2-5), DRC  
727 (6-18), Egypt (19-22), Gambia (23-26), Ghana (27-30), Kenya (31-34), Madagascar (35-36), Mali (37-40), Morocco (41-42), Nigeria (43-  
728 53), Senegal (54-64), South Africa (65-144), Tunisia (145).

729

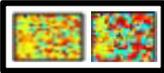
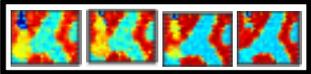
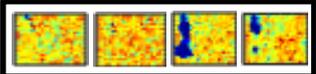
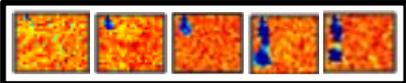
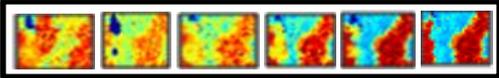
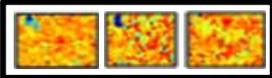
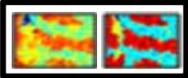
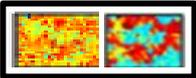
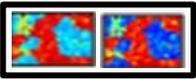
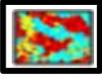
730

731 Discovering transmission routes of a pandemic can be very challenging but could assist both inter- and  
732 intra-country contact tracing. Using the Python programming language, cognitive knowledge was mined to  
733 localize the transmission routes of every isolate and provide appropriate link(s) to similar isolates with identical  
734 pattern(s). We observed multiple inter- and intra-country transmissions, with 10 and 12 sub-strains and their  
735 variants. This further knowledge was filtered out from the SOM component planes visualization of the male and

736 female isolates (Fig. 14), and presented in Table 4 and Table 5, respectively. Although there were noise infested  
 737 genomes (a product of genome sequencing and other unseen defects, which contributed to altering the SOM  
 738 image, causing semblance of dark blue like clots or stains (e.g., clusters 2-7, of Table 4, and clusters 2-8, 10-12,  
 739 of Table 5), they did not however, significantly alter the observed pattern(s).

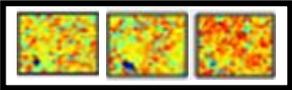
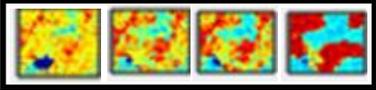
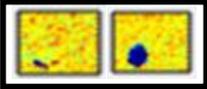
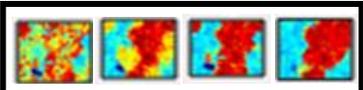
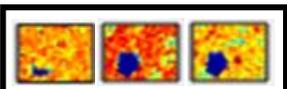
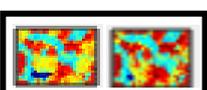
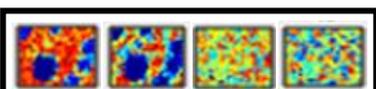
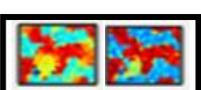
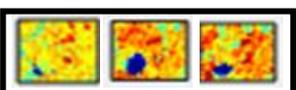
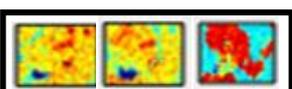
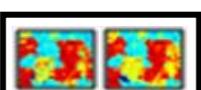
740

741 **Table 4. Discovered SARS-CoV-2 sub-strain clusters and cluster links of male isolates**

Cluster	Filtered pattern	Isolate link cluster
1.		1, 2,68,72
2.		3,29,32,33,34,42,43, 44,45,46,47,65
3.		4,5,6,8,13,14,15,18,19,20,23,24,30, 49,50,56,61,63,64, 69,70,79,82,83,99, 111,117,118,119, 132,135,139, 141, 143,144
4.		7,17,28,35,37,38,39,40,41,48,54,55,58,59, 60,62,66,67, 71,74,89,102,112, 113,115,116,122, 125,128,140,142,145
5.		9,10,11,12,73,76,77,78,81,87,88, 90,91,92,93,95,96, 101,104,106,109, 127,129,134,136,
6.		16,21,22,31,36,75,80,84,94,100,105,107,1 14,120,124,133
7.		25,26,27,51,52,103
8.		53,123,126,131
9.		57,85,86,108,110
10.		97,98,121,130

742

744 **Table 5. Discovered SARS-CoV-2 sub-strain clusters and cluster links of female isolates**

Cluster	Filtered pattern	Isolate link cluster
1.		1,8,22,48,50,51,55,61,81,86,121,124,128,131,137
2.		2,3,5,23,24,25,26,27,28,33,34,39,41,52,53,145
3.		4,10,15,18,19,20,21,31,40,46,57,58,63,66,67,68,70,71,72,78,79,82,9 3,100,103,104,112,127, 133,136,140,143
4.		6,9,11,12,13,14, 16,17,59,60,62,92, 94,96,97,98,130
5.		7,35,37,38,45, 54,56,64,73,77,101,102,105,107,116,117,129,139
6.		29,30,32,42,108,113
7.		36,69,75,76,122 135,138,142
8.		43,44,47,49,74, 80,99,114,126
9.		65,123,134,141,144
10.		83,84,89,90,91,106,109
11.		85,88,110,111,115,119, 125
12.		87,95,118,120,132

746  
747  
748  
749  
750  
751  
752  
753  
754  
755  
756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776

[Table 2, here]

[Table 3, here]

### **Test for Statistical Significance**

We conduct Friedman’s test<sup>[66]</sup>, a non-parametric statistical test similar to the parametric repeated measures ANOVA and used for detecting differences in treatments across multiple test attempts. The procedure ranks each group of isolates (or block) together, and then considers the values of ranks by country (for inter-country analysis) or by states (for intra-country analysis). The Nemenyi’s post hoc test for critical difference (CD)<sup>[67]</sup> was performed where the overall groups significantly differed from the observed characteristics of the isolates, as detected by the Friedman’s test. For the inter-country analysis, we selected countries with up to 13 states (selected randomly), to allow for a balanced block design, hence, resulting in 4 countries (DRC, Nigeria, Senegal and South Africa). For intra-country analysis, we selected the country with the highest number of states have up to 3 samples, to allow for a balanced block design, hence resulting in only 1 country, South Africa.

**Inter-country Analysis:** From the results of the Friedman test, there is a very highly significant difference in the male isolate groups from the selected countries ( $p < 0.01$ ), see Fig. 15a. Moreover, from the Nemenyi post hoc test, the CD plot reveals that isolates in any two countries, from among those listed on the left end of the CD plot (but with the exception of South Africa: Kwazulu Natal isolate) are not significantly different (as evident in the thick horizontal line connecting any pair of lines representing those countries). However, the isolate from South Africa (Kwazulu Natal isolate) is observed to be significantly different from those of Senegal (St. Louis and Diourbel isolates) and Nigeria (Kwara and Oyo isolates), only, among those on the left end of the plot. Similarly, any two countries from among those listed on the right end of the CD plot (with the exception of South Africa: MP isolate) have isolates to be not significantly different. However, Isolates from MP are

777 significantly different from those of Nigeria (unknown) and South Africa (eThekwini), only, among those  
778 countries that also appear on the right end of the plot. Moreover, isolates from South Africa (Kwazulu Natal,  
779 Harry Gwala and Stanger) are observed to be significantly different from those of each country on the right end  
780 of the CD plot, while the isolates from South Africa (Zululand) are significantly different from those of each  
781 country on the right end (except Nigeria: unknown isolate) of the plot. Similarly, isolates from South Africa  
782 (Ilembe, Cape Town and Amajuba) are significantly different from those of each country on the right end  
783 (except Nigeria: unknown isolate, and South Africa: eThekwini isolate) of the plot. Isolates from South Africa  
784 (Berea) are significantly different from those of each country on the right end (with the exception of Nigeria:  
785 unknown isolate, South Africa: eThekwini isolate, and DRC: Haut-Katanga isolate) of the plot, while isolates  
786 from South Africa (King Cetshwayo) are significantly different from those of each country on the right end  
787 (with the exception of Nigeria: unknown isolate, and South Africa: eThekwini isolate, DRC: Haut-Katanga  
788 isolate, and Nigeria: Osun isolate) of the plot. Isolates from Senegal (Kolda) are observed to be significantly  
789 different from each of Daka, South Africa (North West), DRC (Kongo Central) and South Africa (MP), among  
790 the countries on the right end of the CD plot; but isolates from each of Senegal (St. Louis and Diourbel), Nigeria  
791 (Kwara and Oyo) are significantly different from South Africa (MP), only, among the countries on the right end  
792 of the plot.

793 From the results of the Friedman test, there is a very highly significant difference in the female isolates  
794 from the selected countries ( $p < 0.01$ ), see Fig. 15b. Moreover, from the Nemenyi post hoc test, the CD plot  
795 reveals that isolates from any two countries from among those listed on the left end of the CD plot but excluding  
796 South Africa (Ugu, Harry Gwala, Amajuba), Nigeria (Osun and Ogun) are not significantly different (as evident  
797 in the thick horizontal line connecting any pair of lines representing those countries). However, isolates from  
798 any of South Africa (Harry Gwala, Amajuba), Nigeria (Osun and Ogun) are observed to be significantly  
799 different from those of each of Senegal (Thies isolate) and South Africa (Ilembe isolate), while the isolates from  
800 South Africa (Ugu) are significantly different from the isolates from each of Nigeria (unknown isolate), Senegal  
801 (Thies isolate) and South Africa (Ilembe isolate). The isolates from any two countries among those listed on  
802 the right end of the CD plot are not significantly different. Moreover, isolates from each of South Africa (Ugu,  
803 Harry Gwala, Amajuba), Nigeria (Osun and Ogun) are significantly different from those of any country  
804 appearing on the right end of the plot. Similarly, isolates from each of South Africa (Zululand and eThekwini)  
805 are observed to be significantly different from those of each country on the right end (except Senegal: Diourbel  
806 and Nigeria: Ekiti) of the CD plot, while the isolates from DRC (Haut-Katanga) are significantly different from

807 those of each country on the right end (except Senegal: Diourbel isolate, Nigeria: Ekiti isolate, DRC: Kinshasha  
808 isolate and Senegal: Dakar isolate) of the plot. Isolates from South Africa (King Cetshwayo ) are significantly  
809 different from those of each country on the right end (except, Senegal: Diourbel isolate, Nigeria: Ekiti isolate,  
810 DRC: Kinshasha isolate, Senegal: Dakar and Kongo Central isolates) of the CD plot, while the isolates from  
811 each of Nigeria (unknown), Senegal (Thies) and South Africa (Ilembe) differ significantly from those of South  
812 Africa (Free State, Cape Town and North West) and Nigeria (Ondo) only, among those on the right end of the  
813 plot.

814

815

816 **Fig. 15. Inter-country CD plots for male and female patients.** For a significance level  $\alpha$  the Nemenyi's test determines the critical  
817 difference (CD), if the difference between the average ranking of two treatments is greater than CD, then the null hypothesis that the  
818 treatments have the same performance is rejected.

819

820

821 **Intra-country Analysis:** From the results of the Friedman test there is a very highly significant difference in the  
822 male isolates from the various isolate groups from selected states ( $p < 0.01$ ), see Fig. 16a. Moreover, from the  
823 Nemenyi post hoc test, the CD plot reveals that isolates from any two states among Berea, Kwazulu Natal, ,  
824 Ilembe, Harry Gwala, Umgungundlovu, Zululand, Umkhanyakude and eThekweni (which are the states listed  
825 on the left end of the CD plot) are not significantly different (as evident in the thick horizontal line connecting  
826 any pair of lines representing those states). Similarly, any two states from among Amajuba, King Cetshwayo,  
827 Umzinyathi, Uthungulu, Cape Town and Ugu (the first 6 states on the right end of the CD plot) have their  
828 isolates to be not significantly different. However, Isolates from North West are significantly different from  
829 those of each state on the right end (except Uthungulu, Cape Town and Ugu) of the CD plot. Moreover, isolates  
830 from Berea are found to be significantly different from those of each state on the right end (except Amajuba) of  
831 the CD plot, while the isolates from Kwazulu Natal, Ilembe, Harry Gwala and Umgungundlovu are each  
832 significantly different from those of each state on the right end (except Amajuba, King Cetshwayo and  
833 Umzinyathi) of the CD plot. Similarly, isolates from Zululand are significantly different from those of each state  
834 on the right end (except Amajuba, King Cetshwayo, Umzinyathi and Uthungulu) of the plot, while the isolates  
835 from Umkhanyakude and eThekweni are each significantly different from those of North West, only.

836 From the results of the Friedman test there is a very highly significant difference in the female isolate  
837 groups from selected states ( $p < 0.01$ ), see Fig. 16b. Moreover, from the Nemenyi post hoc test, the CD plot

838 reveals that isolates from any two states among Umgungundlovu, Harry Gwala , . . . , Umkhanyakude (which  
839 are the states listed on the left end of the CD plot) are not significantly different (as evident in the thick  
840 horizontal line connecting any pair of lines representing those states). Similarly, any two states from among  
841 Uthungulu, eThekwini, . . . , Ilembe (the first 7 states on the right end of the CD plot) have their isolates to be not  
842 significantly different. However, Isolates from North West are significantly different from those of each state on  
843 the right end (except Capetown, Ilembe and Free State) of the CD plot. Moreover, isolates from  
844 Umgungundlovu and Harry Gwala are each found to be significantly different from those of each state on the  
845 right end (except Uthungulu, eThekwini, Zululand and Kwazulu Natal) of the CD plot, while the isolates from  
846 Berea and Uthukela are each significantly different from those of each of Ilembe, Free State and North West  
847 only (which appear on the right end of the plot). Similarly, isolates from Amajuba, King Cetshawayo,  
848 Umzinyathi, Ugu and Umkhanyakude are each significantly different from those of each of FS and North West  
849 only (both appearing on the right end of the plot).

850

851

852 **Fig. 16. Intra-country CD plots for male and female patients.** For a significance level  $\alpha$  the Nemenyi's test determines the critical  
853 difference (CD), if the difference between the average ranking of two treatments is greater than CD, then the null hypothesis that the  
854 treatments have the same performance is rejected.

855

856

857

## 858 **Discussion**

859 Issues of gender and the human genome present several levels of implications at basic scientific research,  
860 clinical applications and wider societal investigations<sup>[68]</sup>. Excluding those with co-morbidities and the aged,  
861 males have been found to have worse prognosis during COVID-19 infections<sup>[69]</sup> with delayed viral clearance  
862 compared to females, which show evidence of immune tolerance and slow prognosis<sup>[70]</sup>. Hence, characterization  
863 of the sub-strain clusters by gender clearly explains the diversity of SARS-CoV-2. To the best of our  
864 knowledge, research works conducted so far on SARS-CoV-2 genome analysis, aside demographic  
865 classification, have ignored the gender dimension. This paper is therefore the first to consider the extent to  
866 which SARS-CoV-2 sub-strain transmissions impact on gender. The implication of our findings is most likely to  
867 introduce additional information to existing body of knowledge on COVID-19 and aid further research works in  
868 this area that would balance the gender dimension.

869 Understanding the pattern of spontaneous mutation is fundamental in studies of human genome evolution  
870 and genetic disease<sup>[71]</sup>. Mutation diversity therefore appears to be a direct consequence of changing/evolving  
871 sub-strains, as it represents the ultimate source of genetic variation and explains the story behind their evolution.  
872 However, extensive variability exists among different genes or genome regions in between- and within-  
873 species<sup>[72,73]</sup>, and suggests that spontaneous mutation rates are not always constant across the genome—as only a  
874 small subset of the new mutation manifests in disease variants<sup>[74]</sup>. But despite mutation diversity and localized  
875 variation in DNA polymorphism diversity and recombination, the pattern of sub-strains is not affected to cause  
876 confusion in sub-strain(s) identification. Whole-genome alignment predicts evolutionary relationships at the  
877 nucleotide level between two or more genomes<sup>[75]</sup>. In this paper, homologous pairs of (nucleotide) positions  
878 between genome sequences were compared. Aside noise, the genomes are also colinear, as they have not been  
879 broken by rearrangement event. Although local alignments between genomes were realized, our alignment  
880 algorithm is costly (ran in quadratic time) and requires further improvements, but the similarity scores produced  
881 (see Fig. 13) are very useful Knowledge Base component for building intelligent diagnostic systems. Most  
882 whole genome algorithms are restricted in the evolutionary relationships captured, as only a subset of  
883 homologous relationships may be of interest. Large-scale phylogenomic methodology shows high potential, as  
884 application of diverse datasets confirm robustness of the approach.

885 Wittler<sup>[76]</sup> proposed a new whole genome-based approach for inferring aligned- and reference-free  
886 phylogenies. Their method adopted a colored de Bruijn graph to extract common subsequences for deducing  
887 phylogenetic splits, instead of relying on pairwise comparisons to determine distances and tree edges deduction.  
888 Similarity in nucleotide sequence is a diversity indicator that measures the relative closeness of gene or genome  
889 isolates. Identifying similarities between sequences of special interest is one of the most important goals when  
890 working with nucleotide sequences. A distance measure associates the numeric value with a pair of sequences.  
891 Direct nucleotide protein sequencing technology<sup>[77]</sup> have resulted in the explosive growth in the number of  
892 known sequences. The results of the nucleotide similarity analysis revealed isolates with strong evolutionary  
893 relationships between and within countries. Hierarchical clustering that implements agglomerative nesting was  
894 adopted in this research for genome-wide ranking instead of focusing on specific subsets. Sub-strains and  
895 transmission patterns discoveries confirmed multiple strains with some isolates showing identical sub-strains  
896 patterns (with less-diversity), while others showed distinct terminal patterns (without further changes).

897 SOM results showed reduced sub-strain variants—for increased isolates/states, with disproportionate sub-  
898 strains increase or decrease in some states, for male patients (e.g., Gambia=50%, DRC=26.67%,

899 Nigeria=38.46%, Senegal=35.71% and South Africa=10.77%) and female patients (e.g., Gambia=25%,  
900 DRC=23.07%, Nigeria=45.45%, Senegal=45.45% and South Africa=11.25%), hence, establishing a non-linear  
901 relationship between mutation and transmission patterns by gender. The generated cognitive maps (Table 2 and  
902 Table 3) efficiently associates similar isolate clusters for transmission pathway analysis. The practical  
903 implication of this is that early inter- and intra-transmission routes could easily be traced, and immediate contact  
904 tracing commenced. Further, countries/states with high prevalence rate could temporarily be locked until  
905 satisfactory contact tracing is achieved. The discovered sub-strain variants (10 in male patients and 12 in female  
906 patients) reveals the uncertain nature of the SARS-CoV-2 and may be a pointer to a second wave of the virus in  
907 Africa.

908 Finally, the statistical test for significance (Friedman's test) showed highly significant difference for inter-  
909 and intra-country analysis by gender, with the Nemenyi's post hoc test revealing significance difference in the  
910 countries/states selected.

911

912

## 913 **Conclusions**

914 Infectious disease prediction has significantly benefited from the use of genome mining, which is entirely  
915 dependent of computing technology and bioinformatics tools used<sup>[78]</sup>. The World Health Organization  
916 (<https://www.afro.who.int/news/covid-19-genome-sequencing-laboratory-network-launches-africa>) has  
917 underscored the need for application of genome mining in the management of COVID-19 in Africa by  
918 collaborating with the Africa Centers for Disease Control and Prevention (Africa CDC) to launch a network of  
919 twelve specialized laboratories to facilitate genome sequencing of SARS-CoV-2 virus to track the evolution and  
920 mutation of the virus and create an effective mechanism for response to the virus. The grouping of viruses from  
921 different countries into lineages has proved useful in establishing the route of virus importation across countries.

922 In this study, we have introduced a cooperatively inspired open source framework for intelligent mining of  
923 SARS-CoV-2 genomes using the unsupervised self-organizing map (SOM), which takes advantage of similarity  
924 in genetic behavior of the strains of the SARS-CoV-2 virus. The SOM is among the family of machine learning  
925 techniques that facilitate engines that further features probing of genes for precise classification and prediction,  
926 which could be useful for screening and treatment, contact tracing, prediction and forecasting, and drugs/vaccine  
927 discovery. Our open source framework is a hybridized system that helps an in-depth understanding infectious  
928 disease prevalence. Our framework generates phylogenomic trees, pairwise nucleotide similarity

929 matrices/scores, gene diversity plots, genome expression patterns analysis, essential for enriching the genome  
930 datasets, towards intelligent genome characterization and prediction; thus, facilitating community contribution  
931 and replicability.

932 The results of our study show the following: i) Africa countries exhibit varying levels of nucleotide mutation;  
933 for example Congo had a high mutation burden (8.30), while Kenya had the least (5.38); ii) The transition from  
934 the cytosine to the thymine nucleotide (C>T) accounted for the highest level of mutation, followed by the  
935 adenine to guanine (A>G) transition; iii) the average nucleotide count in male and female isolates show  
936 approximately similar ratios ( $A \approx 32$ ,  $C \approx 18$ ,  $G \approx 20$ ,  $T \approx 32$ ); iv) the genome diversity analysis show a  
937 smoother distribution curve with male isolates when compared with female isolates; v) the phylogenomic  
938 analysis suggests independent sub-strain mutant accumulation in various countries; vi) from the excavated data,  
939 various African countries exhibit varying numbers of sub-strain transmissions; South Africa had the highest  
940 (male-7, female-9), followed by Nigeria (male-5, female-5), while Tunisia, Algeria had only one sub-strain for  
941 both male and female isolates; and vii) multiple inter-country transmissions were observed, with 10 and 12 sub-  
942 strains and their variants.

943 The academic and policy implications of this study are as follows: i) it contributes to a better understanding  
944 of the prevalence and transmission of the SARS-CoV-2 virus in Africa; ii) it provides a framework for inter- and  
945 intra-country contact tracing, especially in undocumented infection sources; iii) it provides a basis of revealing  
946 hidden sub-strains, which could be useful in time varying prediction of infection patterns; iv) the computation of  
947 variability in emerging sub-strains by gender isolates could be very useful in the development of appropriate  
948 SARS-CoV-2 vaccines.

949

950

951

952

953

## 954 **List of Abbreviations**

955	AGNES	Agglomerative Nesting
956	AI	Artificial Intelligence
957	COVID-19	Coronavirus Disease 2019
958	DNA	Di-Nucleic Acid

959	GISAID	Global Initiative on Sharing All Influenza Data
960	HCA	Hierarchical Clustering Analysis
961	MATLAB	MATrix LABoratory
962	MERS	Middle East Respiratory Syndrome
963	MERS-CoV	Middle East Respiratory Syndrome Corona Virus
964	ML	Machine Learning
965	NCBI	National Center for Biotechnology Information
966	ORF	Open Reading Frame
967	PANGOLIN	Phylogenetic Assignment of Named Global Outbreak LINEages
968	RDT	Rapid Diagnostic Test
969	RNA	Ribonucleic Acid
970	RT-PCR	Reverse Transcriptase Polymerase Chain Reaction
971	SARS	Severe Acute Respiratory Syndrome
972	SARS-CoV-2	Severe Acute Respiratory Syndrome Corona Virus 2
973	SOM	Self-Organizing Map

974

975

976

977 **Declarations**

978 None declared

979

980

981

982 **Availability of data and materials**

- 983
- The datasets used and/or analyzed during the current study are available at: <https://gisaid.org>
- 984
- All data generated or analyzed during this study are included in this published article [and its
- 985 supplementary information files].

986

987

988

989 **Competing interests**

990 There are no competing interests

991

992

993

994 **Authors' contributions**

995 All authors contributed equally to the final manuscript. Specifically,

996 M.E.\* conceptualized the research idea, contributed to the research methods, preparation of figures,

997 framework/tools design, implementation and interpretation of the results.

998 M.E.1 provided literature materials, performed critical review as well as data validation.

999 U.I. contributed to the research methodology, framework/tools design, preparation of figures, implementation

1000 and interpretation of results.

1001 F-M.U. structurally edited the original draft and contributed to the tools design component and implementation.

1002 I.U. was involved in critical review of literature and research data validation.

1003 N.U. was involved in the analysis and interpretation of statistical related components.

1004 I.E. contributed to the biotechnology and bioinformatics components of the paper.

1005 A.M. was involved in formal analysis of the study data, research methods and implementation.

1006 E.A. was involved in the critical review as well as a formal analysis of the study data.

1007 Y.T. contributed to the biotechnology components of the paper, including editing of the original draft.

1008 G.J. was involved in data curation, collection/excavation and processing of the human SARS-CoV-2 genomes.

1009 E.D. was involved in data curation, collection/excavation and processing of the human SARS-CoV-2 genomes.

1010 J.N. was involved in critical review of literature and research data validation.

1011

1012

1013

1014 **Acknowledgements**

1015 We are grateful to the authors, originating and submitting laboratories of the collected sequences from

1016 GISAID's EpiFlu database on which this research rests.

1017

1018

1019

1020 **References**

- 1021 [1] Fehr AR, Perlman S. Coronaviruses: an overview of their replication and pathogenesis. In Coronaviruses.  
1022 2015; (pp. 1-23). Humana Press, New York, NY. [https://doi.org/10.1007/978-1-4939-2438-7\\_1](https://doi.org/10.1007/978-1-4939-2438-7_1).
- 1023 [2] Tang X, Wu C, Li X, Song Y, Yao X, Wu X, Duan Y, Zhang H, Wang Y, Qian Z, Cui J. On the origin and  
1024 continuing evolution of SARS-CoV-2. *Natl Sci Rev.* 2020; 7: 1012–1023.  
1025 <https://doi.org/10.1093/nsr/nwaa036>.
- 1026 [3] Università di Bologna. The six strains of SARS-CoV-2. ScienceDaily. ScienceDaily, 2020 August 3.  
1027 <https://www.sciencedaily.com/releases/2020/08/200803105246.htm>.
- 1028 [4] Islam H, Rahman A, Masud J, Shweta DS, Araf Y, Ullah MA, Sium SMA, Sarkar, B. A Generalized  
1029 Overview of SARS-CoV-2: Where Does the Current Knowledge Stand? *Electron J Gen Med.* 2020; 17 (6):  
1030 em251. <https://doi.org/10.29333/ejgm/8258>
- 1031 [5] Wiechers IR, Perin NC, Cook-Deegan R. The emergence of commercial genomics: analysis of the rise of a  
1032 biotechnology subsector during the Human Genome Project, 1990 to 2004. *Genome Med.* 2013; 5(83): 1-9.  
1033 <https://doi.org/10.1186/gm487>
- 1034 [6] Giani AM, Gallo, GR, Gianfranceschi L, Formenti G. Long walk to genomics: History and current  
1035 approaches to genome sequencing and assembly. *Comput. Struct. Biotechnol.* 2020; 18. p. 9-19.  
1036 <https://doi.org/10.1016/j.csbj.2019.11.002>.
- 1037 [7] Shey M, Okeibunor JC, Yahaya AA, Herring BL, Tomori O, Coulibaly SO, Gumede-Moeletsi, HN,  
1038 Mwenda JM, Yoti Z, Wiysonge CS, Talisuna AO. Genome sequencing and the diagnosis of novel  
1039 coronavirus (SARS-COV-2) in Africa: how far are we?. *Pan Afr Med J.*2020; 36: 80.  
1040 <https://doi.org/10.11604/pamj.2020.36.80.23723>.
- 1041 [8] Adebisi YA, Oke GI, Ademola PS, Chinemelum, IG, Ogunkola IO, Lucero-Prisno III DE. SARS-CoV-2  
1042 diagnostic testing in Africa: needs and challenges. *Pan Afr Med J.* 2020. 35(2): 4. [https://doi.org/10.11604/](https://doi.org/10.11604/pamj.2020.35.4.22703)  
1043 [pamj.2020.35.4.22703](https://doi.org/10.11604/pamj.2020.35.4.22703).
- 1044 [9] Ekpenyong M, Udo I, Uzoka FM, Attai K. A Spatio-GraphNet Model for Real-time Contact Tracing of  
1045 CoVID-19 Infection in Resource Limited Settings. In Proceedings of the 4th International Conference on  
1046 Medical and Health Informatics 2020 pp. 208-217. <https://doi.org/10.1145/3418094.3418141>.

- 1047 [10] Das S, Ghosh P, Sen B, Mukhopadhyay, I. Critical community size for CoVID-19--a model-based approach  
1048 to provide a rationale behind the lockdown. arXiv preprint arXiv. 2020; 2004.03126.  
1049 <https://arxiv.org/pdf/2004.03126.pdf>
- 1050 [11] Maghdid HS, Ghafoor KZ. A Smartphone enabled Approach to Manage CoVID-19 Lockdown and  
1051 Economic Crisis. arXiv preprint arXiv. 2020; 2004.12240.
- 1052 [12] Sun L, Liu G, Song F, Shi N, Liu F, Li S, Li P, Zhang W, Jiang X, Zhang Y, Sun L, Chen X, Shi Y.  
1053 Combination of four clinical indicators predicts the severe/critical symptom of patients infected COVID-19.  
1054 J Clin Virol. 2020; <https://doi.org/10.1016/j.jcv.2020.104431>.
- 1055 [13] Yan L, Zhang HT, Goncalves J, Xiao Y, Wang M, Guo Y, Sun C, Tang X, Jing L, Zhang M, Huang X,  
1056 Xiao Y, Cao H, Chen Y, Ren T, Wang F, Xiao Y, Huang S, Tan X, Huang N, Jiao B, Cheng C, Zhang Y,  
1057 Luo A, Mombaerts L, Jin J, Cao Z, Li S., Xu H, Yuan Y. An interpretable mortality prediction model for  
1058 COVID-19 patients. Nat Mach Intell. 2020; <https://doi.org/10.1038/s42256-020-0180-7>.
- 1059 [14] Ren, L.L., Wang, Y.M., Wu, Z.Q., Xiang, Z.C., Guo, L., Xu, T., Jiang, Y.Z., Xiong, Y., Li, Ribeiro M. H.  
1060 D. M., da Silva R. G., Mariani V. C., Coelho L. D. S. (2020). Short-term forecasting COVID-19 cumulative  
1061 confirmed cases: perspectives for Brazil. Chaos, Solitons Fractals.  
1062 <https://doi.org/10.1016/j.chaos.2020.109853>.
- 1063 [15] Cabore JW, Karamagi HC, Kipruto H, Asamani JA, Droti B, Seydi ABW, Titi-Ofei R, Impouma B, Yao M,  
1064 Yoti Z, Zawaira F. The potential effects of widespread community transmission of SARS-CoV-2 infection  
1065 in the World Health Organization African Region: a predictive model. BMJ Global Health. 2020;5(5),  
1066 e002647. doi:10.1136/bmjgh-2020-002647.
- 1067 [16] Sun H, Dickens BL, Cook AR, Clapham HE. Importations of COVID-19 into African countries and risk of  
1068 onward spread. BMC Infect Dis. 2020; 20, 598. <https://doi.org/10.1186/s12879-020-05323-w>.
- 1069 [17] Ke YY, Peng TT, Yeh TK, Huang WZ, Chang SE, Wu SH, Hung HC, Hsu TA, Lee SJ, Song JS, Lin WH,  
1070 Chiang TJ, Lin JH, Sytwu HK, Chen CT. Artificial intelligence approach fighting COVID-19 with  
1071 repurposing drugs. Biomed J. 2020; <https://doi.org/10.1016/j.bj.2020.05.001>.
- 1072 [18] Beck BR, Shin B, Choi Y, Park S, Kang K. Predicting commercially available antiviral drugs that may act  
1073 on the novel coronavirus (SARS-CoV-2) through a drug-target interaction deep learning model. Comput  
1074 Struct Biotechnol J. 2020;18: 784–790. <https://doi.org/10.1016/j.csbj.2020.03.025>.

- 1075 [19] Ekins S, Mottin M, Ramos PRP, Sousa BKP, Neves BJ, Foil DH, Zorn KM, Braga RC, Coffee M, Southan  
1076 C, Puhl CA, Andrade CH. Déjà vu: stimulating open drug discovery for SARS-CoV-2. *Drug Discov Today*.  
1077 2020; <https://doi.org/10.1016/j.drudis.2020.03.019>.
- 1078 [20] Mainardes RM, Diedrich C. The potential role of nanomedicine on COVID-19 therapeutics. *Therapeutic*  
1079 *Delivery*. 2020; 11(7): 411-414. <https://doi.org/10.4155/tde-2020-0069>.
- 1080 [21] Chauhan G, Madou MJ, Kalra S, Chopra V, Ghosh D, Martinez-Chapa, SO. Nanotechnology for COVID-  
1081 19: therapeutics and vaccine research. *ACS nano*. 2020; 14(7), 7760-7782.  
1082 <https://doi.org/10.1021/acsnano.0c04006>.
- 1083 [22] Caccuri F, Zani, A, Messali S, Giovanetti M, Bugatti A, Campisi G, Filippini F, Scaltriti E, Ciccozzi M,  
1084 Fiorentini S, Caruso A. A persistently replicating SARS-CoV-2 variant derived from an asymptomatic  
1085 individual. *J Transl Med*. 2020. 18(1), 1-12. <https://doi.org/10.1186/s12967-020-02535-1>.
- 1086 [22] Wu D, Wu T, Liu Q, Yang Z. The SARS-CoV-2 outbreak: what we know. *International Journal of*  
1087 *Infectious Diseases*. 2020; 94 (2020):44-48. <https://doi.org/10.1016/j.ijid.2020.03.004>.
- 1088 [23] Wilder-Smith A, Chiew, CJ, Lee VJ. Can we contain the COVID-19 outbreak with the same measures as  
1089 for SARS?. *The Lancet Infectious Diseases*. 2020; 20: e102–07. [https://doi.org/10.1016/S1473-](https://doi.org/10.1016/S1473-3099(20)30129-8)  
1090 [3099\(20\)30129-8](https://doi.org/10.1016/S1473-3099(20)30129-8).
- 1091 [24] Rabi FA, Al-Zoubi MS, Kasasbeh, GA, Salameh DM, Al-Nasser, AD. SARS-CoV-2 and coronavirus  
1092 disease 2019: what we know so far *Pathogens*. 2020; 9(3), 231. <https://doi.org/10.3390/pathogens9030231>
- 1093 [25] Kamatenesi-Mugisha M, Oryem-Origa, H. Traditional herbal remedies used in the management of sexual  
1094 impotence and erectile dysfunction in western Uganda. *Afr. Health Sci*. 2005;5(1): 40–49.  
1095 <https://www.ajol.info/index.php/ahs/article/view/6896>.
- 1096 [26] Nordling L. Unproven herbal remedy against COVID-19 could fuel drug-resistant malaria, scientist warn.  
1097 *Science*. 2020; <https://doi.org/10.1126/science.abc6665>.
- 1098 [27] Li X, Giorgi EE, Marichannegowda, MH, Foley B, Xiao C, Kong XP, Chen Y, Gnanakaran S, Korber B,  
1099 Gao F. Emergence of SARS-CoV-2 through recombination and strong purifying selection. *Sci. Adv*. 2020;  
1100 6(27). <https://doi.org/10.1126/sciadv.abb9153>.
- 1101 [28] Lu R, Zhao X, Li J, Niu P, Yang B, Wu H, Wang W, Song, H, Huang B, Zhu, N, Bi, Y. Genomic  
1102 characterisation and epidemiology of 2019 novel coronavirus: implications for virus origins and receptor  
1103 binding. *Lancet Glob Health*. 2020; 395(10224): 565-574.

- 1104 [29] Zhang L, Yang, JR, Zhang Z, Lin Z. Genomic variations of SARS-CoV-2 suggest multiple outbreak  
1105 sources of transmission. medRxiv. 2020; <https://doi.org/10.1101/2020.02.25.20027953>.
- 1106 [30] Korber B, Fischer W, Gnanakaran SG, Yoon H, Theiler J, Abfalterer W, Foley B, Giorgi EE, Bhattacharya  
1107 T, Parker MD, Partridge DG. Spike mutation pipeline reveals the emergence of a more transmissible form  
1108 of SARS-CoV-2. bioRxiv. 2020; <https://doi.org/10.1101/2020.04.29.069054>
- 1109 [31] Wyllie AL, Fournier J, Casanovas-Massana A, et al. Saliva is more sensitive for SARS-CoV-2 detection in  
1110 CoVID-19 patients than nasopharyngeal swabs. medRxiv.2020;  
1111 <https://doi.org/10.1101/2020.04.16.20067835>
- 1112 [32] Chan KH, Poon LL, Cheng VC, Guan Y, Hung IF, Kong J, Yam LY, Seto WH, Yuen KY, Peiris JS.  
1113 Detection of SARS coronavirus in patients with suspected SARS. *Emerg. Infect. Dis.* 2004; 10(2), 294-299.  
1114 <https://doi.org/10.3201/eid1002.030610>
- 1115 [33] Tabibzadeh A, Zamani F, Laali A, Esghaei M, Tameshkel FS, Keyvani H, Makiani MJ, Panahi M,  
1116 Motamed N, Perumal D, Khoonsari M. SARS-CoV-2 molecular and phylogenetic analysis in COVID-19  
1117 patients: a preliminary report from Iran. *Infection, Genetics and Evolution.* 2020; p.104387.  
1118 <https://doi.org/10.1016/j.meegid.2020.104387>.
- 1119 [34] van der Made CI, Simons A, Schuurs-Hoeijmakers J, van den Heuvel G, Mantere T, Kersten S, van Deuren,  
1120 RC, Steehouwer M, van Reijmersdal, SV, Jaeger M, Hofste T. Presence of Genetic Variants Among Young  
1121 Men With Severe COVID-19. *JAMA.* 2020; 324(7):663-673. <https://doi:10.1001/jama.2020.13719>.
- 1122 [35] Torti C, Mazzitelli M, Treçarichi EM, Darius O. Potential implications of SARS-CoV-2 epidemic in  
1123 Africa: where are we going from now?. *BMC Infect Dis.* 2020; 20 (412). [https://doi.org/10.1186/s12879-](https://doi.org/10.1186/s12879-020-05147-8)  
1124 [020-05147-8](https://doi.org/10.1186/s12879-020-05147-8).
- 1125 [36] Sekizuka T, Kuramoto S, Nariai E, Taira M, Hachisu Y, Tokaji A, Shinohara M, Kishimoto T, Itokawa K,  
1126 Kobayashi Y, Kadokura K. SARS-CoV-2 Genome Analysis of Japanese Travelers in Nile River Cruise.  
1127 *Frontiers in Microbiology.* 2020; 11(1316). <https://doi.org/10.3389/fmicb.2020.01316>.
- 1128 [37] Kandeil A, Mostafa A, El-Shesheny R, Shehata M, Roshdy WH, Ahmed SS, Gomaa M, El Taweel A,  
1129 Kayed A E, Mahmoud SH, Moatasim Y, Kutkat O, Kamel MN, Mahrous N, El Sayes M, El Guindy NM,  
1130 Naguib A, Ali MA. Coding complete genome sequences of two SARSCoV-2 isolates from Egypt.  
1131 *Microbiol. Resour. Announc.* 2020; 9: e00489-20. <https://doi.org/10.1128/MRA.00489-20>.

1132 [38] Laamarti M, Kartti S, Alouane T, Laamarti R, Allam L, Ouadghiri M, Chemaou-Elfihri MW, Smyej I,  
1133 Rahoui J, Benrahma H, Diawara, I. Genetic analysis of SARS-CoV-2 strains collected from North Africa:  
1134 viral origins and mutational spectrum. *bioRxiv*, 2020; <https://doi.org/10.1101/2020.06.30.181123>.

1135 [39] Badaoui B, Sadki K, Talbi C, Tazi L, Salah D. Genetic diversity and genomic epidemiology of sars-cov-2  
1136 in Morocco. *bioRxiv*. 2020; <https://doi.org/10.1101/2020.06.23.165902>.

1137 [40] Laamarti M, Chemaou-Elfihri MW, Kartti S, Laamarti R, Allam L, Ouadghiri M, Smyej I, Rahoui J,  
1138 Benrahma H, Diawara I, Alouane T. Genome Sequences of Six SARS-CoV-2 Strains Isolated in Morocco,  
1139 Obtained Using Oxford Nanopore MinION Technology. *Microbiol. Resour. Announc.* 2020; 9(32): 1-4.  
1140 <https://doi.org/10.1128/MRA.00767-20>.

1141 [41] Jouali F, Marchoudi N, El Ansari FZ, Kasmi Y, Chenaoui M, El Aliani A, Azami N, Loukman S, Ennaji  
1142 MM, Benhida R, Fekkak J. SARS-CoV-2 Genome Sequence from Morocco, Obtained Using Ion AmpliSeq  
1143 Technology. *Microbiol. Resour. Announc.* 2020; 9(31): 1-3. <https://doi.org/10.1128/MRA.00690-20>.

1144 [42] Awoyelu EH, Oladipo EK, Adetuyi BO, Senbadejo TY, Oyawoye OM, Oloke J K. Phyloevolutionary  
1145 analysis of SARS-CoV-2 in Nigeria. *New microbes and new infections.* 2020; 36(100717).  
1146 <https://doi.org/10.1016/j.nmni.2020.100717>.

1147 [43] Happi C, Ihekweazu C, Oluniyi PE, Olawoye I. New SARS-CoV-2 Genomes from Nigeria Reveals  
1148 Dominance of Viruses with Spike Protein Mutation (D614G), and Additional Virus Lineages in Circulation.  
1149 *Genome Reports.* 2020; [https://virological.org/t/new-sars-cov-2-genomes-from-nigeria-reveals-dominance-  
1150 of-viruses-with-spike-protein-mutation-d614g-and-additional-virus-lineages-in-circulation/527](https://virological.org/t/new-sars-cov-2-genomes-from-nigeria-reveals-dominance-of-viruses-with-spike-protein-mutation-d614g-and-additional-virus-lineages-in-circulation/527).

1151 [44] Allam M, Ismail A, Khumalo ZTH, Kwenda S, van Heusden P, Cloete R, Wibmer CK, Mtshali PS,  
1152 Mnyameni F, Mohale T, Subramoney K, Walaza S, Ngubane W, Govender N, Motaze NV, Bhiman J.N.  
1153 Genome sequencing of a severe acute respiratory syndrome coronavirus 2 isolate obtained from a South  
1154 African patient with coronavirus disease 2019. *Microbiol. Resour. Announc.* 2020; 9:e00572-20.  
1155 <https://doi.org/10.1128/MRA.00572-20>.

1156 [45] Pillay S, Giandhari J, Tegally H, Wilkinson E, Chimukangara B, Lessells R, Mattison S, Moosa Y, Gazy I,  
1157 Fish, M, Singh L. Whole Genome Sequencing of SARS-CoV-2: Adapting Illumina Protocols for Quick and  
1158 Accurate Outbreak Investigation During a Pandemic. *bioRxiv*.2020;  
1159 <https://doi.org/10.1101/2020.06.10.144212>.

- 1160 [46] Bugembe, DL, Kayiwa J, Phan MV, Tushabe P, Balinandi, S, Dhaala, B, Lexow, J, Mwebesa, H, Aceng, J,  
1161 Kyobe, H, Ssemwanga, D. Main Routes of Entry and Genomic Diversity of SARS-CoV-2, Uganda.  
1162 *Emerging Infectious Diseases*. 2020; 26(10): 2411–2415. <https://doi.org/10.3201/eid2610.202575>.
- 1163 [47] Hamidouche M. COVID-19 outbreak in Algeria: A mathematical model to predict the incidence. medRxiv.  
1164 2020; <https://doi.org/10.1101/2020.03.20.20039891>.
- 1165 [48] Kouriba B, Duerr A, Rehn A, Sangare AK, Traoure BY, Bestehorn-Willmann MS, Ouedraogo JL, Heitzer  
1166 A, Sogodogo E, Maiga A, Walter MC. First phylogenetic analysis of Malian SARS-CoV-2 sequences  
1167 provide molecular insights into the genomic diversity of the Sahel region. medRxiv. 2020;  
1168 <https://doi.org/10.1101/2020.09.23.20165639>
- 1169 [49] TIBA COVID-19 Pandemic Response Unit (2020). SARS-CoV-2 genomes report for WHO Africa Region  
1170 [21/08/2020]. [http://tiba-partnership.org/tiba/sites/sbsweb2.bio.ed.ac.uk.tiba/files/pdf/SARS-CoV-](http://tiba-partnership.org/tiba/sites/sbsweb2.bio.ed.ac.uk.tiba/files/pdf/SARS-CoV-2%20Genome%20Report%2024.08.2020.pdf)  
1171 [2%20Genome%20Report%2024.08.2020.pdf](http://tiba-partnership.org/tiba/sites/sbsweb2.bio.ed.ac.uk.tiba/files/pdf/SARS-CoV-2%20Genome%20Report%2024.08.2020.pdf)
- 1172 [50] Muyembe-Tamfum, JJ, Ahuka-Mundeke S, Mbala-Kingebeni P, Nkwembe-Mgabana E, Kinganda-  
1173 Lusamaki E, Amuri-Aziza A, Muyembe-Mawete F, Lokilo-Lofiko E, Claude CCJ, Marceline AO, Bibiche  
1174 NM. Phylogenetic analysis of SARS-CoV-2 in DRC. ARTIC Network. 2020;  
1175 <https://virological.org/t/phylogenetic-analysis-of-sars-cov-2-in-drc/528>
- 1176 [51] KEMRI-CGMRC Introduction and local transmission of SARS-CoV-2 cases in Kenya. *Genome Reports*.  
1177 2020; <https://virological.org/t/introduction-and-local-transmission-of-sars-cov-2-cases-in-kenya/497>.
- 1178 [52] Giandhari J, Pillay S, Wilkinson E, Tegally H, Sinayskiy I, Schuld M, Lourenço J, Chimukangara B,  
1179 Lessells, RJ, Moosa Y, Gazy, I, Early transmission of SARS-CoV-2 in South Africa: An epidemiological  
1180 and phylogenetic report. medRxiv. 2020; <https://doi.org/10.1101/2020.05.29.20116376>.
- 1181 [53] Fares WA, Kais C, Ghedira, M Dorra G, Sondes R, Imene HB, Henda BD Walid T, Zina H, Amel M,  
1182 Nahed S, Imen H, Aurelia A, Veronique K, Guillain H, Valerie M, Jean-Claude C, Nissaf M, Alaya B, Triki  
1183 H. 2020. First whole genome sequences and phylogenetic analysis of SARS-CoV-2 virus isolates during  
1184 COVID-19 outbreak in Tunisia, North Africa. *Authorea Preprints*.  
1185 <https://doi.org/10.22541/au.159137642.26983355>.
- 1186 [54] Yin C. Genotyping coronavirus SARS-CoV-2: methods and implications. *Genomics*. 2020; 112(5): 3588-  
1187 3596. <https://doi.org/10.1016/j.ygeno.2020.04.016>.

- 1188 [55] Wang C, Liu Z, Chen Z, Huang X, Xu M, He T, Zhang, Z. The establishment of reference sequence for  
1189 SARS-CoV-2 and Fvariation analysis. *Journal of J Med Virol*. 2020; 92(6):667-674.  
1190 <https://doi.org/10.1002/jmv.25762>.
- 1191 [56] Wu F, Zhao S, Yu B, Chen YM, Wang W, Song ZG, Hu Y, Tao, ZW, Tian JH, Pei,YY, Yuan ML. A new  
1192 coronavirus associated with human respiratory disease in China. *Nature*. 2020;579(7798), 265–269.  
1193 <https://doi.org/10.1038/s41586-020-2008-3>.
- 1194 [57] Wang Y, Mao JM, Wang GD, Qiu Z, Yao Q, Chen KP. Human SARS-CoV-2 has evolved to reduce CG  
1195 dinucleotide in its open reading frames. *Sci. Rep*. 2020;10(12331). [https://doi.org/10.1038/s41598-020-](https://doi.org/10.1038/s41598-020-69342-y)  
1196 [69342-y](https://doi.org/10.1038/s41598-020-69342-y).
- 1197 [58] Mercatelli D, Giorgi FM. Geographic and Genomic Distribution of SARS-CoV-2 Mutations. *Front*.  
1198 *Microbiol*. 2020; <https://doi.org/10.3389/fmicb.2020.01800>.
- 1199 [59] Sjaarda CP, Rustom N, Huang D, Perez-Patrigeon S, Hudson ML, Wong H, Guan TH, Ayub M, Soares CN,  
1200 Colautti RI, Evans GA. Chasing the origin of SARS-CoV-2 in Canada’s COVID-19 cases: A genomics  
1201 study. *bioRxiv*. 2020; <https://doi.org/10.1101/2020.06.25.171744>.
- 1202 [60] De Maio N, Walker C, Borges R, Weilguny L, Slodkowicz G, Goldman, N. Issues with SARS-CoV-2  
1203 sequencing data. 2020; <https://virological.org/t/issues-with-sars-cov-2-sequencing-data/473>.
- 1204 [61] Kuehn BM. Genetic Analysis Tracks SARS-CoV-2 Mutations in Human Hosts. *JAMA*. 2020; 323(23),  
1205 2363-2363. <https://doi.org/10.1001/jama.2020.9825>.
- 1206 [62] Ackermann MR, Blömer J, Kuntze D, Sohler C. Analysis of agglomerative clustering. *Algorithmica*. 2014  
1207 69(1), 184-215. <https://doi.org/10.1007/s00453-012-9717-4>.
- 1208 [63] Inyang UG, Akpan EE, Akinyokun OC. A Hybrid Machine Learning Approach for Flood Risk Assessment  
1209 and Classification. *J. Comput*. 2020; 19(2), 1-20. <https://doi.org/10.1142/S1469026820500121>.
- 1210 [64] Abascal F, Zardoya R, Telford M J. TranslatorX: multiple alignment of nucleotide sequences guided by  
1211 amino acid translations. *Nucleic Acids Res. Spec. Publ*. 2010; 38(suppl\_2), W7-W13.  
1212 <https://doi.org/10.1093/nar/gkq291>.
- 1213 [65] Vesanto J, Ahola J. Hunting for Correlations in Data Using the Self-Organizing Map. *Proceeding of the*  
1214 *International ICSC Congress on Computational Intelligence Methods and Applications*.1999; p. 279–285.
- 1215 [66] Friedman M. A comparison of alternative tests of significance for the problem of m rankings. *The Annals*  
1216 *of Mathematical Statistics*. 1940. 11(1):86-92. <https://www.jstor.org/stable/2235971>.
- 1217 [67] Pohlert T. The pairwise multiple comparison of mean ranks package (PMCMR). *R package*. 2014. 1-9.

1218 [68] Chadwick, R. Gender and the human genome. *Mens Sana Monographs*.2009; 7(1), 10-19.  
1219 <https://doi.org/10.4103 /0973 -1229.44075>.

1220 [69] Chen, X., Hu, W., Ling, J., Mo, P., Zhang, Y., Jiang, Q., Ma, Z., Cao, Q., Deng, L., Song, S. and Zheng, R.  
1221 Hypertension and diabetes delay the viral clearance in COVID-19 patients. *medRxiv*. (2020).  
1222 <https://doi.org/10.1101/2020.03.22.20040774>.

1223 [70] Shastri, A., Wheat, J., Agrawal, S., Chaterjee, N., Pradhan, K., Goldfinger, M., Kornblum, N., Steidl, U.,  
1224 Verma, A. and Shastri, J. Delayed clearance of SARS-CoV2 in male compared to female patients: High  
1225 ACE2 expression in testes suggests possible existence of gender-specific viral reservoirs. (2020). *medRxiv*.  
1226 <https://doi.org/10.1101/2020.04.16.20060566>.

1227 [71] Schaibley VM, Zawistowski, M., Wegmann, D., Ehm, M. G., Nelson, M. R., Jean, P. L. S., Abecasis, G. R.,  
1228 Novembre, J., Zöllner, S. and Li, J. Z. (2013). The influence of genomic context on mutation patterns in the  
1229 human genome inferred from rare variants. *Genome Res*. 23(12): 1974-1984.  
1230 <https://doi.org/doi/10.1101/gr.154971.113>.

1231 [72] Nachman MW, Crowell SL. Estimate of the mutation rate per nucleotide in humans. *Genetics*. 2000;156(1),  
1232 297-304. <https://www.genetics.org/content/genetics/156/1/297.full.pdf>.

1233 [73] Hodgkinson A, Ladoukakis E, Eyre-Walker A. Cryptic variation in the human mutation rate. *PLoS Biol*.  
1234 2009; 7(2), e1000027. <https://doi.org/10.1371/journal.pbio.1000027>

1235 [74] Nachman MW. Haldane and the first estimates of the human mutation rate. *J. Genet*. 2004; 83(3), 231-233.  
1236 <http://doi.org/10.1007/bf02717891>.

1237 [75] Armstrong J, Fiddes IT, Diekhans M, Paten B. Whole-genome alignment and comparative annotation.  
1238 *Annu Rev Anim Biosci*. 2019; 7: 41-64. <https://doi.org/10.1038/nmeth.1179>.

1239 [76] Wittler, R. Alignment-and reference-free phylogenomics with colored de Bruijn graphs. *Algorithms Mol*  
1240 *Biol*. 2020; 15, 1-12. <https://doi.org/10.1186/s13015-020-00164-3>.

1241 [77] Feng S, Zhongxi M. Similarity among nucleotide sequences. *Acta Biotheor*. 2002;50(2):95-99.  
1242 <https://doi.org/10.1023/a:1016376910987>.

1243 [78] Amagata T. Natural Products Structural Diversity-II Secondary Metabolites: Sources, Structures and  
1244 Chemical Biology. *Comprehensive Natural Products II*. 2010; 2:581-621.

1245 [79] Koyama T, Platt D, Parida L. Variant analysis of SARS-CoV-2 genomes. *Bull World Health Organ*. 2020;  
1246 98(7): 495–504. <https://doi.org/10.2471/BLT.20.253591>.

1247

1249 **Table 1. Summary of data source, transmission history and identified intra-country sub-strains of African genome isolates**

Author and year	Country/ countries	Isolates considered and source	Transmission history	Identified intra-country sub-strains	Additional information
[47]	Algeria	Information on cases with confirmed COVID-19 infection based on official reports from governmental institutes in Algeria	First case of COVID-19 was reported on when an Italian national tested positive in Ouargla region in the south of the country, later, 2 cases were reported in Blida region in the North of Algeria, following their contacts with two Algerian nationals who came from France.	No information	-
[48]	Mali	21 whole genome sequences were generated from 38 positive isolates	First two COVID-19 cases living in Bamako and Kayes, both returning from France were confirmed. Further 2 lineages, 1 (Asia, Europe, Mali, Oceania) and 2 (Europe, USA, Canada, Northern Africa), were identified	Patients returning from Tunisia, were identified.	The presence of 2 lineages strongly suggests at least two different and independent introduction points of the SARS-CoV-2 infection in the Sahel region.
[49]	15 African countries (Algeria, Benin, Cameroon, Democratic Republic of the Congo (DRC), Gambia, Ghana, Kenya, Madagascar, Mali, Nigeria, Senegal, Sierra Leone, South Africa, Uganda, Zambia.)	SARS-CoV-2 genomes (n=1340) from 15 countries in WHO Africa region (out of 47), excavated from GSAID database on August 17, 2020, and representing ~2% of publicly available sequences globally. Majority of genomes were from South Africa (51%), DRC (20%) and Senegal (10%).	Multiple separate introductions into Africa from other continents, of which 72% came from Europe, 19% from Asia, 6% from South America and 3% from North America.	An estimated 123 introductions came from other continents while 14 introductions were between African countries.	93% of African SARS-CoV-2 genomes have the D614G mutation in spike protein.

Author and year	Country/ countries	Isolates considered and source	Transmission history	Identified intra-country sub-strains	Additional information
[38]	North Africa (Morocco, Tunisia, Algeria, Egypt)	40 SARS-CoV-2 genomes (28 from Morocco, 7 from Tunisia, 3 from Algeria, 2 from Egypt), were downloaded from GISAID between March 03, 2020 and May 15, 2020.	Phylogenetic analysis showed that the Moroccan and Tunisian SARS-CoV-2 strains were closely related to those from different origins (Asia, Europe, North and South America) and distributed in different distinct subclades.	Moroccan and Tunisian strains were closely related to those from different continents, which could indicate different sources of infection with no single dominant strain circulating in Morocco.	2 waves of SARS-CoV-2 infections were recorded, with the first one mostly imported from Europe and the second one dominated by local infections.
[40]	Morocco	6 genome sequences of SARS-CoV-2 strains deposited in DDBJ/ENA/GenBank and GenBank.	The mutation in Morocco was associated with the observed transmission increase in the United States	No information	Mutation analysis revealed the presence of the spike D614G mutation in all six genomes, which is widely present in several genomes around the world.
[39]	Morocco	22 genome sequences reported by three different laboratories in Morocco up to June 7, 2020, as well as 40366 genomes from all around the world.	Introduction of SARS-CoV-2 strains to Morocco came from Belgium, Spain and France at the beginning of the epidemic. Later, strains from USA and Vietnam were noticed after the lockdown, possibly, through sea trades.	Spread in Morocco did not show a predominant SARS-CoV-2 route during analysis.	Different SARS-CoV-2 strains, with different mutation patterns, coexist in Morocco.
[41]	Morocco	1 isolate obtained from Moroccan patient infected in Casablanca.	Sequence belonged a clade which is similar to cases in Florida	No information	Currently circulating strains in Morocco came from different countries with a local evolution.
[50]	DRC	127 SARS-CoV-2 genome sequences excavated from GISAID database	7 distinct lineages which diversity originated from China, are distributed across 4 continents as follows: Asia (China, India, Jordan, Iran);	Ghana constituted one of the large lineage group.	Sampled cases from the DRC are the result of repeated introduction of the virus from a range of locations

Author and year	Country/ countries	Isolates considered and source	Transmission history	Identified intra-country sub-strains	Additional information
			Europe (Italy, Austria, UK); Australia; North America (USA); Africa (Ghana)		followed by local transmission.
[37]	Egypt	2 isolates obtained from Egyptians residing in Upper and Lower Egypt	Egyptian clades fell into A2a denoting strains from Asia, Europe, Africa, Australia and USA	Senegal and DRC were the closely associated strains from Africa.	Sequence analysis showed mutations that differentiate Egyptian strains from the reference strain 2019-nCoV WHU01.
[51]	Kenya	122 genome sequences from Nairobi (n=102) and Coastal Kenya (n=20) collected between 12th March and 30th April 2020.	Evidence of 10 global lineages in China with several global exports to Iran, France, UK and Italy, with multiple introduction of European-centric lineages into the country.	Local transmission, between Nairobi and Mombasa, were found. Infection from the Coastal Kenya was imported from Nairobi	Sequencing of additional SARS-CoV-2 genomes in Kenya will provide a more detailed picture of local transmission patterns.
[43]	Nigeria	24 genome sequences (18 full, 6 partial) were assembled for experiment	7 lineages circulating in Nigeria are from Asia, Europe, USA, Middle East, Australia and Africa.	Strains from Egypt, Senegal and Ghana were revealed through Phylogenetic analysis	Four of the new sequences clustered closely together and formed a separate clade which strongly suggests local community transmission.
[42]	Nigeria	39 complete genomes of SARS-CoV-2 were retrieved from GISAID	First confirmed case in Nigeria came from an infected traveler from Italy.	The strain from Nigeria was found in the Wuhan subclade 3 together with some strains from Congo.	Strain in Nigeria clustered in a monophyletic clade with a Wuhan sublineage.
[52]	South Africa	21 SARS-CoV-2 whole genome, sampled in the first port of entry, KwaZulu-Natal (KZN), during the first month of the epidemic.	First COVID-19 case was a South African citizen returning home from a skiing holiday in Italy.	No information	13 independent introductions in KZN, which lineages revealed imported infection from Europe and North America.

Author and year	Country/ countries	Isolates considered and source	Transmission history	Identified intra-country sub-strains	Additional information
[44]	South Africa	Complete genome sequence of a SARS-CoV-2 isolate obtained from a South African patient.	The patient returned from Italy	No information	This sequence has been deposited in GenBank
[53]	Tunisia	SARS-CoV-2 strains were obtained from imported and locally transmission cases	Patients had travel history from Vietnam, Turkey and France	The Turkey sequence grouped with sequences showing lineage from Egypt	272 imported and 799 locally transmitted cases, have been diagnosed and tested positive for SARS-CoV-2
[46]	Uganda	20 genomic sequences from 14 persons entering Uganda	Miami to Istanbul, UK to Netherlands to Rwanda, Kenya, Tanzania	Kenya, Tanzania, and South Sudan	6 lineages among were imported into Uganda.

1250

1251

**Table 2. Cognitive map for male isolates**

Algeria	Algeria	Benin	Benin	Benin	Benin	Benin	Cameroon	DRC	DRC
Blida (1)	Blida (2)	Cotonou (3)	Cotonou (4)	Cotonou (5)	Cotonou (6)	Oueme (7)	Yaounde (8)	Ht-Katanga (9)	Ht-Katanga (10)
2,68,72	1,68,72	29,32,33,34,42,43, 44,45,46,47,65	5,6,8,13,14,15, 18,19,20,23,24,30, 49,50,56,61,63,64, 69,70,79,82,83,99, 111,117,118,119, 132,135,139, 141, 143,144	4,6,8,13,14,15, 18,19,20,23,24,30, 49,50,56,61,63,64, 69,70,79,82,83,99, 111,117,118,119, 132,135,139, 141, 143,144	4,5,8,13,14,15, 18,19,20,23,24,30, 49,50,56,61,63,64, 69,70,79,82,83,99, 111,117,118,119, 132,135,139, 141, 143,144	7,17,28,35,37,38, 39,40,41,48,54,55, 58,59,60,62,66,67, 71,74,89,102,112, 113,115,116,122, 125,128,140,142,145	4,5,6,13,14,15, 18,19,20,23,24,30, 49,50,56,61,63,64, 69,70,79,82,83,99, 111,117,118,119, 132,135,139, 141, 143,144	10,11,12,73, 76,77,78,81,87,88, 90,91,92,93,95,96, 101,104,106,109, 127,129,134,136, 137,138	9,11,12,73, 76,77,78,81,87,88, 90,91,92,93,95,96, 101,104,106,109, 127,129,134,136, 137,138
DRC	DRC	DRC	DRC	DRC	DRC	DRC	DRC	DRC	DRC
Ht-Katanga (11)	Kinshasha (12)	Kinshasha (13)	Kinshasha (14)	Kinshasha (15)	Kongo Central (16)	Kongo Central (17)	Kongo Central (18)	Kongo Central (19)	Nord Kivu (20)
9,10,12,73, 76,77,78,81,87,88, 90,91,92,93,95,96, 101,104,106,109, 127,129,134,136, 137,138	9,10,11,73, 76,77,78,81,87,88, 90,91,92,93,95,96, 101,104,106,109, 127,129,134,136, 137,138	4,5,6,8,14,15, 18,19,20,23,24,30, 49,50,56,61,63,64, 69,70,79,82,83,99, 111,117,118,119, 132,135,139, 141, 143,144	4,5,6,8,14,15, 18,19,20,23,24,30, 49,50,56,61,63,64, 69,70,79,82,83,99, 111,117,118,119, 132,135,139, 141, 143,144	4,5,6,8,13,14, 18,19,20,23,24,30, 49,50,56,61,63,64, 69,70,79,82,83,99, 111,117,118,119, 132,135,139, 141, 143,144	21,22,31,36,75, 80,84,94,100,105, 107,114,120,124,133	7,28,35,37,38, 39,40,41,48,54,55, 58,59,60,62,66,67, 71,74,89,102,112, 113,115,116,122, 125,128,140,142,145	4,5,6,8,13,14,15, 19,20,23,24,30, 49,50,56,61,63,64, 69,70,79,82,83,99, 111,117,118,119, 132,135,139, 141, 143,144	4,5,6,8,13,14,15, 18,20,23,24,30, 49,50,56,61,63,64, 69,70,79,82,83,99, 111,117,118,119, 132,135,139, 141, 143,144	4,5,6,8,13,14,15, 18,19,23,24,30, 49,50,56,61,63,64, 69,70,79,82,83,99, 111,117,118,119, 132,135,139, 141, 143,144
DRC	DRC	DRC	Egypt	Egypt	Egypt	Egypt	Egypt	Gambia	Gambia
Sud Kivu (21)	Sud Kivu (22)	Sud Kivu (23)	Cairo (24)	Cairo (25)	Cairo (26)	Cairo (27)	Kalyoubia (28)	Kombo (29)	Kombo (30)
16,22,31,36,75, 80,84,94,100,105, 107,114,120,124,133	16,21,31,36,75, 80,84,94,100,105, 107,114,120,124,133	4,5,6,8,13,14,15, 18,19,20,24,30, 49,50,56,61,63,64, 69,70,79,82,83,99,	4,5,6,8,13,14,15, 18,19,20,23,30, 49,50,56,61,63,64, 69,70,79,82,83,99,	26,27,51,52,103, 25,27,51,52,103	25,27,51,52,103, 25,26,51,52,103	25,26,51,52,103, 25,26,51,52,103	7,17,35,37,38, 39,40,41,48,54,55, 58,59,60,62,66,67, 71,74,89,102,112,	3,32,33,34,42,43, 44,45,46,47,65	4,5,6,8,13,14,15, 18,19,20,23,24, 49,50,56,61,63,64, 69,70,79,82,83,99,

		111,117,118,119,	111,117,118,119,				113,115,116,122,		111,117,118,119,
		132,135,139, 141,	132,135,139, 141,				125,128,140,142,145		132,135,139, 141,
		143,144	143,144						143,144
Gambia	Gambia	Gambia	Gambia	Ghana	Ghana	Ghana	Ghana	Kenya	Kenya
Kombo (31)	Kombo (32)	W.Coast Reg (33)	W.Coast Reg (34)	Greater Ghana (35)	Greater Ghana (36)	Greater Ghana (37)	Greater Ghana (38)	Unknown (39)	Unknown (40)
16,21,22,36,75,	3,29,33,34,42,43,	3,29,32,34,42,43,	3,29,32,33,42,43,	7,17,28,37,38,	16,21,22,31,75,	7,17,28,35,38,	7,17,28,35,37	7,17,28,35,37,38,	7,17,28,35,37,38,
80,84,94,100,105,	44,45,46,46,65	44,45,46,46,65	44,45,46,46,65	39,40,41,48,54,55,	80,84,94,100,105,	39,40,41,48,54,55,	39,40,41,48,54,55,	40,41,48,54,55,	39,41,48,54,55,
107,114,120,124,133				58,59,60,62,66,67,	107,114,120,124,133	58,59,60,62,66,67,	58,59,60,62,66,67,	58,59,60,62,66,67,	58,59,60,62,66,67,
				71,74,89,102,112,		71,74,89,102,112,	71,74,89,102,112,	71,74,89,102,112,	71,74,89,102,112,
				113,115,116,122,		113,115,116,122,	113,115,116,122,	113,115,116,122,	113,115,116,122,
				125,128,140,142,145		125,128,140,142,145	125,128,140,142,145	125,128,140,142,145	125,128,140,142,145
Kenya	Kenya	Mali	Mali	Mali	Mali	Mali	Morocco	Morocco	Morocco
Unknown (41)	Unknown (42)	Bamako (43)	Bamako (44)	Bamako (45)	Bamako (46)	Mopti (47)	Casablanca (48)	Ouarzazate (49)	Rabat (50)
7,17,28,35,37,38,	3,29,32,33,34,43,	3,29,32,33,34,42,	3,29,32,33,34,42,	3,29,32,33,34,42,	3,29,32,33,34,42,	3,29,32,33,34,42,	7,17,28,35,37,38,	4,5,6,8,13,14,15,	4,5,6,8,13,14,15,
39,40,48,54,55,	44,45,46,47,65	44,45,46,47,65	43,45,46,47,65	43,44,46,47,65	43,44,45,47,65	43,44,45,46,65	39,40,41,54,55,	18,19,20,23,24,30,	18,19,20,23,24,30,
58,59,60,62,66,67,							58,59,60,62,66,67,	50,56,61,63,64,	49,56,61,63,64,
71,74,89,102,112,							71,74,89,102,112,	69,70,79,82,83,99,	69,70,79,82,83,99,
113,115,116,122,							113,115,116,122,	111,117,118,119,	111,117,118,119,
125,128,140,142,145							125,128,140,142,145	132,135,139, 141,	132,135,139, 141,
								143,144	143,144
Morocco	Morocco	Nigeria	Nigeria	Nigeria	Nigeria	Nigeria	Nigeria	Nigeria	Nigeria
Rabat (51)	Rabat (52)	Kwara (53)	Kwara (54)	Kwara (55)	Lagos (56)	Osun (57)	Osun (58)	Osun (59)	Osun (60)
25,26,27,52,103	25,26,27,51,103	123,126,131	7,17,28,35,37,38,	7,17,28,35,37,38,	4,5,6,8,13,14,15,	85,86,108,110	7,17,28,35,37,38,	7,17,28,35,37,38,	7,17,28,35,37,38,
			39,40,41,48,55,	39,40,41,48,54,	18,19,20,23,24,30,		39,40,41,48,54,55,	39,40,41,48,54,55,	39,40,41,48,54,55,
			58,59,60,62,66,67,	58,59,60,62,66,67,	49,50,61,63,64,		59,60,62,66,67,	58,60,62,66,67,	58,59,62,66,67,

			71,74,89,102,112,	71,74,89,102,112,	69,70,79,82,83,99,		71,74,89,102,112,	71,74,89,102,112,	71,74,89,102,112,
			113,115,116,122,	113,115,116,122,	111,117,118,119,		113,115,116,122,	113,115,116,122,	113,115,116,122,
			125,128,140,142,145	125,128,140,142,145	132,135,139, 141,		125,128,140,142,145	125,128,140,142,145	125,128,140,
					143,144				142,145
Nigeria	Nigeria	Nigeria	Nigeria	Nigeria	Senegal	Senegal	Senegal	Senegal	Senegal
Oyo (61)	Oyo (62)	Oyo (63)	Oyo (64)	Unknown (65)	Dakar (66)	Dakar (67)	Dakar (68)	Dakar (69)	Diourbel (70)
4,5,6,8,13,14,15,	7,17,28,35,37,38,	4,5,6,8,13,14,15,	4,5,6,8,13,14,15,	3,29,32,33,34,42,	7,17,28,35,37,38,	7,17,28,35,37,38,	1,2,72	4,5,6,8,13,14,15,	4,5,6,8,13,14,15,
18,19,20,23,24,30,	39,40,41,48,54,55,	18,19,20,23,24,30,	18,19,20,23,24,30,	43,44,45,46,47	39,40,41,48,54,55,	39,40,41,48,54,55,		18,19,20,23,24,30,	18,19,20,23,24,30,
49,50,56,63,64,	58,59,60,66,67,	49,50,56,61,64,	49,50,56,61,63,		58,59,60,62,67,	58,59,60,62,66,		49,50,56,61,63,64,	49,50,56,61,63,64,
69,70,79,82,83,99,	71,74,89,102,112,	69,70,79,82,83,99,	69,70,79,82,83,99,		71,74,89,102,112,	71,74,89,102,112,		70,79,82,83,99,	69,79,82,83,99,
111,117,118,119,	113,115,116,122,	111,117,118,119,	111,117,118,119,		113,115,116,122,	113,115,116,122,		111,117,118,119,	111,117,118,119,
132,135,139, 141,	125,128,140,142,145	132,135,139, 141,	132,135,139, 141,		125,128,140,142,145	125,128,140,142,145		132,135,139, 141,	132,135,139, 141,
143,144		143,144	143,144					143,144	143,144
Senegal	Senegal	Senegal	Senegal	Senegal	Senegal	Senegal	Senegal	Senegal	South Africa
Diourbel (71)	Diourbel (72)	Kolda (73)	Pikine (74)	St Louis (75)	St Louis (76)	St Louis (77)	Thies (78)	Thies (79)	Amajuba (80)
7,17,28,35,37,38,	1,2,68	9,10,11,12,	7,17,28,35,37,38,	16,21,22,31,36,	9,10,11,12,73,	9,10,11,12,73,	9,10,11,12,73,	4,5,6,8,13,14,15,	16,21,22,31,36,75,
39,40,41,48,54,55,		76,77,78,81,87,88,	39,40,41,48,54,55,	80,84,94,100,105,	77,78,81,87,88,	76,78,81,87,88,	76,77,81,87,88,	18,19,20,23,24,30,	84,94,100,105,
58,59,60,62,66,67,		90,91,92,93,95,96,	58,59,60,62,66,67,	107,114,120,124,133	90,91,92,93,95,96,	90,91,92,93,95,96,	90,91,92,93,95,96,	49,50,56,61,63,64,	107,114,120,124,
74,89,102,112,		101,104,106,109,	71,89,102,112,		101,104,106,109,	101,104,106,109,	101,104,106,109,	69,70,79,82,83,99,	133
113,115,116,122,		127,129,134,136,	113,115,116,122,		127,129,134,136,	127,129,134,136,	127,129,134,136,	111,117,118,119,	
125,128,140,142,145		137,138	125,128,140,142,145		137,138	137,138	137,138	132,135,139, 141,	
								143,144	
South Africa	South Africa	South Africa	South Africa	South Africa	South Africa	South Africa	South Africa	South Africa	South Africa
Amajuba (81)	Amajuba (82)	Amajuba (83)	Berea (84)	Berea (85)	Berea (86)	Cape Town (87)	Cape Town (88)	Cape Town (89)	Cape Town (90)

9,10,11,12,73,	4,5,6,8,13,14,15,	4,5,6,8,13,14,15,	16,21,22,31,36,75,	57,86,108,110	57,85,108,110	9,10,11,12,73,	9,10,11,12,73,	7,17,28,35,37,38,	9,10,11,12,73,
76,77,78,87,88,	18,19,20,23,24,30,	18,19,20,23,24,30,	80,94,100,105,			76,77,78,88,	76,77,78,81,87,	39,40,41,48,54,55,	76,77,78,81,87,88,
90,91,92,93,95,96,	49,50,56,61,63,64,	49,50,56,61,63,64,	107,114,120,124,133			90,91,92,93,95,96,	90,91,92,93,95,96,	58,59,60,62,66,67,	91,92,93,95,96,
101,104,106,109,	69,70,79,83,99,	69,70,79,82,99,				101,104,106,109,	101,104,106,109,	71,74,102,112,	101,104,106,109,
127,129,134,136,	111,117,118,119,	111,117,118,119,				127,129,134,136,	127,129,134,136,	113,115,116,122,	127,129,134,136,
137,138	132,135,139, 141,	132,135,139, 141,				137,138	137,138	125,128,140,142,145	137,138
	143,144	143,144							
South Africa	South Africa	South Africa	South Africa	South Africa	South Africa	South Africa	South Africa	South Africa	South Africa
EC (91)	eThekwini (92)	eThekwini (93)	eThekwini (94)	eThekwini (95)	Harry Gwala (96)	Harry Gwala (97)	Harry Gwala (98)	Harry Gwala (99)	Ilembe (100)
9,10,11,12,73,	9,10,11,12,73,	9,10,11,12,73,	16,21,22,31,36,75,	9,10,11,12,73,	9,10,11,12,73,	98,121,130	97,121,130	4,5,6,8,13,14,15,	16,21,22,31,36,75,
76,77,78,81,87,88,	76,77,78,81,87,88,	76,77,78,81,87,88,	80,84,100,105,	76,77,78,81,87,88,	76,77,78,81,87,88,			18,19,20,23,24,30,	80,84,94,105,
90,92,93,95,96,	90,91,93,95,96,	90,91,92,95,96,	107,114,120,124,133	90,91,92,93,96,	90,91,92,93,95,			49,50,56,61,63,64,	107,114,120,124,
101,104,106,109,	101,104,106,109,	101,104,106,109,		101,104,106,109,	101,104,106,109,			69,70,79,82,83,	133
127,129,134,136,	127,129,134,136,	127,129,134,136,		127,129,134,136,	127,129,134,136,			111,117,118,119,	
137,138	137,138	137,138		137,138	137,138			132,135,139, 141,	
								143,144	
South Africa	South Africa	South Africa	South Africa	South Africa	South Africa	South Africa	South Africa	South Africa	South Africa
Ilembe (101)	Ilembe (102)	Ilembe (103)	Kg Cetshwayo (104)	Kg Cetshwayo (105)	Kg Cetshwayo (106)	Kg Cetshwayo (107)	KZN (108)	KZN (109)	KZN (110)
9,10,11,12,73,	7,17,28,35,37,38,	25,26,27,51,52	9,10,11,12,73,	16,21,22,31,36,75,	9,10,11,12,73,	16,21,22,31,36,75,	57,85,86,110	9,10,11,12,73,	57,85,86,108
76,77,78,81,87,88,	39,40,41,48,54,55,		76,77,78,81,87,88,	80,84,94,100,	76,77,78,81,87,88,	80,84,94,100,105,		76,77,78,81,87,88,	
90,91,92,93,95,96,	58,59,60,62,66,67,		90,91,92,93,95,96,	107,120,124,133	90,91,92,93,95,96,	120,124,133		90,91,92,93,95,96,	
104,106,109,	71,74,89,112,		101,106,109,		101,104,109,			101,104,106,	
127,129,134,136,	113,115,116,122,		127,129,134,136,		127,129,134,136,			127,129,134,136,	
137,138	125,128,140,142,145		137,138		137,138			137,138	
South Africa	South Africa	South Africa	South Africa	South Africa	South Africa	South Africa	South Africa	South Africa	South Africa

KZN (111)	MP (112)	North West (113)	North West (114)	North West (115)	North West (116)	Stanger (117)	Ugu (118)	Ugu (119)	Ugu (120)
4,5,6,8,13,14,15,	7,17,28,35,37,38,	7,17,28,35,37,38,	16,21,22,31,75,	7,17,28,35,37,38,	7,17,28,35,37,38,	4,5,6,8,13,14,15,	4,5,6,8,13,14,15,	4,5,6,8,13,14,15,	16,21,22,31,36,75,
18,19,20,23,24,30,	39,40,41,48,54,55,	39,40,41,48,54,55,	80,84,94,100,105,	39,40,41,48,54,55,	39,40,41,48,54,55,	18,19,20,23,24,30,	18,19,20,23,24,30,	18,19,20,23,24,30,	80,84,94,100,105,
49,50,56,61,63,64,	58,59,60,62,66,67,	58,59,60,62,66,67,	107,120,124,133	58,59,60,62,66,67,	58,59,60,62,66,67,	49,50,56,61,63,64,	49,50,56,61,63,64,	49,50,56,61,63,64,	107,124,133
69,70,79,82,83,99,	71,74,89,102,	71,74,89,102,112,115,		71,74,89,102,112,	71,74,89,102,112,	69,70,79,82,83,99,	69,70,79,82,83,99,	69,70,79,82,83,99,	
117,118,119,	113,115,116,122,	116,122,125,128,140,		113,116,122,	113,115,122, 125,	111,118,119,	111,117,119,	111,117,118,	
132,135,139, 141,	125,128,140,142,145	142,145		125,128,140,142,145	128,140,142,145	132,135,139, 141,	132,135,139, 141,	132,135,139, 141,	
143,144						143,144	143,144	143,144	
South Africa	South Africa	South Africa	South Africa	South Africa	South Africa	South Africa	South Africa	South Africa	South Africa
Ugu (121)	Umbilo (122)	Umgungundl (123)	Umgungundl (124)	Umgungundl (125)	Umgungundl (126)	Umkhanyak (127)	Umkhanyak (128)	Umkhanyak (129)	Umkhanyak (130)
97,98,130	7,17,28,35,37,38,	53,126,131	16,21,22,31,36,75,	7,17,28,35,37,38,	53,123,131	9,10,11,12,73,	7,17,28,35,37,38,	9,10,11,12,73,	97,98,121
	39,40,41,48,54,55,		80,84,94,100,105,	39,40,41,48,54,55,		76,77,78,81,87,88,	39,40,41,48,54,55,	76,77,78,81,87,88,	
	58,59,60,62,66,67,		107,120,133	58,59,60,62,66,67,		90,91,92,93,95,96,	58,59,60,62,66,67,	90,91,92,93,95,96,	
	71,74,89,102,112,			71,74,89,102,112,		101,104,106,109,	71,74,89,102,112,	101,104,106,109,	
	113,115,116,125,			113,115,116,122,		129,134,136,	113,115,116,122,	127,134,136,	
	128,140,142,145			128,140,142,145		137,138	125,140,142,145	137,138	
South Africa	South Africa	South Africa	South Africa	South Africa	South Africa	South Africa	South Africa	South Africa	South Africa
Umzinyathi (131)	Umzinyathi (132)	Umzinyathi (133)	Umzinyathi (134)	Uthukela (135)	Uthukela (136)	Uthukela (137)	Uthukela (138)	Uthungulu (139)	Uthungulu (140)
53,123,126	4,5,6,8,13,14,15,	16,21,22,31,36,75,	9,10,11,12,73,	4,5,6,8,13,14,15,	9,10,11,12,73,	9,10,11,12,73,	9,10,11,12,73,	4,5,6,8,13,14,15,	7,17,28,35,37,38,
	18,19,20,23,24,30,	80,84,94,100,105,	76,77,78,81,87,88,	18,19,20,23,24,30,	76,77,78,81,87,88,	76,77,78,81,87,88,	76,77,78,81,87,88,	18,19,20,23,24,30,	39,40,41,48,54,55,
	49,50,56,61,63,64,	107,120,124	90,91,92,93,95,96,	49,50,56,61,63,64,	90,91,92,93,95,96,	90,91,92,93,95,96,	90,91,92,93,95,96,	49,50,56,61,63,64,	58,59,60,62,66,67,
	69,70,79,82,83,99,		101,104,106,109,	69,70,79,82,83,99,	101,104,106,109,	101,104,106,109,	101,104,106,109,	69,70,79,82,83,99,	71,74,89,102,112,
	111,117,118,119,		127,129,136,	111,117,118,119,	127,129,134,	127,129,134,136,	127,129,134,136,	111,117,118,119,	113,115,116,122,
	135,139, 141,		137,138	132, 139, 141,	137,138	138	137	132,135, 141,	125,128,142,145
	143,144			143,144				143,144	

South Africa	South Africa	South Africa	South Africa	Tunisia
Zululand (141)	Zululand (142)	Zululand (143)	Zululand (144)	Bizerte (145)
4,5,6,8,13,14,15,	7,17,28,35,37,38,	4,5,6,8,13,14,15,	4,5,6,8,13,14,15,	7,17,28,35,37,38,
18,19,20,23,24,30,	39,40,41,48,54,55,	18,19,20,23,24,30,	18,19,20,23,24,30,	39,40,41,48,54,55,
49,50,56,61,63,64,	58,59,60,62,66,67,	49,50,56,61,63,64,	49,50,56,61,63,64,	58,59,60,62,66,67,
69,70,79,82,83,99,	71,74,89,102,112,	69,70,79,82,83,99,	69,70,79,82,83,99,	71,74,89,102,112,
111,117,118,119,	113,115,116,122,	111,117,118,119,	111,117,118,119,	113,115,116,122,
132,135,139,	125,128,140,145	132,135,139, 141,144	132,135,139, 141,	125,128,140,142
143,144			143	

**Table 3. Cognitive map for female isolates**

Algeria	Benin	Benin	Benin	Benin	DRC	DRC	DRC	DRC	DRC
Blida (1)	Cotonou (2)	Cotonou (3)	Cotonou (4)	Cotonou (5)	Haut Katanga (6)	Kinshasha (7)	Kinshasha (8)	Kinshasha (9)	Kinshasha (10)
8,22,48,50,51,55,	3,5,23,24,25,26,	2,5,23,24,25,26,	10,15,18,19,20,21,	2,3,23,24,25,26,	9,11,12,13,14,	35,37,38,45,	1,22,48,50,51,55,	6,11,12,13,14,	4,15,18,19,20,21,
61,81,86,121,124,	27,28,33,34,39,	27,28,33,34,39,	31,40,46,57,58,63,	27,28,33,34,39,	16,17,59,60,62,92,	54,56,64,73,77,	61,81,86,121,124,	16,17,59,60,62,92,	31,40,46,57,58,63,
128,131,137	41,52,53,145	41,52,53,145	66,67,68,70,71,72,	41,52,53,145	94,96,97,98,130	101,102,105,107,	128,131,137	94,96,97,98,130	66,67,68,70,71,72,
			78,79,82,93,100,			116,117,129,139			78,79,82,93,100,
			103,104,112,127,						103,104,112,127,
			133,136,140,143						133,136,140,143
DRC	DRC	DRC	DRC	DRC	DRC	DRC	DRC	Egypt	Egypt
Kongo Central (11)	Kongo Central (12)	Kongo Central (13)	Kongo Central (14)	Sud Kivu (15)	Sud Kivu (16)	Sud Kivu (17)	Sud Kivu (18)	Cairo (19)	Cairo (20)

6,9,12,13,14,	6,9,11,13,14,	6,9,11,12,14,	6,9,11,12,13,	4,10,18,19,20,21,	6,9,11,12,13,14,	6,9,11,12,13,14,	4,10,15,19,20,21,	4,10,15,18,20,21,	4,10,15,18,19,21,
16,17,59,60,62,92,	16,17,59,60,62,92,	16,17,59,60,62,92,	16,17,59,60,62,92,	31,40,46,57,58,63,	17,59,60,62,92,	16,59,60,62,92,	31,40,46,57,58,63,	31,40,46,57,58,63,	31,40,46,57,58,63,
94,96,97,98,130	94,96,97,98,130	94,96, 97,98,130	94,96,97,98,130	66,67,68,70,71,72,	94,96,97,98,130	94,96,97,98,130	66,67,68,70,71,72,	66,67,68,70,71,72,	66,67,68,70,71,72,
				78,79,82,93,100,			78,79,82,93,100,	78,79,82,93,100,	78,79,82,93,100,
				103,104,112,127,			103,104,112,127,	103,104,112,127,	103,104,112,127,
				133,136,140,143			133,136,140,143	133,136,140,143	133,136,140,143
Egypt	Egypt	Gambia	Gambia	Gambia	Gambia	Ghana	Ghana	Ghana	Ghana
Cairo (21)	Cairo (22)	Kombo (23)	Kombo (24)	Kombo (25)	West Coast Reg (26)	Greater Ghana (27)	Greater Ghana (28)	Greater Ghana (29)	Greater Ghana (30)
4,10,15,18,19,20,	1,8,48,50,51,55,	2,3,5,24,25,26,	2,3,5,23,25,26,	2,3,5,23,24,26,	2,3,5,23,24,25,	2,3,5,23,24,25,26,	2,3,5,23,24,25,26,	30,32,42,108,113	29,32,42,108,113
31,40,46,57,58,63,	61,81,86,121,124,	27,28,33,34,39,	27,28,33,34,39,	27,28,33,34,39,	27,28,33,34,39,	28,33,34,39,	27,33,34,39,		
66,67,68,70,71,72,	128,131,137	41,52,53,145	41,52,53,145	41,52,53,145	41,52,53,145	41,52,53,145	41,52,53,145		
78,79,82,93,100,									
103,104,112,127,									
133,136,140,143									
Kenya	Kenya	Kenya	Kenya	Madagascar	Madagascar	Mali	Mali	Mali	Mali
Unknown (31)	Unknown (32)	Unknown (33)	Unknown (34)	Antananarivo (35)	Fenoarivo (36)	Bamako (37)	Bamako (38)	Bamako (39)	Bamako (40)
4,10,15,18,19,20,	29,30,42,108,113	2,3,5,23,24,25,26,	2,3,5,23,24,25,26,	7,37,38,45,	69,75,76,122	7,35,38,45,	7,35,37,45,	2,3,5,23,24,25,26,	4,10,15,18,19,20,
21,40,46,57,58,63,		27,28,34,39,	27,28,33,39,	54,56,64,73,77,	135,138,142	54,56,64,73,77,	54,56,64,73,77,	27,28,33,34,	21,31,46,57,58,63,
66,67,68,70,71,72,		41,52,53,145	41,52,53,145	101,102,105,107,		101,102,105,107,	101,102,105,107,	41,52,53,145	66,67,68,70,71,72,
78,79,82,93,100,				116,117,129,139		116,117,129,139	116,117,129,139		78,79,82,93,100,
103,104,112,127,									103,104,112,127,
133,136,140,143									133,136,140,143
Morocco	Morocco	Nigeria	Nigeria	Nigeria	Nigeria	Nigeria	Nigeria	Nigeria	Nigeria
Rabat (41)	Rabat (42)	Ekiti (43)	Ekiti (44)	Ekiti (45)	Ogun (46)	Ondo (47)	Osun (48)	Osun (49)	Osun (50)

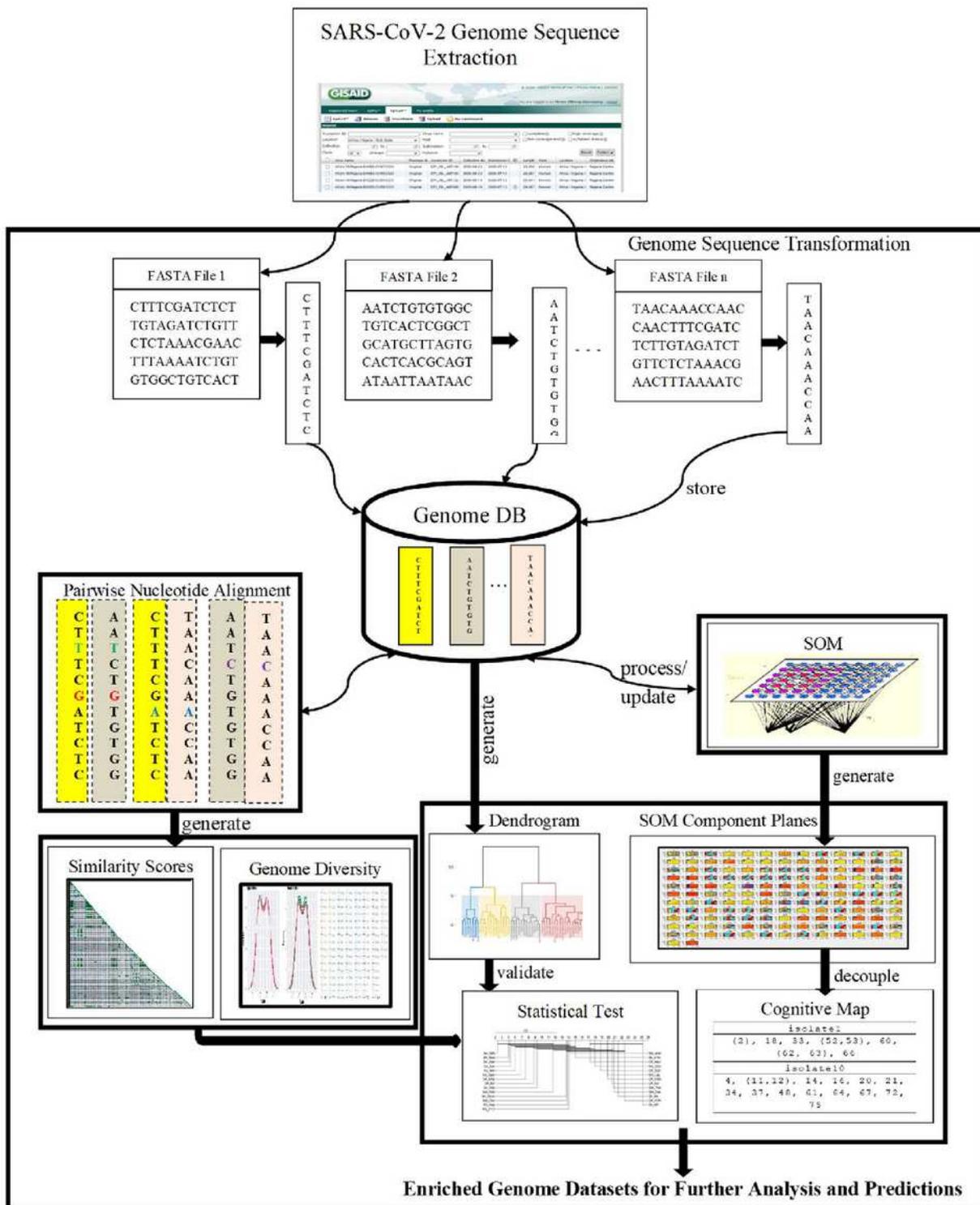
2,3,5,23,24,25,26,	29,30,32,108,113	44,47,49,74,	43,47,49,74,	7,35,37,38,	4,10,15,18,19,20,	43,44,49,74,	1,8,22,50,51,55,	43,44,47,74,	1,8,22,48,51,55,
27,28,33,34,39,		80,99,114,126	80,99,114,126	54,56,64,73,77,	21,31,40,57,58,63,	80,99,114,126	61,81,86,121,124,	80,99,114,126	61,81,86,121,124,
52,53,145				101,102,105,107,	66,67,68,70,71,72,		128,131,137		128,131,137
				116,117,129,139	78,79,82,93,100,				
					103,104,112,127,				
					133,136,140,143				
Nigeria	Nigeria	Nigeria	Senegal	Senegal	Senegal	Senegal	Senegal	Senegal	Senegal
Osun (51)	Unknown (52)	Unknown (53)	Dakar (54)	Dakar (55)	Dakar (56)	Dakar (57)	Diourbel (58)	Diourbel (59)	Diourbel (60)
1,8,22,48,50,55,	2,3,5,23,24,25,26,	2,3,5,23,24,25,26,	7,35,37,38,45,	1,8,22,48,50,51,	7,35,37,38,45,	4,10,15,18,19,20,	4,10,15,18,19,20,	6,9,11,12,13,14,	6,9,11,12,13,14,
61,81,86,121,124,	27,28,33,34,39,	27,28,33,34,39,	56,64,73,77,	61,81,86,121,124,	54,64,73,77,	21,31,40,46,58,63,	21,31,40,46,57,63,	16,17,60,62,92,	16,17,59,62,92,
128,131,137	41,53,145	41,52,145	101,102,105,107,	128,131,137	101,102,105,107,	66,67,68,70,71,72,	66,67,68,70,71,72,	94,96, 97,98,130	94,96, 97,98,130
			116,117,129,139		116,117,129,139	78,79,82,93,100,	78,79,82,93,100,		
						103,104,112,127,	103,104,112,127,		
						133,136,140,143	133,136,140,143		
Senegal	Senegal	Senegal	Senegal	South Africa	South Africa	South Africa	South Africa	South Africa	South Africa
Diourbel (61)	Thies (62)	Thies (63)	Thies (64)	Amajuba (65)	Amajuba (66)	Amajuba (67)	Amajuba (68)	Berea (69)	Berea (70)
1,8,22,48,50,51,55,	6,9,11,12,13,14,	4,10,15,18,19,20,	7,35,37,38,45,	123,134,141,144	4,10,15,18,19,20,	4,10,15,18,19,20,21,	4,10,15,18,19,20,21,	36,75,76,122	4,10,15,18,19,20,21,
81,86,121,124,	16,17,59,60,92,	21,31,40,46,57,58,	54,56,73,77,		21,31,40,46,57,58,	31,40,46,57,58, 63,	31,40,46,57,58,63,	135,138,142	31,40,46,57,58,63,
128,131,137	94,96, 97,98,130	66,67,68,70,71,72,	101,102,105,107,		63,67,68,70,71,72,	66,68,70,71,72,	66,67,70,71,72,		66,67,68,71,72,
		78,79,82,93,100,	116,117,129,139		78,79,82,93,100,	78,79,82,93,100,	78,79,82,93,100,		78,79,82,93,100,
		103,104,112,127,			103,104,112,127,	103,104,112,127,	103,104,112,127,		103,104,112,127,
		133,136,140,143			133,136,140,143	133,136,140,143	133,136,140,143		133,136,140,143
South Africa	South Africa	South Africa	South Africa	South Africa	South Africa	South Africa	South Africa	South Africa	South Africa
Berea (71)	Berea (72)	Cape Town (73)	Cape Town (74)	Cape Town (75)	Cape Town (76)	Eastern Cape (77)	eThekwini (78)	eThekwini (79)	eThekwini (80)

4,10,15,18,19,20,21, 31,40,46,57,58,63, 66,67,68,70,72, 78,79,82,93,100, 103,104,112,127, 133,136,140,143	4,10,15,18,19,20,21, 31,40,46,57,58,63, 66,67,68,70,71, 78,79,82,93,100, 103,104,112,127, 133,136,140,143	7,35,37,38,45, 54,56,64,77, 101,102,105,107, 116,117,129,139	43,44,47,49, 80,99,114,126	36,69,76, 135,138,142	36,69,75,122, 135,138,142	7,35,37,38,45, 54,56,64,73, 101,102,105,107, 116,117,129,139	4,10,15,18,19,20,21, 31,40,46,57,58,63, 66,67,68,70,71,72, 79,82,93,100, 103,104,112,127, 133,136,140,143	4,10,15,18,19,20,21, 31,40,46,57,58,63, 66,67,68,70,71,72, 78,82,93,100, 103,104,112,127, 133,136,140,143	43,44,47,49,74, 99,114,126
South Africa eThekweni (81)	South Africa Free State (82)	South Africa Free State (83)	South Africa Free State (84)	South Africa Free State (85)	South Africa Harry Gwala (86)	South Africa Harry Gwala (87)	South Africa Harry Gwala (88)	South Africa Harry Gwala (89)	South Africa Ilembe (90)
1,8,22,48,50,51,55, 61,86,121,124, 128,131,137	4,10,15,18,19,20,21, 31,40,46,57,58,63, 66,67,68,70,71,72, 78,79,93,100, 103,104,112,127, 133,136,140,143	84,89,90,91,106,109	83,89,90,91,92, 106,109	88,110,111,115,119, 125	1,8,22,48,50,51,55, 61,81,121,124, 128,131,137	95,118,120,132	85,110,115,119, 125	83,84,90,91,106,109	83,84,89,91,92, 106,109
South Africa Ilembe (91)	South Africa Ilembe (92)	South Africa Ilembe (93)	South Africa Kg Cetshwayo (94)	South Africa Kg Cetshwayo (95)	South Africa Kg Cetshwayo (96)	South Africa Kg Cetshwayo (97)	South Africa KZN (98)	South Africa KZN (99)	South Africa KZN (100)
83,84,89,90,106,109	6,9,11,12,13,14, 16,17,59,60,62, 94,96,97,98,130	4,10,15,18,19,20,21, 31,40,46,57,58,63, 66,67,68,70,71,72, 78,79,82,100, 103,104,112,127, 133,136,140,143	6,9,11,12,13,14, 16,17,59,60,62,92, 96,97,98,130	87,118,120,132	6,9,11,12,13,14, 16,17,59,60,62,92, 94,97,98,130	6,9,11,12,13,14, 16,17,59,60,62,92, 94,96,98,130	6,9,11,12,13,14, 16,17,59,60,62,92, 94,96,97,130	43,44,47,49,74, 80,114,126	4,10,15,18,19,20,21, 31,40,46,57,58,63, 66,67,68,70,71, 72,78,79,82,93 103,104,112,127, 133,136,140,143
South Africa KZN (101)	South Africa LP (102)	South Africa Mmorningside (103)	South Africa Mmorningside (104)	South Africa North West (105)	South Africa North West (106)	South Africa North West (107)	South Africa North West (108)	South Africa Overport (109)	South Africa Sisonke (110)

7,35,37,38,45,	7,35,37,38,45,	4,10,15,18,19,20,21,	4,10,15,18,19,20,21,	7,35,37,38,45,	83,84,89,90,91,	7,35,37,38,45,	29,30,32,42,113	83,84,89,90,91,92,	85,88,111,115,119,
54,56,64,73,77,	54,56,64,73,77,	31,40,46,57,58,63,	31,40,46,57,58,63,	54,56,64,73,77,	106,109	54,56,64,73,77,		106	125
102,105,107,	101,105,107,	66,67,68,70,71,	66,67,68,70,71,	101,102,107,		101,102,105,			
116,117,129,139	116,117,129,139	72,78,79,82,93,	72,78,79,82,93,	116,117,129,139		116,117,129,139			
		100,104,112,127,	100,103,112,127,						
		133,136,140,143	133,136,140,143						
South Africa	South Africa	South Africa	South Africa	South Africa	South Africa	South Africa	South Africa	South Africa	South Africa
Sisonke (111)	Stanger (112)	Ugu (113)	Ugu (114)	Ugu (115)	Ugu (116)	Umbilo (117)	Umbilo (118)	Umbilo (119)	Umbilo (120)
85,88,110,115,119,	4,10,15,18,19,20,21,	29,30,32,42,108	43,44,47,49,74,	85,88,110,119, 125	7,35,37,38,45,	7,35,37,38,45,	87,95,120,132	85,88,110,115,125	87,95,118,132
125	31,40,46,57,58,63,		80,99,126		54,56,64,73,77,	54,56,64,73,77,			
	66,67,68,70,71,				101,102,105,107,	101,102,105,107,			
	72,78,79,82,93,				117,129,139	116,129,139			
	100,103,104,127,								
	133,136,140,143								
South Africa	South Africa	South Africa	South Africa	South Africa	South Africa	South Africa	South Africa	South Africa	South Africa
Umgungundl (121)	Umgungundl (122)	Umgungundl (123)	Umgungundl (124)	Umkhanyak (125)	Umkhanyak (126)	Umkhanyak (127)	Umkhanyak (128)	Umzinyathi (129)	Umzinyathi (130)
1,8,22,48,50,51,55,	36,69,75,76,	65,134,141,144	1,8,22,48,50,51,55,	85,88,110,115,119	43,44,47,49,74,	4,10,15,18,19,20,21,	1,8,22,48,50,51,55,	7,35,37,38,45,	6,9,11,12,13,14,
61,81,86,124,	135,138,142		61,81,86,121,		80,99,114	31,40,46,57,58,63,	61,81,86,121,124,	54,56,64,73,77,	16,17,59,60,62,92,
128,131,137			128,131,137			66,67,68,70,71,72,	131,137	101,102,105,107,	94,96, 97,98
						78,79,82,93,100,		116,117,139	
						103,104,112,127,			
						133,136,140,143			
South Africa	South Africa	South Africa	South Africa	South Africa	South Africa	South Africa	South Africa	South Africa	South Africa
Umzinyathi (131)	Umzinyathi (132)	Uthukela (133)	Uthukela (134)	Uthukela (135)	Uthukela (136)	Uthungulu (137)	Uthungulu (138)	Uthungulu (139)	Uthungulu (140)

1,8,22,48,50,51,55,	87,95,118,120	4,10,15,18,19,20,21,	65,123,141,144	36,69,75,76,	4,10,15,18,19,20,21,	1,8,22,48,50,51,55,	36,69,75,76,	7,35,37,38,45,	4,10,15,18,19,20,21,
61,81,86,121,124,		31,40,46,57,58,63,		122,138,142	31,40,46,57,58,63,	61,81,86,121,124,	122,135,142	54,56,64,73,77,	31,40,46,57,58,63,
128,137		66,67,68,70,71,72,			66,67,68,70,71,72,	128,131		101,102,105,107,	66,67,68,70,71,72,
		78,79,82,93,100,			78,79,82,93,100,			116,117,129	78,79,82,93,100,
		103,104,112,127,			103,104,112,127,				103,104,112,127,
		136,143			133,143				133,136,143
South Africa	South Africa	South Africa	South Africa	Tunisia					
Zululand (141)	Zululand (142)	Zululand (143)	Zululand (144)	Ben Arous (145)					
65,123,134,144	36,69,75,76,	4,10,15,18,19,20,21,	65,123,134,141	2,3,5,23,24,25,26,					
	122,135,138	31,40,46,57,58,63,		27,28,33,34,39,					
		66,67,68,70,71,72,		41,52,53					
		78,79,82,93,100,							
		103,104,112,127,							
		133,136,140							

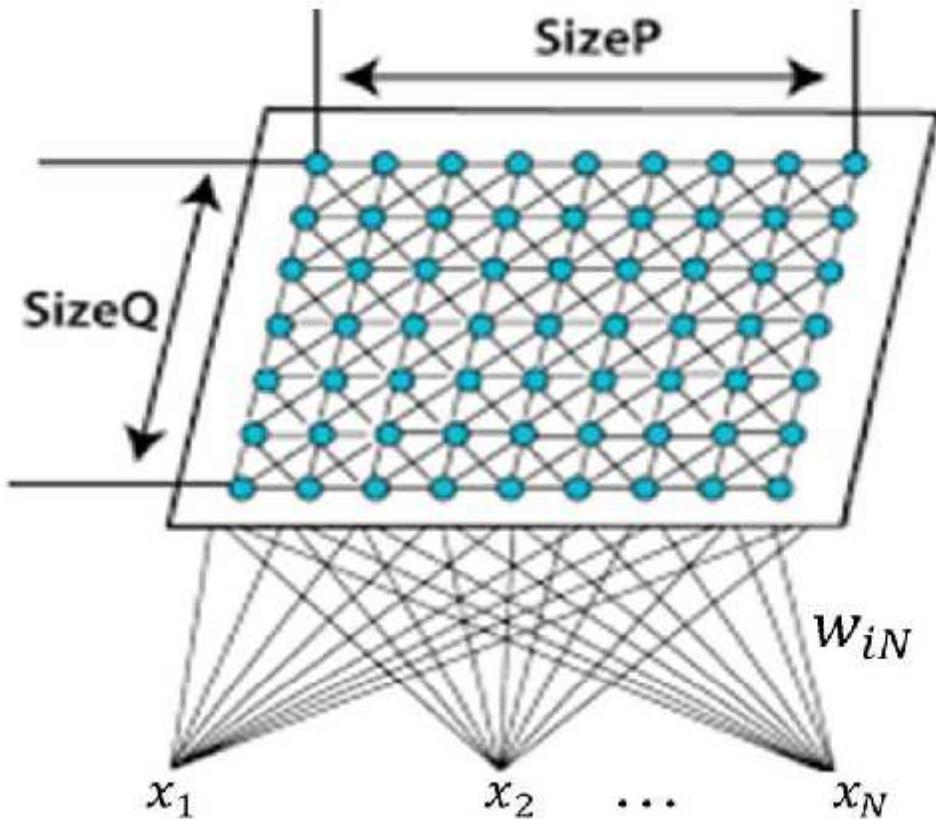
# Figures



**Figure 1**

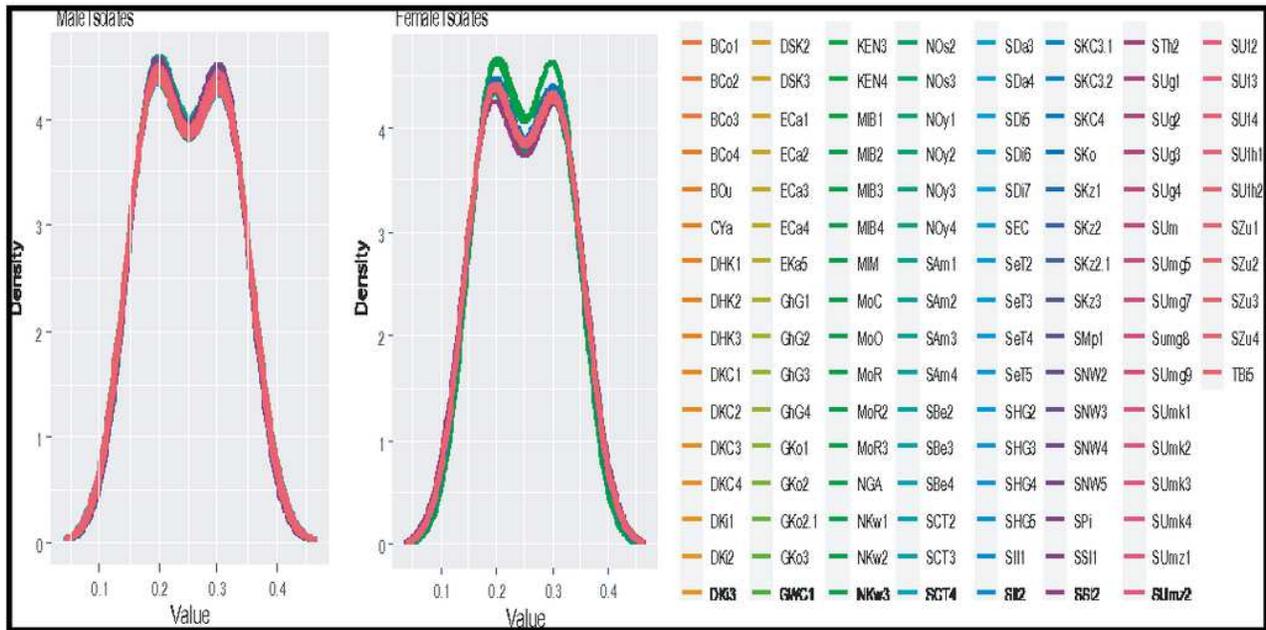
Workflow describing the proposed system framework. The workflow begins with the excavation of FASTA files of human SARS-CoV-2 genome sequences from GISAID. These files are stripped and processed into a genome database (DB) as multiple columns of nucleotide sequence. A series of AI/ML techniques are

applied to extract knowledge from the genome datasets as follows: Compute dis(similarities) scores between the various pairs of genome sequences and obtain a genomic tree of highly dis(similar) isolates grouped in the form of a dendrogram/phylogenomic tree. Determine the optimal number of natural clusters—to provide additional knowledge. Separate the viral sub-strains using self-organizing map (SOM) component planes—for transmission pathways visualization. Perform direct pairwise nucleotide alignment of the entire genome sequences—to yield a nucleotide similarity matrix. Generate cognitive map—for intelligent sub-strains contact tracing and prediction.



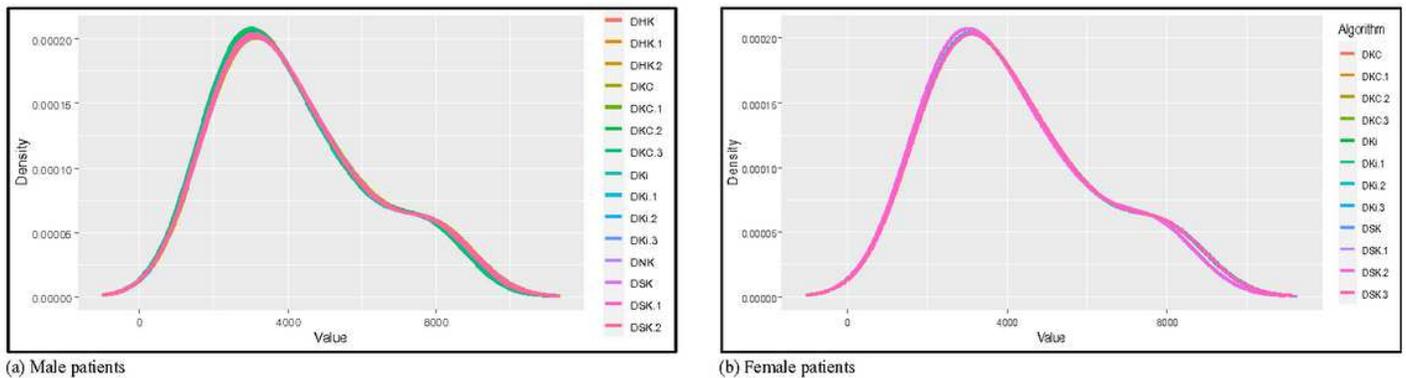
**Figure 2**

SOM showing the map topology and interactions between nodes. Each neuron is assigned a vector of weights ( $w=w_{i1},w_{i2},..w_{iN}$ ) with dimension similar to the input vector  $i$  ( $i=1,2,..,L$ ); where  $L$  is the total number of neurons in the network. The input nodes have  $p$  features, and the output nodes,  $q$  prototypes, with each prototype connected to all features. The weight vector of the connections consumes the prototype of each neuron and has same dimension as the input vector. SOMs differ from other artificial neural networks as they apply competitive learning, against error correction learning such as backpropagation, and the fact that they preserve the topological properties of the input space using a neighborhood function.



**Figure 3**

Density plots of entire genome datasets. A Density plot visualizes the distribution of data over a continuous scale. It is a variation of the histogram that applies kernel smoothing to plot values, enabling smoother distributions by smoothening out the noise. The peaks of the density plot help display where values are concentrated over the interval.

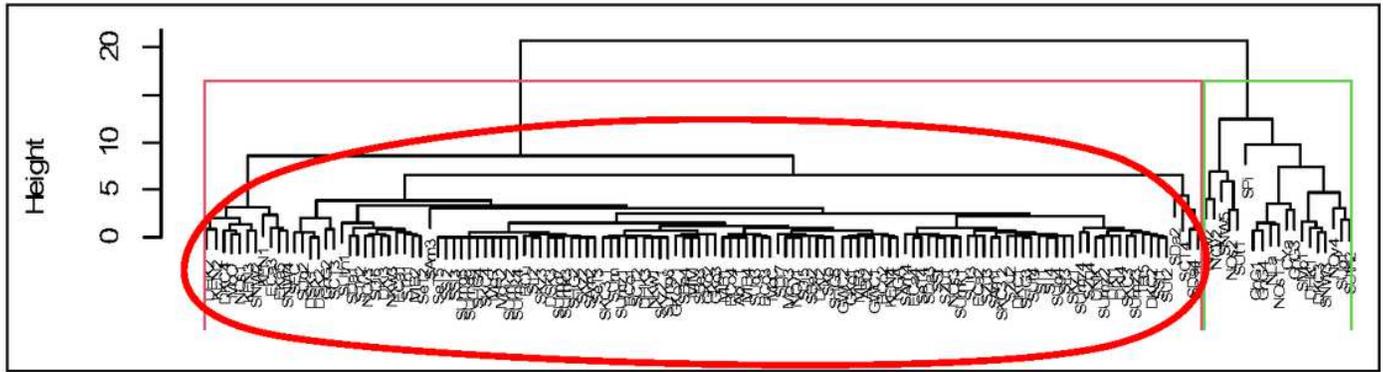


**Figure 4**

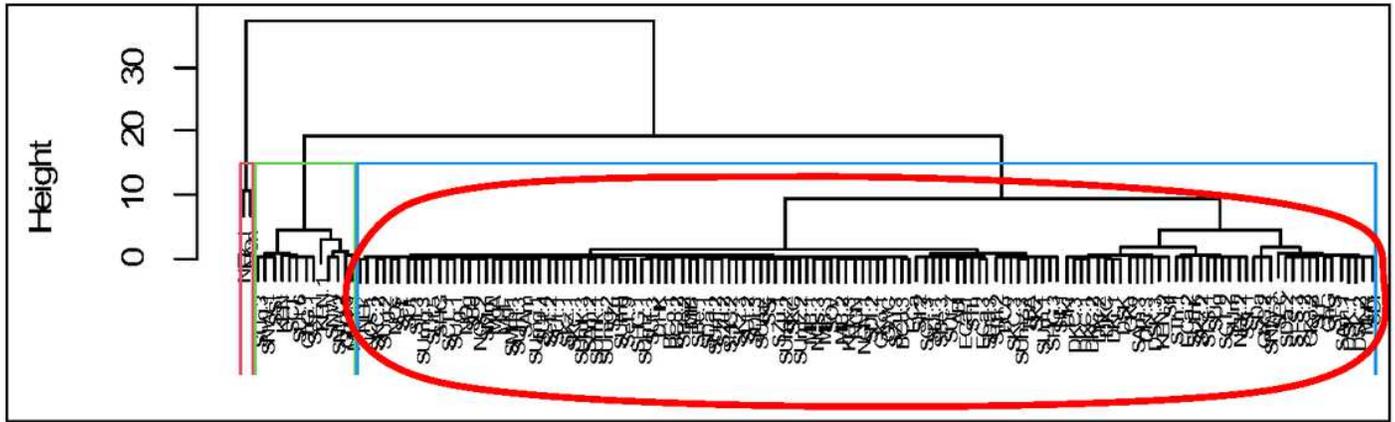
Density plots of DRC's isolates. Both plots exhibit similar curve pattern with near aligned isolates, indicating identical genome expression between the isolates.



Density plots of South African isolates. Both plots maintain same patterns with shades of similar isolates dominating the plot area.



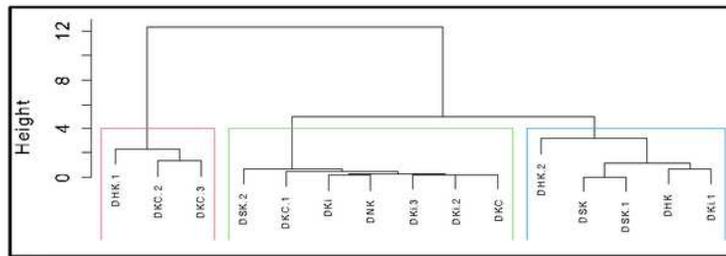
(a) Male patients



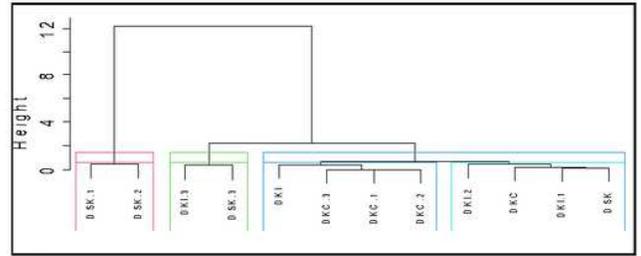
(b) Female patients

**Figure 8**

Phylogenomic trees for African male and female isolates. For full names of country codes, see Additional file 1: SupplData1\_1.xlsx.



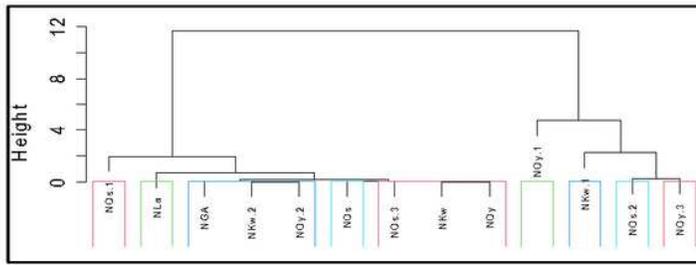
(a) Male patients



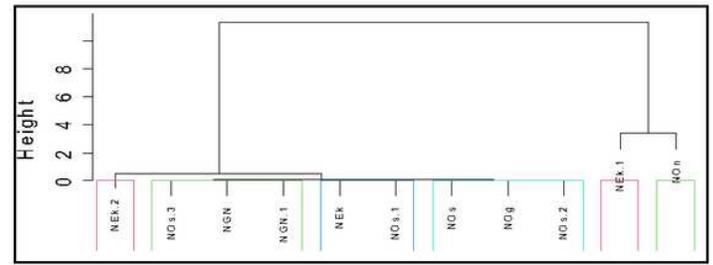
(b) Female patients

**Figure 9**

Phylogenomic trees for DRC's male and female isolates. For full names of country codes, see Additional file 1: SupplData1\_1.xlsx.



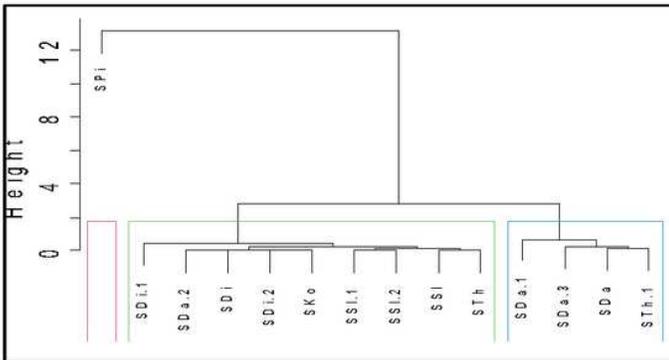
(a) Male patients



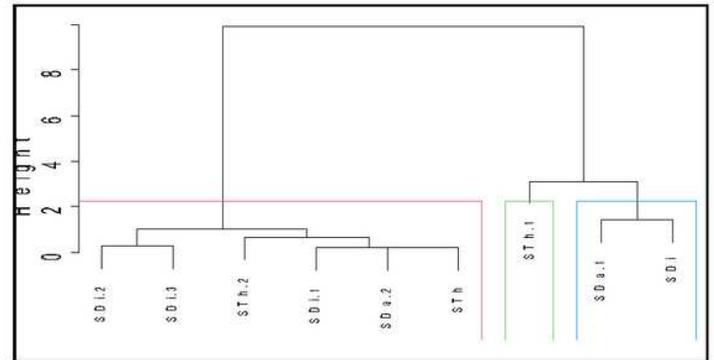
(b) Female patients

**Figure 10**

Phylogenomic trees for Nigerian male and female isolates. For full names of country codes, see Additional file 1: SupplData1\_1.xlsx.



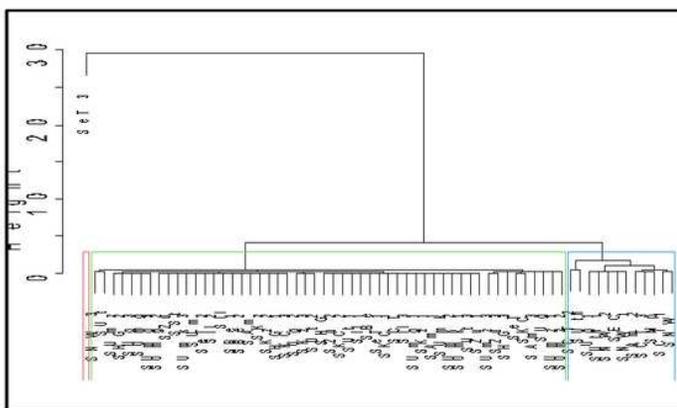
(a) Male patients



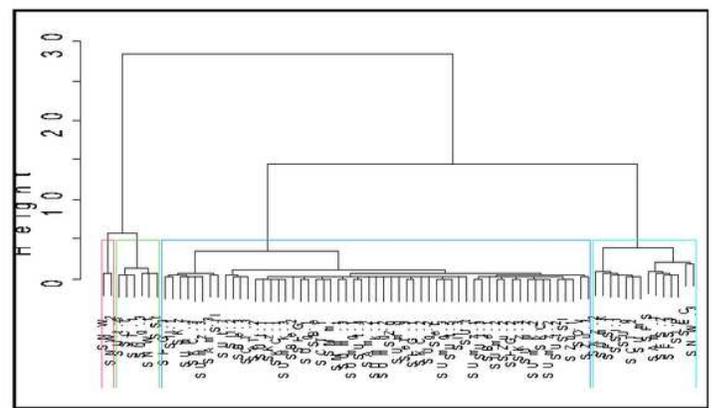
(b) Female patients

**Figure 11**

Phylogenomic trees for Senegalese male and female isolates. For full names of country codes, see Additional file 1: SupplData1\_1.xlsx.



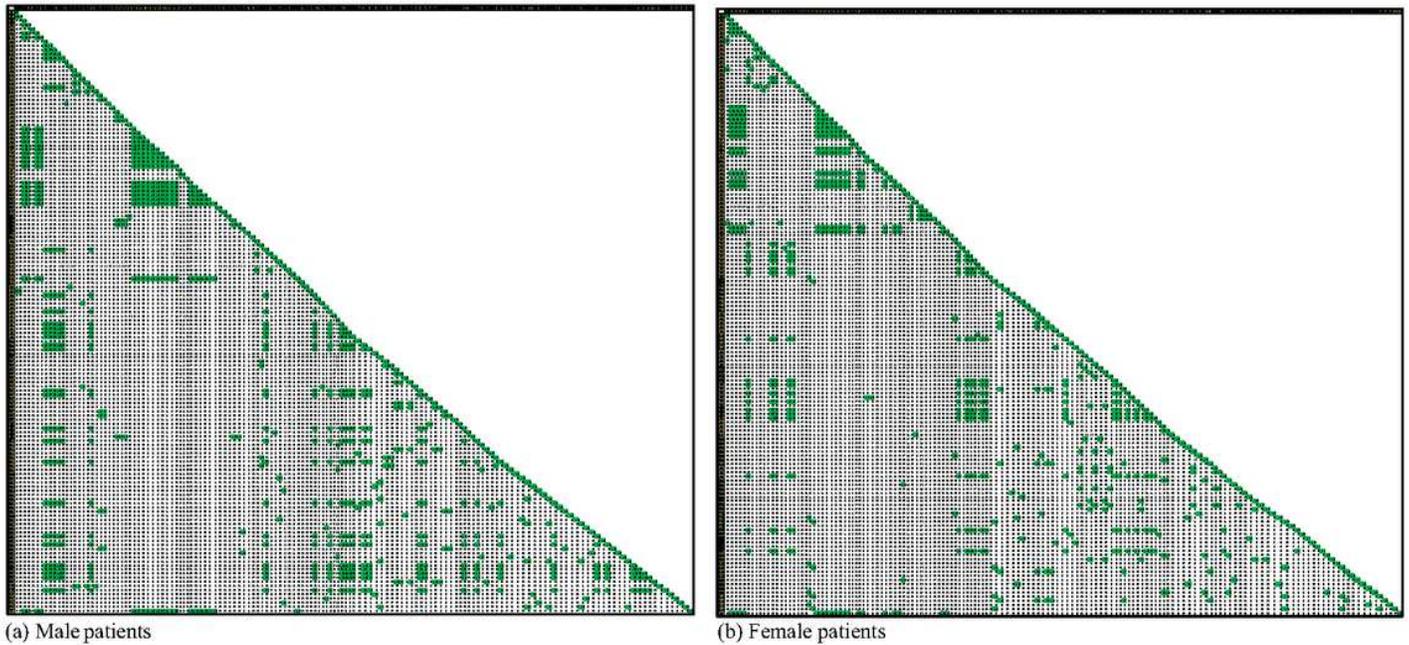
(a) Male patients



(b) Female patients

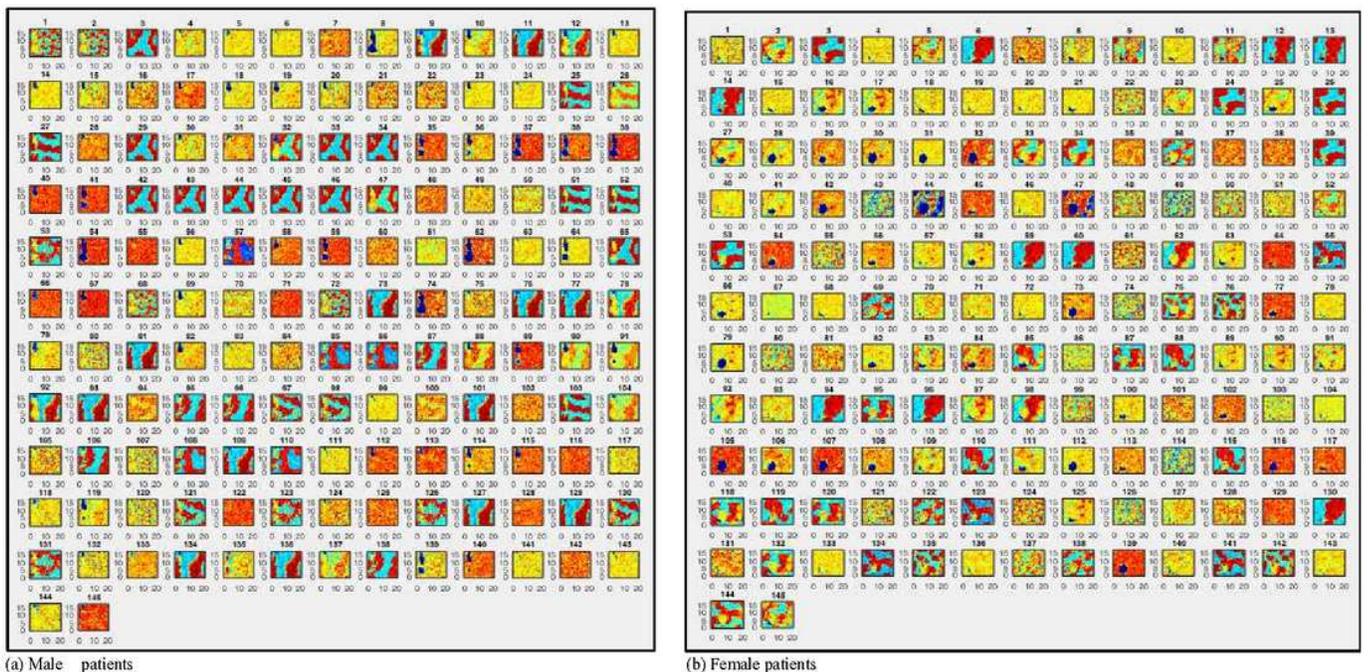
**Figure 12**

Phylogenomic trees for South African male and female isolates. For full names of country codes, see Additional file 1: SupplData1\_1.xlsx.



**Figure 13**

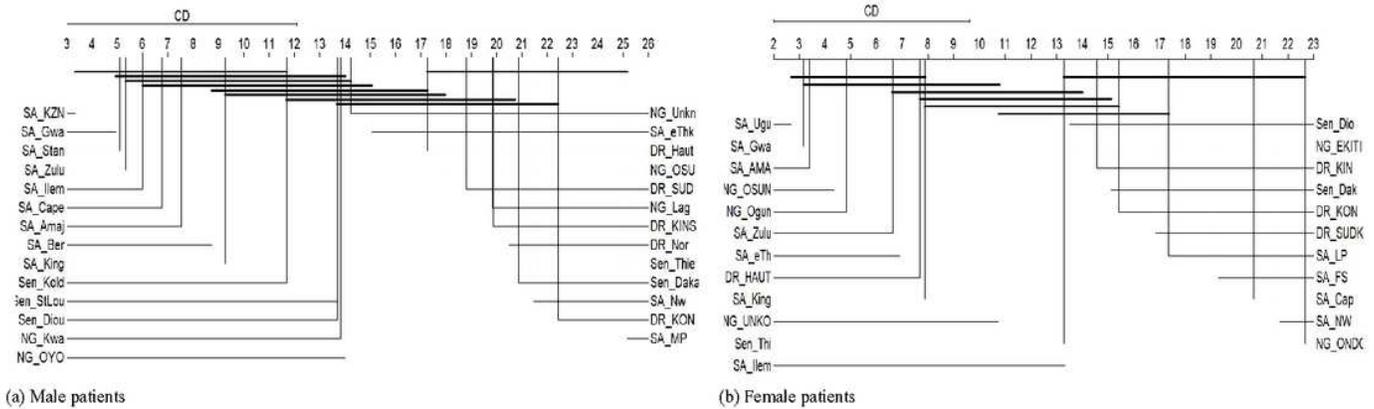
Nucleotide similarity matrices. Green colored cells are regions of high similarity that may indicate functional, structural and/or evolutionary relationships between nucleotide sequences.



**Figure 14**

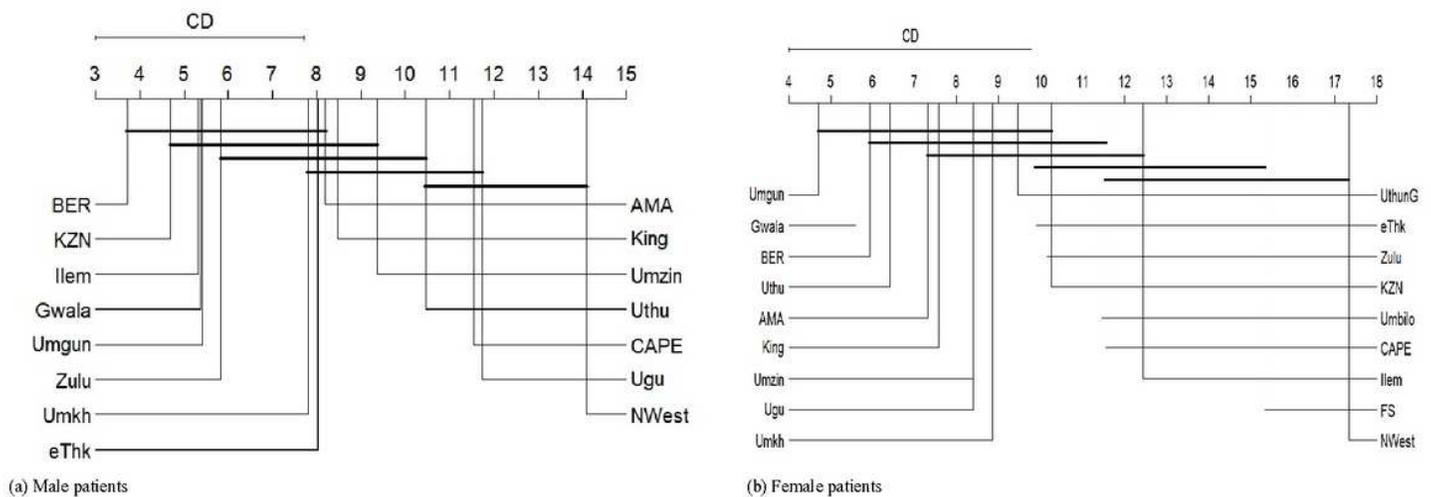
SOM component planes visualization. Maps are ordered by countries, with at least 1 isolate per country. The isolate numbers (1-145) represent the various states of the country excavated from GSAID (see SupplData1\_1.xlsx). The isolates are distributed by countries as follows. Male: Algeria (1-2), Benin (3-7), Cameroon (8), DRC (9-23), Egypt (24-28), Gambia (29-34), Ghana (35-38), Kenya (39-42), Mali (43-47),

Morocco (48-52), Nigeria (53-65), Senegal (66-79), South Africa (80-144), Tunisia (145). Female: Algeria (1), Benin (2-5), DRC (6-18), Egypt (19-22), Gambia (23-26), Ghana (27-30), Kenya (31-34), Madagascar (35-36), Mali (37-40), Morocco (41-42), Nigeria (43-53), Senegal (54-64), South Africa (65-144), Tunisia (145).



**Figure 15**

Inter-country CD plots for male and female patients. For a significance level  $\alpha$  the Nemenyi's test determines the critical difference (CD), if the difference between the average ranking of two treatments is greater than CD, then the null hypothesis that the treatments have the same performance is rejected.



**Figure 16**

Intra-country CD plots for male and female patients. For a significance level  $\alpha$  the Nemenyi's test determines the critical difference (CD), if the difference between the average ranking of two treatments is greater than CD, then the null hypothesis that the treatments have the same performance is rejected.

## Supplementary Files

This is a list of supplementary files associated with this preprint. [Click to download.](#)

- SupplData11.xlsx
- SupplData21.xlsx