

1 **A SARS-CoV-2 lineage A variant (A.23.1) with altered spike has emerged and is**
2 **dominating the current Uganda epidemic**

3 Daniel Lule Bugembe^{1*}, My V.T.Phan^{1*}, Isaac Ssewanyana², Patrick Semanda², Hellen
4 Nansumba², Beatrice Dhaala¹, Susan Nabadda², Áine Niamh O'Toole³, Andrew Rambaut³,
5 Pontiano Kaleebu^{1,4}, Matthew Cotten^{1,5}

6
7 **Affiliations**

- 8 1. MRC/UVRI & LSHTM Uganda Research Unit, Entebbe, Uganda
9 2. Central Public Health Laboratories of the Republic of Uganda, Kampala, Uganda
10 3. Institute for Evolutionary Biology, University of Edinburgh
11 4. Uganda Virus Research Institute, Entebbe, Uganda
12 5. MRC-University of Glasgow Centre for Virus Research, Glasgow, UK

13 *These authors contribute equally.

14
15 Correspondence: Matthew Cotten (Matthew.Cotten@lshtm.ac.uk)

16
17 **Keywords:** SARS-CoV-2, new variant, A.23; A.23.1, spike protein changes

18
19 **Introductory paragraph**

20
21 SARS-CoV-2 genomic surveillance in Uganda provides an opportunity to provide a focused
22 description of the virus evolution in a small landlocked East African country. Here we show a
23 recent shift in the local epidemic with a newly emerging lineage A.23 evolving into A.23.1
24 which is now dominating the Uganda cases and has spread to 26 other countries. Although
25 the precise changes in A.23.1 as it has adapted are different from the changes in the
26 variants of concern (VOC), the evolution shows convergence on a similar set of proteins.
27 The A.23.1 spike protein coding region has accumulated changes that resemble many of the
28 changes seen in VOC including a change at position 613, a change in the furin cleavage site
29 that extends the basic amino acid motif, and multiple changes in the immunogenic N-

30 **NOTE:** This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.

31 other VOC show (nsp6, ORF8 and ORF9). The clinical impact of the A.23.1 variant is not yet
32 clear, however it is essential to continue careful monitoring of this variant, as well as rapid
33 assessment of the consequences of the spike protein changes for vaccine efficacy.
34

35 **Main Text (2324 words)**

36 The novel Severe Acute Respiratory Syndrome Coronavirus 2 (SARS-CoV-2)(1) and
37 the associated disease Coronavirus Disease 2019 (COVID-19)(2)(3) continue to spread
38 throughout the world, causing >120 million infections and >2.6 million deaths (16 Mar 2021,
39 [Johns Hopkins COVID-19 Dashboard](#)). Genomic surveillance has played a key role in the
40 response to the pandemic; sequence data from SARS-CoV-2 provides information on the
41 transmission patterns and the evolution of the virus as it enters new regions and spreads. As
42 COVID-19 vaccines become available and are implemented, monitoring SARS-CoV-2
43 genetic changes, especially changes at epitopes with implications for immune escape is
44 crucial. A detailed classification system has been defined to help monitor SARS-CoV-2 as it
45 evolves (4) with virus sequences classified into 2 main phylogenetic lineages (Pango
46 lineages) A and B, representing the earliest divergence of SARS-CoV-2 in the pandemic and
47 then into sub-lineages within these. Several Variants of Concern (VOC) have emerged
48 showing increased transmission patterns and reduced susceptibility to vaccine and/or
49 therapeutic antibody treatments. These VOC include lineage B.1.1.7 first identified in the UK
50 (5), B.1.351 in South Africa (6) and lineage P.1 (B.1.1.28.1) in Brazil (7).

51

52 **Status of the SARS-CoV-2 epidemic in Uganda**

53 SARS-CoV-2 infection was first detected in Uganda in March 2020, initially among
54 international travellers until passenger flights were stopped in late March 2020. A second
55 route of virus entry with truck drivers from adjacent countries then became apparent (8).
56 Since August 2020, community transmission dominated the Uganda case numbers. By
57 March 2021 total cases in Uganda were 40,535, with 334 deaths attributed to the virus. We

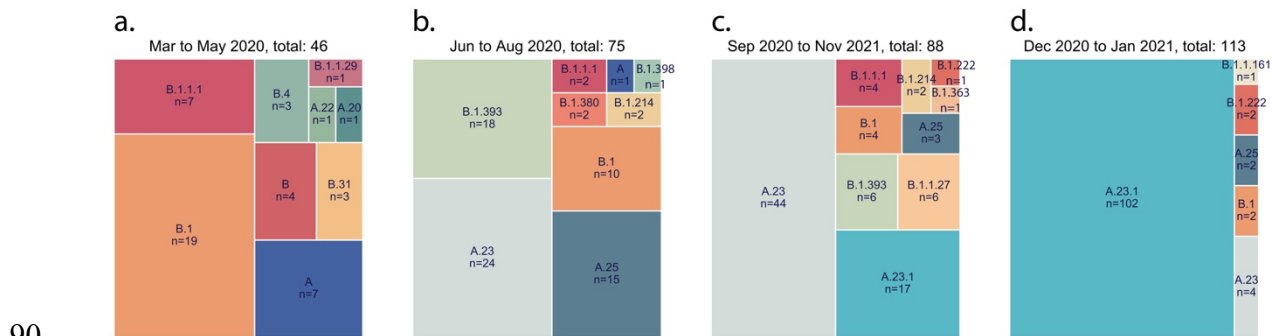
58 have continued our efforts to generate SARS-CoV-2 genomic sequence data to monitor
59 virus movement and genetic changes and we report here on a novel sub-lineage A (A.23.1)
60 that emerged and is dominating the local epidemic. The A.23.1 variant encodes multiple
61 changes in the spike protein as well as in nsp6, ORF8 and ORF9, some predicted to be
62 functionally similar to those observed in VOC in lineage B.

63

64 **Changes in prevalence of lineage A viruses**

65 The genomes generated here were classified into Pango lineages(4) using the
66 Pangolin module pangoleARN (<https://github.com/cov-lineages/pangolin>) and into
67 NextStrain clades using NextClade (9) (<https://clades.nextstrain.org/>). The distribution of
68 virus lineages circulating in Uganda changed dramatically over the course of the year. A
69 clear feature of the earlier COVID-19 epidemic in the country was the diversity of viruses
70 found throughout the country attributed to frequent flights into Uganda from Europe, UK, US
71 and Asia; this is reflected in the 9 lineages seen from March to May 2020 with a mixture of
72 both lineage A and B viruses (Figure 1, panel a). After passenger flights were limited in
73 March 2020, the virus entered via land travel with truck drivers. Uganda is landlocked
74 country, characterised by its important geographical position, i.e. the crossing of two main
75 routes of the Trans-Africa Highway in East Africa. The essential nature of produce and
76 goods transport allowed virus movement from/to Kenya, South Sudan, DRC, Rwanda and
77 Tanzania. In the period of June to August 2020 lineage B.1 and B.1.393 strains were
78 abundant, similar to patterns observed in Kenya (10) (Figure 1, panel b) although lineage A
79 viruses did not decline as seen in US and Europe. Lineage A.23 strains were first observed
80 in two prison outbreaks in Amuru and Kitgum, Uganda in August 2020 and by the
81 September-November period, the A.23 was the major lineage circulating throughout the
82 country (Figure 1, panel c). The A.23 virus continued to evolve into the lineage A.23.1, first
83 observed in late October 2020. Given the diversity of virus lineages found in the country

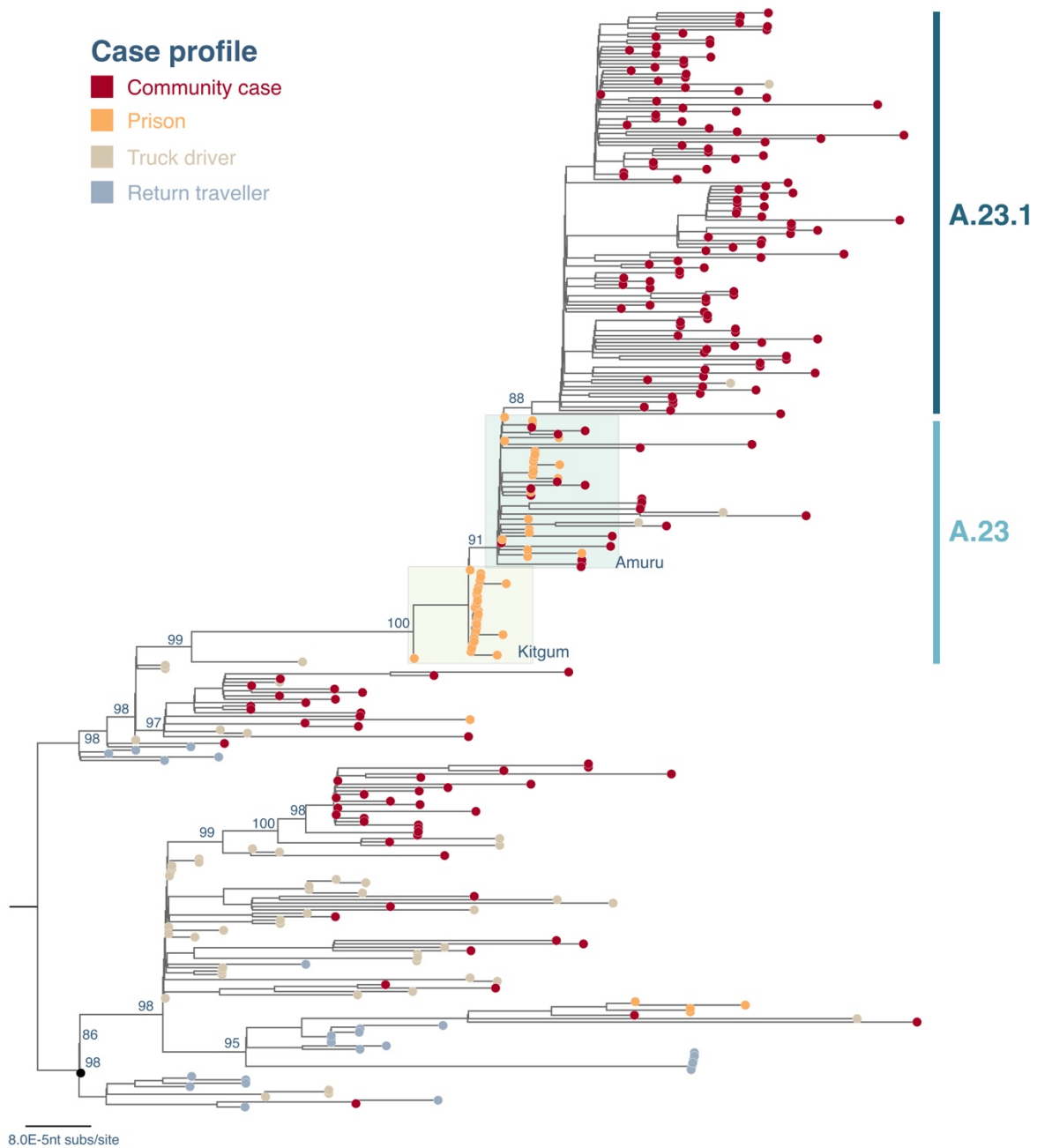
84 from March until November 2020, it was unexpected that by late December 2020 to January
 85 2021, lineage A.23.1 viruses represented 90% (102 of 113 genomes) of all viruses
 86 observed in Uganda (Figure 1, panel d). In all time periods, the SARS-CoV-2 positive
 87 sample were obtained from multiple clinical and surveillance locations throughout Uganda
 88 indicating that the differences are unlikely to be due to sampling different subpopulations in
 89 the country at different times.



91 **Figure 1 SARS-CoV-2 lineage diversity in Uganda.** All high-coverage complete from
 92 Uganda (n=322) were lineage typed using the pangolin resource ([https://github.com/cov-](https://github.com/cov-lineages/pangolin)
 93 [lineages/pangolin](https://github.com/cov-lineages/pangolin)) Lineage counts were stratified into four periods (**panel a:** March-May
 94 2020, **panel b:** June-August 2020, **panel c:** September to November 2020 and **panel d:**
 95 December 2020 to January 2021). The percentage of each lineage within each set was
 96 plotted as a treemap using squarify (<https://github.com/laserson/squarify>) with the size of
 97 each sector proportional to the number of genomes, genomes numbers are listed with "n=".
 98

99 **Virus sequence diversity including fatal cases**

100 All newly and previously generated Uganda genomes that were complete and high-
101 coverage (n=322) were used to construct a maximum-likelihood phylogenetic tree (Figure 2).



102

103 **Figure 2. Maximum-likelihood phylogenetic tree comparing all available complete and**
104 **high-coverage Uganda sequences (N=322).** Strain names are coloured according to the
105 case profile (cases from the community: dark red, prison: orange, truck driver: light brown,
106 return traveller: light blue). The case clusters from prisons in Kitgum and Amuru are
107 highlighted in colour boxes in light yellow and light green, respectively. The lineages A.23
108 and A.23.1 are indicated. The tree was rooted where lineages A and B were split. The
109 branch length is drawn to the scale of number of nucleotide substitutions per site, indicated
110 in lower left, and only bootstrap values of major nodes were shown.

111 A number of A and B variant lineages were observed briefly at low frequencies and
112 may have undergone extinction, similar to patterns observed in the UK (11) and Scotland
113 (12). Genomes identified from a truck driver are often observed basal to community clusters
114 (Figure 2), suggesting the importance of this route in the introduction and spread of the virus
115 into Uganda. Most of genomes from truck drivers sampled at points of entry (POEs)
116 bordering Kenya belonged to lineage B.1 and B.1.393 consistent with the pattern reported in
117 Kenya (10). However, genomes identified from truck drivers from Tanzania, and from the
118 Elegu POE bordering South Sudan, albeit small numbers, belonged to both A and B.1
119 lineages. Continued monitoring of truck drivers coming in and out of the Uganda provides a
120 useful description of the inland circulation of strains in this part of world, where genomic
121 surveillance is not as detailed as in other parts of the world.

122

123 **Emergence of A.23 and A.23.1**

124 Outbreaks of SARS-CoV-2 infections were reported in the Amuru and Kitgum prisons in
125 August 2020 (13)(14). The SARS-CoV-2 genome sequences from individuals in the prisons
126 were exclusively belonging to lineage A (Figure 2) with three amino acid (aa) changes
127 encoded in the spike protein (F157L, V367F and Q613H, Figure 3) that now define lineage
128 A.23 (see below). By October 2020, lineage A.23 viruses were also found outside of the
129 prisons in a community sample from Lira (a town 140 km from Amuru), in two samples from
130 the Kitgum hospital, in several community samples from Kampala, Jinja, Mulago, Tororo,
131 Soroti as well as in 2 truck drivers collected at POE bordering Kenya. By November 2020,
132 the A.23 viruses spread further to northern Uganda in Gulu and Adjumani. Lineage A.23
133 viruses were not seen in Uganda (or anywhere in the world) before August 2020 (Figure 3
134 panel c), yet the A.23 viruses were attributed to 32% of the viruses in Uganda (Figure 1)
135 from June to August 2020 and 50% of the observed viruses in September to November
136 2020. In late October, the A.23.1, a variant evolving from A.23, with additional change in the

137 spike protein (P681R) was observed (Figure 3, panel b, c) and by December 2020-January
138 2021, 90% of identified genomes (102 out of 113) belonged to the new A.23.1 lineage
139 (Figure 1 and 2). The mutations in A.23.1 were consistent with evolution from an original
140 A.23 virus observed in Amuru/Kitgum cluster (Figure 2 and Supplemental Figure 1) as well
141 as changes nsp6 and ORF9 (Supplemental Figure 2 and 4).

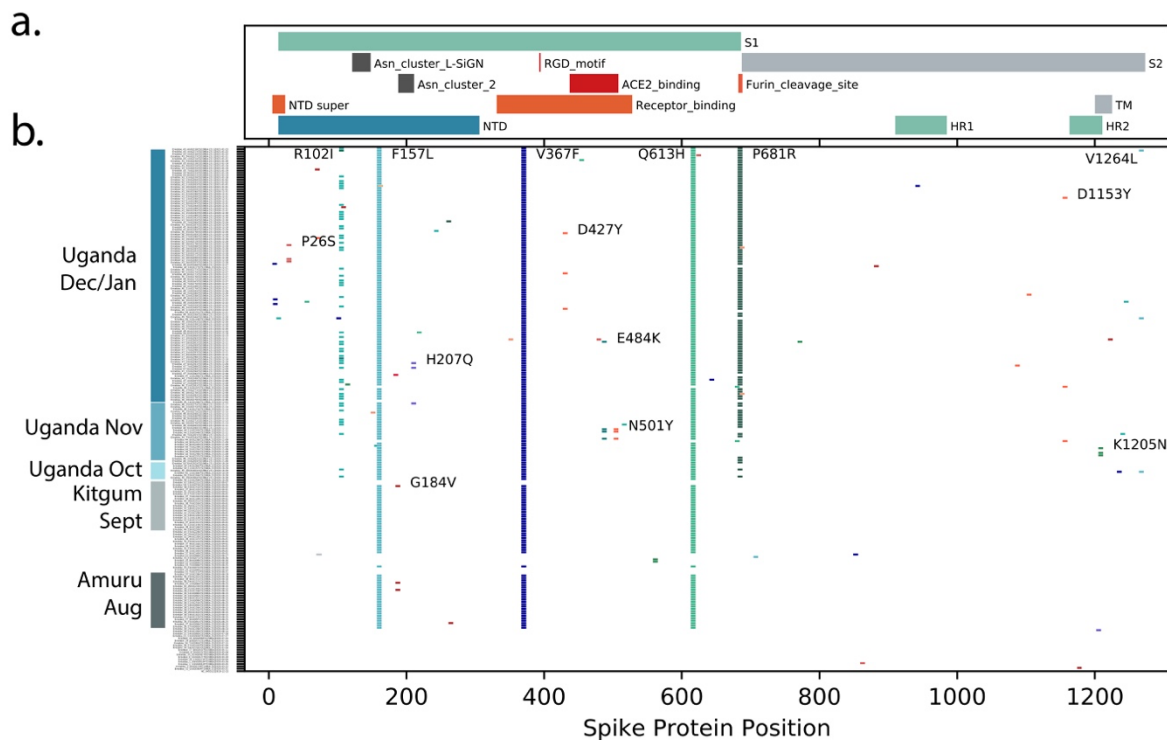
142

143 **Important changes observed in the spike protein**

144 The spike protein is crucial for virus entry into host cells, for tropism, and is a critical
145 component of COVID-19 vaccine development and monitoring. The changes in spike protein
146 observed in Uganda and global A.23 and A.23.1 viruses are shown in Figure 3 panel b.
147 Many amino acid (aa) changes were single events with no apparent transmission observed.
148 However, the initial lineage A.23 genomes from Amuru and Kitgum encoded three amino
149 acid changes in the exposed S1 domain of spike (F157L, V367F and Q613H, Figure 3 panel
150 b). The V367F change is reported to modestly increase infectivity(15), the Q613H change
151 may have similar consequences as the D614G change observed in the B.1 lineage found
152 predominantly in Europe and USA; in particular, D614G was reported to increase infectivity,
153 spike trimer stability and furin cleavage (15),(16),(17),(18). These changes were not
154 observed in previously reported genomes from Uganda (8). Of some concern, the mutations
155 E484K and N501Y amino acid changes in the receptor binding domain (RBD) were
156 observed in the A.23 viruses identified in Adjumani cases on 9th to 11th November 2020
157 (Figure 3, panel b). These two amino acid changes are shown to substantially compromise
158 the vaccine efficacy as well as antibody treatments.

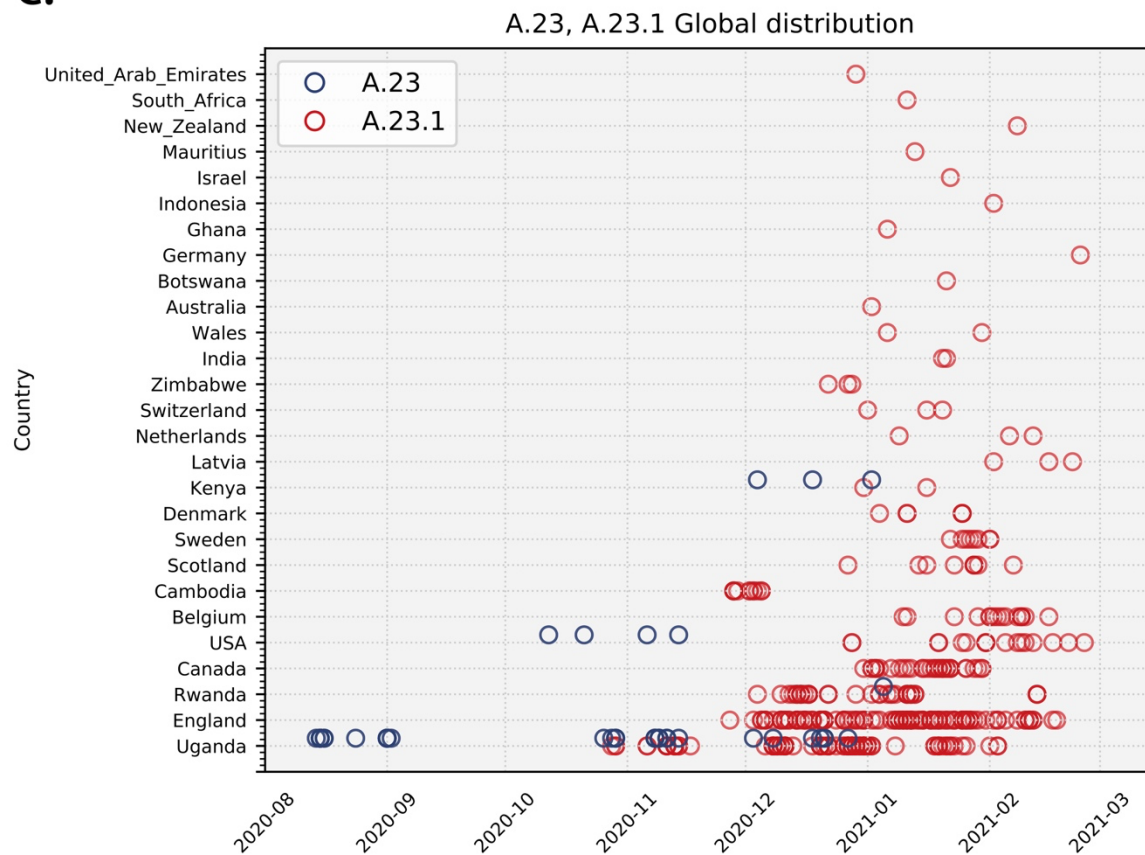
159

160



161

C.



162

163 **Figure 3. Spike protein changes in lineage A.23 and A.23.1 relative to the SARS-CoV-2**
 164 **reference strain (NC_045512) encoded protein are documented. Panel a: The locations**

165 of important spike protein features are indicated. NTD: N-terminal domain, RBD: receptor-
166 binding domain, S1: spike 1, S1: Spike 2, TM: transmembrane domain, HR1: helical repeat
167 1, HR2: helical repeat 2, NTD super: N-terminal domain supersite. **Panel b:** Each line
168 represents the encoded spike protein sequence from a single genome, ordered by date of
169 samples collection (bottom earliest, top most recent). Sequences from Amuru in August
170 2020, Kitgum in September 2020 and Uganda in October, November, December
171 2020/January 2021 are indicated. Markers indicating the positions of amino acid (aa)
172 differences from the reference strain, changes observed in multiple genomes are annotated
173 with the annotation (original aa position new aa). **Panel c. Current global distribution of**
174 **A.23 and A.23.1.** All available SARS-CoV-2 complete genomes annotated as complete and
175 lineage A from GISAID were retrieved on Feb 4 2021 and lineage typed using Pangolin(19).
176 and confirmed as A.23 and A.23.1 by extracting examining the encoded spike protein. All
177 novel Uganda A.23 and A.23.1 reported here were also included. Genomes were plotted by
178 country and sample collection date.
179
180

181 Of concern, the recent Kampala and global A.23.1 virus sequences from December
182 2020-January 2021 now encoded 4 or 5 amino acid changes in the spike protein (now
183 defining lineage A.23.1, see below) plus additional protein changes in nsp3, nsp6, ORF8
184 and ORF9 (Figure 3 panel b, Figure 4). The P681R spike change adds a basic amino acid
185 adjacent to the spike furin cleavage site. This same change has been shown *in vitro* to
186 enhance the fusion activity of the SARS-CoV-2 spike protein, likely due to increased
187 cleavage by the cellular furin protease (20); importantly, a similar change (P681H) is
188 encoded by the recently emerging VOC B.1.1.7 that is now spreading globally across 75
189 countries as of 5 February 2021 (5) (21). There are also changes in the spike N-terminal
190 domain (NTD), a known target of immune selection, observed in samples from Kampala
191 A.23.1 lineage, including P26S and R102I (Figure 3 panel b). Additionally and importantly, a
192 A.23.1 strain identified in Kampala on 11th December 2020 carried the E484K change in the
193 RBD, which may add further concern of this particular variant as it gains higher
194 transmissibility and enhanced resistance to vaccine and therapeutics. Outside of the spike
195 protein, a single nucleotide change (G27870T) leading to early termination of the ORF7b
196 (E39*) was observed in the A.23.1 from the community cases in Tororo in late December
197 2020. Although the clinical implication of this change is yet to be determined, it is important
198 to document such change for further follow-up.

199 **New lineage A designations.**

200 The viruses detected in Amuru and Kitgum met the criteria for a new SARS-CoV-2
201 lineage (22)(23) by clustering together on a global phylogenetic tree, sharing epidemiological
202 history and source from a single geographical origin, and encoding multiple defining SNPs.
203 These features including especially the three spike changes F157L, Q613H and V367F
204 define the new lineage A.23. Continued circulation and evolution of A.23 in Uganda was
205 observed and two additional changes in spike R102I and P681R were observed in
206 December 2020 in Kampala; these SNPs define the sub-lineage A.23.1. Additional changes
207 in non-spike regions also define the A.23 and A.23.1, including nsp3: E95K, nsp6: M86I,
208 L98F, ORF 8: L84S, E92K and ORF9 N: S202N, Q418H. These new lineages can be
209 assigned since pangolin version v2.1.10 and pangoLEARN data release 2021-02-01.

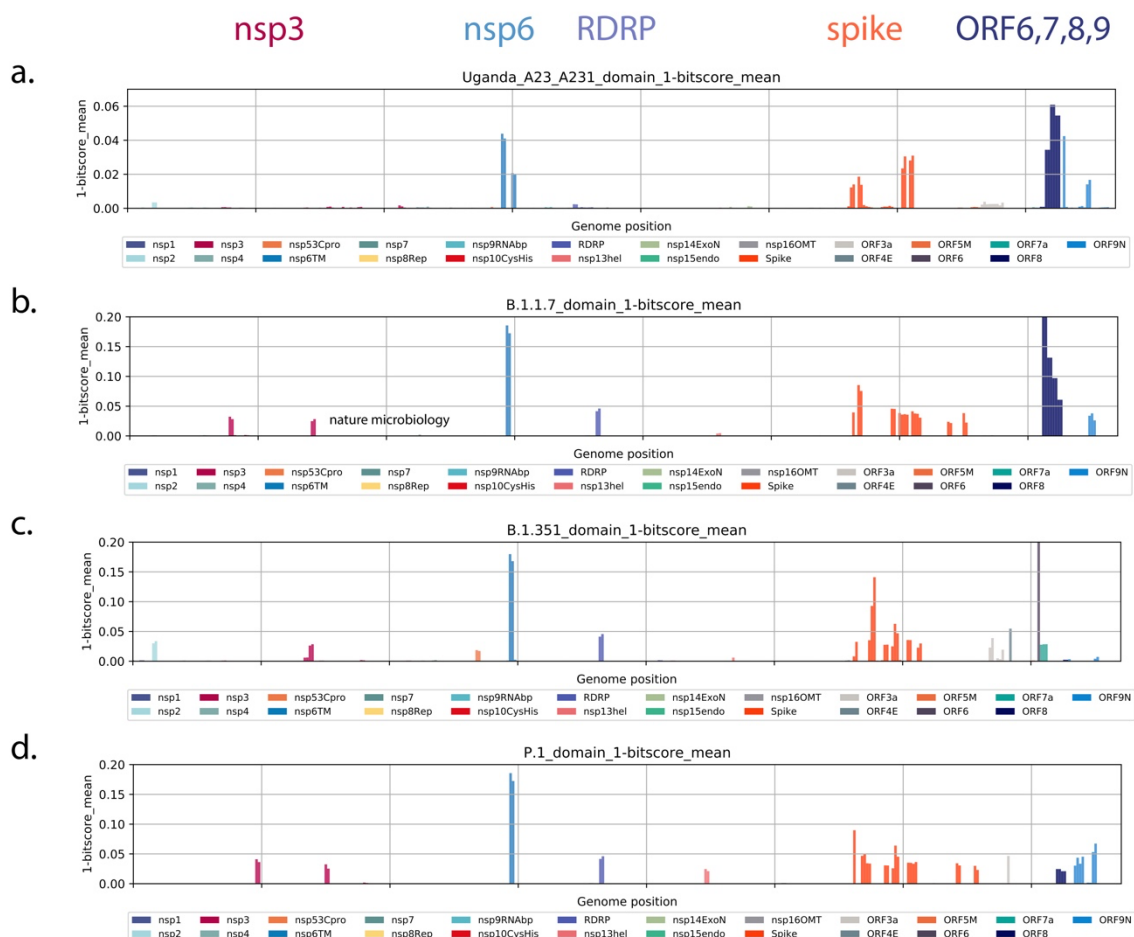
210 Screening SARS-CoV-2 genomic data from GISAID (March 12, 2021), A.23 and
211 A.23.1 viruses are now found in 26 countries outside of Uganda (Figure 3 panel c). The A.23
212 was first observed in Uganda in August 2020, subsequently in USA in October and Kenya
213 and Rwanda in December (Figure 3 panel c). The A.23.1 was first seen in Uganda in the
214 community cases in Mbale on 28th October 2020 and in Jinja on 29th October 2020, and
215 soon spreading across the country in early November 2020. Outside of Uganda, the A.23.1
216 was found in England and Cambodia from the end of November, in Rwanda from the
217 beginning of December. Of note, the international flights out of Uganda were restarted on 1
218 October 2020 with flights to Europe, Asia and USA. Phylogenetic analysis supports the close
219 evolution of A.23 to A.23.1 (Supplementary Figure 1).

220

221 **Additional changes in Ugandan A.23 and A.23.1 genomes compared to other VOC**
222 **genomes**

223 Although a main focus has been on spike protein changes, there are changes in
224 other genomic regions of the SARS-CoV-2 virus accompanying the adaptation to human

225 infection. We employed profile Hidden Markov Models (pHMMs) prepared from 44 amino
 226 acid peptides across the SARS-CoV-2 proteome (24) to detect and visualize protein
 227 changes from the early lineage B reference strain NC_045512. Measuring the identity score
 228 (bit-score) of each pHMMs across a query genome provides a measure of protein changes
 229 in 44 amino acid steps across the viral genome (Figure 4 panel a). This method applied to
 230 A.23 and A.23.1 genome sequences revealed the changes in spike (discussed above) as
 231 well as changes in the transmembrane protein nsp6 and the interferon modulators ORF8
 232 and 9 (Figure 5 panel a).



233

234 **Figure 4. All protein changes across lineage new variants.** All forward open reading
 235 frames from the 35 early lineage B SARS-CoV-2 genomes were translated, and processed
 236 into 44 aa peptides (with 22 aa overlap), clustered at 0.65 identity using uclust (25), aligned
 237 with MAAFT (26) and converted into pHMMs using HMMER-3 (27). The presence of each
 238 domains and its bit-score (a measure of the similarity between the query sequence and the
 239 sequences used for the pHMM(27)) was sought in each set of SARS-CoV-2 VOC genomes

240 and the 1-mean of the normalized domain bit-scores was plotted across the genome (e.g. 1 -
241 the similarity of the identified query domain to the reference lineage B SARS-CoV-2
242 domain). Domains were coloured by the proteins from which they were derived with the
243 colour code indicated below the figure. **Panel a.** Query set are 49 most Uganda lineage
244 A.23.1 genomes. **Panel b.** All B.1.1.7 full genomes lacking Ns deposited in GISAID on Jan
245 26 2021, **Panel c.** All B.1.351 full genomes lacking Ns present in GISAID on Jan 26 2021,
246 **Panel d.** All P.1 full genomes lacking Ns present in GISAID on Jan 26 2021.
247

248 We asked if a similar pattern of evolution was appearing in VOC as SARS-CoV-2
249 adapted to human infection. We gathered the sets of genomes described in the initial
250 published descriptions of these VOC (B.1.1.7 (5), B.1.351 (28) or P.1 (7)) and applied the
251 same profileHMM analysis. Similar to A.23/A.23.1, the B.1.1.7 lineage encodes nsp6, spike,
252 ORF8 and 9 changes as well as changes in nsp3 and the RNA-dependent RNA polymerase
253 (RRP, Figure 4 panel b). Lineage B.1.351 encodes nsp3, nsp6, RDRP, spike and ORF6
254 changes (Figure 4 panel c) and lineage P.1 encodes nsp3, nsp6, RDRP, nsp13, spike and
255 ORF8 and 9 changes (Figure 4 panel d). Although the exact amino acid and positions of
256 change within the proteins differ in each lineage, there are some striking similarities in the
257 common proteins that have been altered. Of interest, the nsp6 change present in B.1.1.7,
258 B.1.351 and P.1 is a 3 amino acid deletion (106, 107 and 108) in a protein loop of nsp6
259 predicted to be on exterior of the autophagy vesicles on which the protein accumulates
260 (29).The three amino acid nsp6 changes of lineage A.23.1 are L98F in the same exterior
261 loop region, and the M86I and M183I changes predicted to be in intramembrane regions but
262 adjacent to where the protein exits the membrane (29), (Supplementary Figure 2). A
263 compilation of the amino acid changes in A.23.1 and the VOC lineages is found in
264 Supplementary Table 2 with proteins that are altered in all 4 lineages marked in red.

265

266 **Discussion**

267 We report the emergence and spread of a new SARS-CoV-2 variant of the A lineage
268 (A.23.1) with multiple protein changes throughout the viral genome. A similar phenomenon
269 recently occurred with the B.1.1.7 lineage, detected first in the southeast of England (5) and

270 now globally and with the B.1.351 lineage in South Africa (6), and P.1 lineage in Brazil (30)
271 suggesting that local evolution (perhaps to avoid the initial population immune responses)
272 and spread may be a common feature of SARS-CoV-2. Importantly, lineage A.23.1 shares
273 many features found in the lineage B VOC including: alteration of key spike protein regions,
274 especially ACE2 binding region which is exposed and immunogenic, the furin cleavage site
275 and the 613/614 change that may increase spike multimer formation. The VOC and A.23.1
276 strains also encode changes in similar region of the nsp6 protein which may be important for
277 altering cellular autophagy pathways that promote replication. Changes or disruption of
278 ORF7,8 and 9 are also present in the VOC and A.23.1. The ORF8 changes or deletion
279 probably indicates this protein is unnecessary for human replication, similar deletions
280 accompanied SARS-CoV-2 adaption to humans(31),(32).

281 We suspect that emerging SARS-CoV-2 lineages may be adjusting to infection and
282 replication in humans and it is notable that the VOCs and lineage A.23.1 share some
283 common features in their evolution. The spike changes are best understood due to the
284 massive global effort to define the receptor and develop vaccines against the infection. The
285 analysis reported in Figure 4 reveals common functions of SARS-CoV-2 that have been
286 altered in all four variants, especially nsp6 and the ORFs 8 and 9. The functional
287 consequences of the additional non-spike changes warrant additional studies and the
288 current analysis may focus efforts of the proteins that are commonly changed in the variant
289 lineages. Finally the susceptibility of A.23.1 to vaccine immune responses is of great
290 importance to determine as vaccines become available in this part of Africa.

291 **Methods**

292 **Sample collection, whole genome MinION sequencing and genome assembly**

293 Residual nucleic acid extract from SARS-CoV-2 RT-PCR positive samples were
294 obtained from Central Public Health Laboratory (Kampala, Uganda). The nucleic acid was
295 converted to cDNA and amplified using SARS-CoV specific 1500bp-amplicon spanning the

296 entire genome as previously described(33).The resulting DNA amplicons were used to
297 prepare sequencing libraries, barcoded individually and then pooled to sequence on MinION
298 R.9.4.1 flowcells, following the standard manufacturer’s protocol.

299 The genome assemblies were performed as previously described (8). Briefly, reads
300 from fast5 files were base-called and demultiplexed using Guppy 3.6 running on the UMIC
301 HPC. Adapters and primers sequences were removed using Porechop
302 (<https://github.com/rrwick/Porechop>) and the resulting reads were mapped to the reference
303 genome Wuhan-1 (GenBank NC_045512.2) using minimap2(34) and consensus genomes
304 were generated in Geneious (Biomatters Ltd). Genome polishing was performed in Medaka,
305 and SNPs and mismatches were checked and resolved by consulting raw reads.

306

307 **Phylogenetic analyses**

308 For the local Uganda virus comparison, all newly and previously generated genomes
309 from Uganda (N=322) were aligned using MAFFT (26) and manually checked in AliView
310 (35). The 5’ and 3’ untranslated regions (UTRs) were trimmed. Maximum-likelihood (ML)
311 phylogenetic tree was constructed using RAxML-NG (36) under the GTR+I+G4 model as
312 best-fitted substitution model according to Akaike Information Criterion (AIC) determined by
313 ModelTest-NG (37) and run for 100 pseudo-replicates. Resulting tree was visualised in
314 Figtree(38) and rooted at the point of splitting lineage And B.

315 For phylogenetic analyses of Uganda lineage A.23 and A.23.1 strains comparing to
316 global A.23/A.23.1 strains, the global SARS-CoV-2 lineage A.23 (N=8) and A.23.1 (N=38)
317 genomes were retrieved from GISAID on 12 March 2021. These global A.23/A.23.1
318 genomes combining with Ugandan A.23/A.23.1 genomes (N=191) were aligned using
319 MAFFT and manually checked in AliView, followed by trimming 5’ and 3’ UTRs. The global
320 and Ugandan A.23/A.23.1 genomes were used to construct a ML tree under the GTR+I+G4
321 model as best-fitted substitution model according to AIC determined by ModelTest-NG (37)

322 and run for 100 pseudo-replicates using RAxML-NG. Resulting tree was visualised in Figtree
323 and rooted using the A.23 lineage.

324 Profile Hidden Markov Model (profileHMM) domain analysis of A.23/A.23.1 and VOC
325 genomes was performed as previously described (24) with some changes. A database of
326 profileHMMs was generated from the first 65 lineage B SARS-CoV-2 genome sequences.
327 All 3 forward open reading frames of each genome were translated computationally and
328 then sliced into 44 amino acid segment with overlapping with 22 amino acids. All 44 amino
329 acid query peptides were then clustered with uclust (25) and their original identity and
330 coordinates determined by blastp search against a protein database made from the
331 NC_045512 reference strain.

332 Query sets of genomes were processed to remove any genomes containing Ns
333 (which disrupt the HMM scoring process). The hmmscan function from HMMER-3 (27) was
334 used with the early B database. Query matches were identified using an E-value cutoff of
335 0.0001 and the bit-score values for each hit (a measure of the distance between the query
336 44 amino acid peptide and the B-lineage reference) was collected. Bit-scores for each
337 domain were normalized by dividing each query score by the maximum score for that
338 domain (x/x_{max}). In all analyses the original B lineage NC_045512 reference genome was
339 included to define the maximum bit-score.

340

341 **Ethical approvals**

342 This study was approved by the Uganda Virus Research Institute- Research and
343 Ethics Committee (UVRI-REC Federalwide Assurance [FWA] FWA No. 00001354, study
344 reference. GC/127/20/04/771) and by the Uganda National Council for Science and
345 Technology, reference number HS936ES. The novel reported SARS-CoV-2 genomes are
346 available on GISAID (<https://www.gisaid.org/>) under the accession numbers

347 EPI_ISL_954226-EPI_ISL_954300 and EPI_ISL_955136. A second tranche of genomes has
348 been submitted and is awaiting accession numbers.

349

350 **Acknowledgements**

351 We thank all global SARS-CoV-2 sequencing groups for their open and rapid sharing
352 of sequence data and GISAID for providing an effective platform for making these data
353 available. We are grateful to the Oxford Nanopore Technologies and the ARTIC Network for
354 their support and we thank Pope Moseley for his constructive comments on the manuscript.
355 The SARS-CoV2 diagnostic and sequencing award is jointly funded by the UK Medical
356 Research Council (MRC/UKRI) and the UK Department for International Development
357 (DFID) under the MRC/DFID Concordat agreement (grant agreement number
358 NC_PC_19060) and is also part of the EDCTP2 programme supported by the European
359 Union. The UMIC high performance computer was supported by MRC (grant number
360 MC_EX_MR/L016273/1) to PK. A.R. acknowledges the support of the Wellcome Trust
361 (Collaborators Award 206298/Z/17/Z [ARTIC network](#)) and the European Research Council
362 (grant agreement no. 725422 – ReservoirDOCS). The study is additionally funded by the
363 Wellcome, DFID - Wellcome Epidemic Preparedness – Coronavirus (grant agreement
364 number 220977/Z/20/Z) awarded to MC.

365

366 **Author contribution statement**

367 All authors contributed to the work presented in this paper.

368

369 **Competing Interests statement**

370 The authors declare no competing interests.

371

372 **References**

- 373 1. Edward C. Holmes, Yong-Zhen Zhang EC. Initial genome release of novel coronavirus.
374 Virological.org [Internet]. 2020 [cited 2021 Jan 24]; Available from:
375 <http://virological.org/t/319>
- 376 2. Li Q, Guan X, Wu P, Wang X, Zhou L, Tong Y, et al. Early Transmission Dynamics in
377 Wuhan, China, of Novel Coronavirus–Infected Pneumonia. *N Engl J Med*. 2020 Jan
378 29;NEJMoa2001316.
- 379 3. Yang X, Yu Y, Xu J, Shu H, Xia J, Liu H, et al. Clinical course and outcomes of critically
380 ill patients with SARS-CoV-2 pneumonia in Wuhan, China: a single-centered,
381 retrospective, observational study. *Lancet Respir Med*. 2020
382 Feb;S2213260020300795.
- 383 4. Rambaut A, Holmes EC, O’Toole Á, Hill V, McCrone JT, Ruis C, et al. A dynamic
384 nomenclature proposal for SARS-CoV-2 lineages to assist genomic epidemiology. *Nat*
385 *Microbiol*. 2020 Nov;5(11):1403–7.
- 386 5. Volz E, Mishra S, Chand M, Barrett JC, Johnson R, Geidelberg L, et al. Transmission of
387 SARS-CoV-2 Lineage B.1.1.7 in England: Insights from linking epidemiological and
388 genetic data [Internet]. *Infectious Diseases (except HIV/AIDS)*; 2021 Jan [cited 2021
389 Jan 29]. Available from: <http://medrxiv.org/lookup/doi/10.1101/2020.12.30.20249034>
- 390 6. Tegally H, Wilkinson E, Giovanetti M, Iranzadeh A, Fonseca V, Giandhari J, et al.
391 Emergence and rapid spread of a new severe acute respiratory syndrome-related
392 coronavirus 2 (SARS-CoV-2) lineage with multiple spike mutations in South Africa
393 [Internet]. *Epidemiology*; 2020 Dec [cited 2021 Jan 6]. Available from:
394 <http://medrxiv.org/lookup/doi/10.1101/2020.12.21.20248640>
- 395 7. Voloch CM, da Silva Francisco R, de Almeida LGP, Cardoso CC, Brustolini OJ, Gerber
396 AL, et al. Genomic characterization of a novel SARS-CoV-2 lineage from Rio de
397 Janeiro, Brazil. *J Virol*. 2021 Mar 1;
- 398 8. Bugembe DL, Kayiwa J, Phan MVT, Tushabe P, Balinandi S, Dhaala B, et al. Main
399 Routes of Entry and Genomic Diversity of SARS-CoV-2, Uganda. *Emerg Infect Dis*.
400 2020 Oct;26(10):2411–5.
- 401 9. Hadfield J, Megill C, Bell SM, Huddleston J, Potter B, Callender C, et al. Nextstrain: real-
402 time tracking of pathogen evolution. Kelso J, editor. *Bioinformatics*. 2018 Dec
403 1;34(23):4121–3.
- 404 10. Githinji G, de Laurent ZR, Mohammed KS, Omuoyo DO, Macharia PM, Morobe JM, et
405 al. Tracking the introduction and spread of SARS-CoV-2 in coastal Kenya [Internet].
406 *Epidemiology*; 2020 Oct [cited 2020 Dec 7]. Available from:
407 <http://medrxiv.org/lookup/doi/10.1101/2020.10.05.20206730>
- 408 11. Page AJ, Mather AE, Le Viet T, Meader EJ, Alikhan N-FJ, Kay GL, et al. Large scale
409 sequencing of SARS-CoV-2 genomes from one region allows detailed epidemiology
410 and enables local outbreak management [Internet]. *Epidemiology*; 2020 Sep [cited
411 2020 Oct 9]. Available from:
412 <http://medrxiv.org/lookup/doi/10.1101/2020.09.28.20201475>
- 413 12. Filipe ADS, Shepherd J, Williams T, Hughes J, Aranday-Cortes E, Asamaphan P, et al.
414 Genomic epidemiology of SARS-CoV-2 spread in Scotland highlights the role of

- 415 European travel in COVID-19 emergence [Internet]. *Infectious Diseases (except*
416 *HIV/AIDS)*; 2020 Jun [cited 2020 Dec 14]. Available from:
417 <http://medrxiv.org/lookup/doi/10.1101/2020.06.08.20124834>
- 418 13. Daily Monitor. Amuru prison closed as 153 test positive for Covid-19. 2020 Aug 23;
419 Available from: [https://www.monitor.co.ug/uganda/news/national/amuru-prison-closed-](https://www.monitor.co.ug/uganda/news/national/amuru-prison-closed-as-153-test-positive-for-covid-19-1924660)
420 [as-153-test-positive-for-covid-19-1924660](https://www.monitor.co.ug/uganda/news/national/amuru-prison-closed-as-153-test-positive-for-covid-19-1924660)
- 421 14. Penelope Nankunda. COVID-19: Uganda registers 318 new cases in a single day. *MSN*
422 [Internet]. 2020 Aug 22; Available from: [https://www.msn.com/en-xl/news/other/covid-](https://www.msn.com/en-xl/news/other/covid-19-uganda-registers-318-new-cases-in-a-single-day/ar-BB18gprA)
423 [19-uganda-registers-318-new-cases-in-a-single-day/ar-BB18gprA](https://www.msn.com/en-xl/news/other/covid-19-uganda-registers-318-new-cases-in-a-single-day/ar-BB18gprA)
- 424 15. Li Q, Wu J, Nie J, Zhang L, Hao H, Liu S, et al. The Impact of Mutations in SARS-CoV-2
425 Spike on Viral Infectivity and Antigenicity. *Cell*. 2020 Sep 3;182(5):1284-1294.e9.
- 426 16. Nguyen HT, Zhang S, Wang Q, Anang S, Wang J, Ding H, et al. Spike glycoprotein and
427 host cell determinants of SARS-CoV-2 entry and cytopathic effects. *J Virol*. 2020 Dec
428 11;
- 429 17. Gobeil SM-C, Janowska K, McDowell S, Mansouri K, Parks R, Manne K, et al. D614G
430 Mutation Alters SARS-CoV-2 Spike Conformation and Enhances Protease Cleavage at
431 the S1/S2 Junction. *Cell Rep*. 2021 Jan 12;34(2):108630.
- 432 18. Volz E, Hill V, McCrone JT, Price A, Jorgensen D, O'Toole Á, et al. Evaluating the
433 Effects of SARS-CoV-2 Spike Mutation D614G on Transmissibility and Pathogenicity.
434 *Cell*. 2020 Nov;S0092867420315373.
- 435 19. Áine O'Toole et al. Phylogenetic Assignment of Named Global Outbreak LINEages
436 (PANGOLIN). 2020; Available from: <https://github.com/cov-lineages/pangolin>
- 437 20. Hoffmann M, Kleine-Weber H, Pöhlmann S. A Multibasic Cleavage Site in the Spike
438 Protein of SARS-CoV-2 Is Essential for Infection of Human Lung Cells. *Mol Cell*. 2020
439 May;78(4):779-784.e5.
- 440 21. Áine O'Toole et al. B.1.1.7 report 2021-02-05. 2021; Available from: [https://cov-](https://cov-lineages.org/global_report_B.1.1.7.html)
441 [lineages.org/global_report_B.1.1.7.html](https://cov-lineages.org/global_report_B.1.1.7.html)
- 442 22. Áine O'Toole, JT McCrone, Verity Hill and Andrew Rambaut. Pangolin COVID-19
443 Lineage Assigner. Available from: <https://pangolin.cog-uk.io/>
- 444 23. Rambaut A, Holmes EC, Hill V, O'Toole Á, McCrone J, Ruis C, et al. A dynamic
445 nomenclature proposal for SARS-CoV-2 to assist genomic epidemiology [Internet].
446 *Microbiology*; 2020 Apr [cited 2020 Apr 27]. Available from:
447 <http://biorxiv.org/lookup/doi/10.1101/2020.04.17.046086>
- 448 24. Phan MVT, Ngo Tri T, Hong Anh P, Baker S, Kellam P, Cotten M. Identification and
449 characterization of Coronaviridae genomes from Vietnamese bats and rats based on
450 conserved protein domains. *Virus Evol* [Internet]. 2018 Jul 1 [cited 2020 Jan 12];4(2).
451 Available from: <https://academic.oup.com/ve/article/doi/10.1093/ve/vey035/5250438>
- 452 25. Edgar RC. Search and clustering orders of magnitude faster than BLAST.
453 *Bioinformatics*. 2010 Oct 1;26(19):2460–1.

- 454 26. Katoh K, Standley DM. MAFFT Multiple Sequence Alignment Software Version 7:
455 Improvements in Performance and Usability. *Mol Biol Evol.* 2013 Apr 1;30(4):772–80.
- 456 27. Eddy SR. Accelerated Profile HMM Searches. *PLOS Comput Biol.* 2011 Oct
457 20;7(10):e1002195.
- 458 28. Tegally H, Wilkinson E, Giovanetti M, Iranzadeh A, Fonseca V, Giandhari J, et al.
459 Emergence of a SARS-CoV-2 variant of concern with mutations in spike glycoprotein.
460 *Nature* [Internet]. 2021 Mar 9 [cited 2021 Mar 12]; Available from:
461 <http://www.nature.com/articles/s41586-021-03402-9>
- 462 29. Benvenuto D, Angeletti S, Giovanetti M, Bianchi M, Pascarella S, Cauda R, et al.
463 Evolutionary analysis of SARS-CoV-2: how mutation of Non-Structural Protein 6
464 (NSP6) could affect viral autophagy. *J Infect.* 2020 Jul;81(1):e24–7.
- 465 30. Voloch CM, Silva F R da, de Almeida LGP, Cardoso CC, Brustolini OJ, Gerber AL, et al.
466 Genomic characterization of a novel SARS-CoV-2 lineage from Rio de Janeiro, Brazil
467 [Internet]. *Genetic and Genomic Medicine*; 2020 Dec [cited 2021 Jan 30]. Available
468 from: <http://medrxiv.org/lookup/doi/10.1101/2020.12.23.20248598>
- 469 31. Su YCF, Anderson DE, Young BE, Linster M, Zhu F, Jayakumar J, et al. Discovery and
470 Genomic Characterization of a 382-Nucleotide Deletion in ORF7b and ORF8 during the
471 Early Evolution of SARS-CoV-2. Schultz-Cherry S, editor. *mBio.* 2020 Jul
472 21;11(4):e01610-20, /mbio/11/4/mBio.01610-20.atom.
- 473 32. The Chinese SARS Molecular Epidemiology Consortium. Molecular Evolution of the
474 SARS Coronavirus During the Course of the SARS Epidemic in China. *Science.* 2004
475 Mar 12;303(5664):1666–9.
- 476 33. Cotten M, Bugembe DL, Kaleebu P, Phan MVT. Alternate primers for whole-genome
477 SARS-CoV-2 sequencing [Internet]. *Genomics*; 2020 Oct [cited 2020 Nov 30]. Available
478 from: <http://biorxiv.org/lookup/doi/10.1101/2020.10.12.335513>
- 479 34. Li H. Minimap2: pairwise alignment for nucleotide sequences. Birol I, editor.
480 *Bioinformatics.* 2018 Sep 15;34(18):3094–100.
- 481 35. Larsson A. AliView: a fast and lightweight alignment viewer and editor for large datasets.
482 *Bioinformatics.* 2014 Nov 15;30(22):3276–8.
- 483 36. Kozlov AM, Darriba D, Flouri T, Morel B, Stamatakis A. RAxML-NG: a fast, scalable and
484 user-friendly tool for maximum likelihood phylogenetic inference. *Bioinforma Oxf Engl.*
485 2019 Nov 1;35(21):4453–5.
- 486 37. Darriba D, Posada D, Kozlov AM, Stamatakis A, Morel B, Flouri T. ModelTest-NG: A
487 New and Scalable Tool for the Selection of DNA and Protein Evolutionary Models. *Mol*
488 *Biol Evol.* 2020 Jan 1;37(1):291–4.
- 489 38. Rambaut A. FigTree <http://tree.bio.ed.ac.uk/software/figtree>. 2019;
- 490 39. Josh B. Singer, Gifford R, Cotten M, Robertson DL. CoV-GLUE project. 2020; Available
491 from: <http://cov-glue.cvr.gla.ac.uk/>
- 492 40. Flower TG, Buffalo CZ, Hooy RM, Allaire M, Ren X, Hurley JH. Structure of SARS-CoV-
493 2 ORF8, a rapidly evolving coronavirus protein implicated in immune evasion [Internet].

494 Biophysics; 2020 Aug [cited 2021 Mar 14]. Available from:
495 <http://biorxiv.org/lookup/doi/10.1101/2020.08.27.270637>

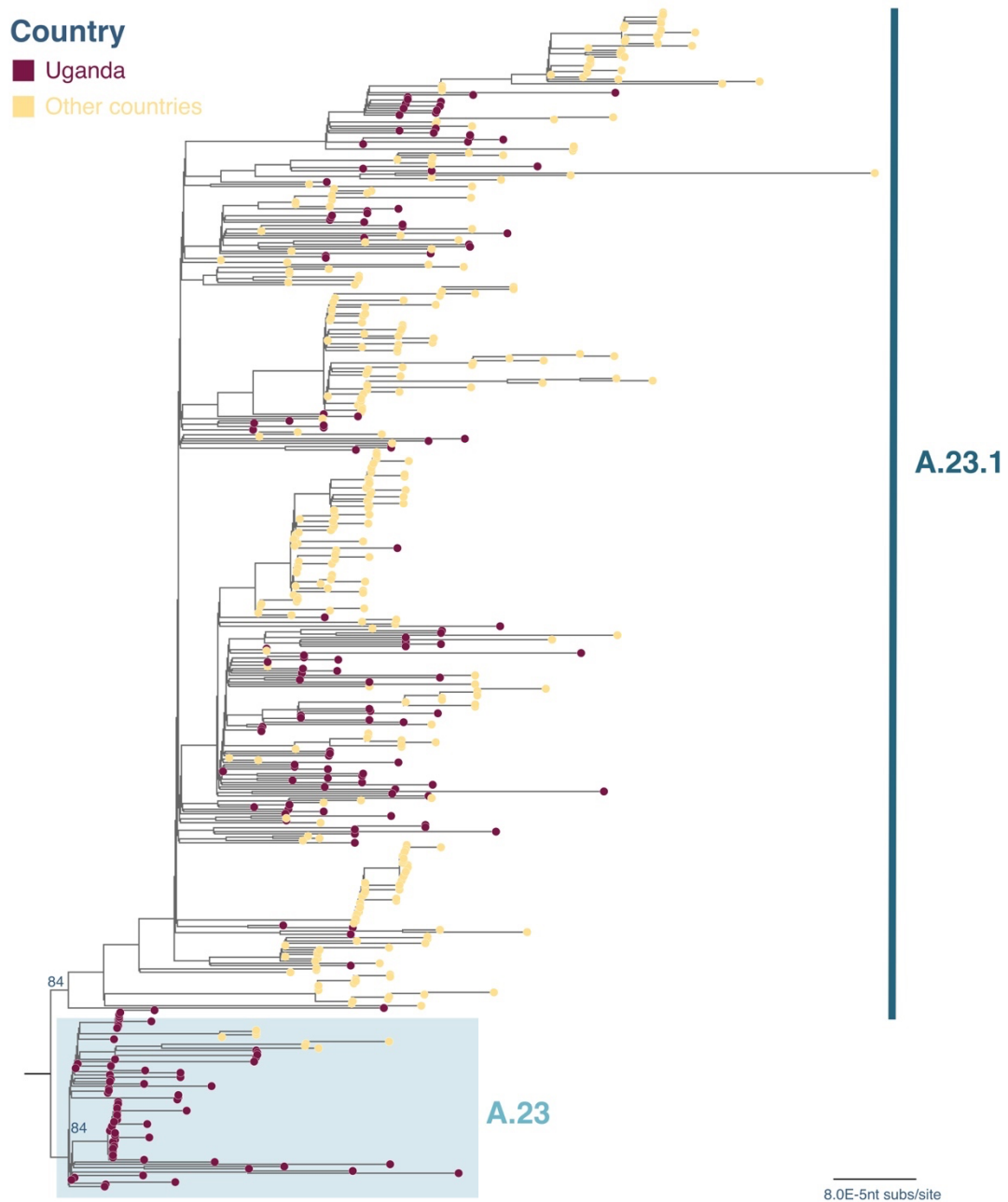
496 41. Chang C, Hou M-H, Chang C-F, Hsiao C-D, Huang T. The SARS coronavirus
497 nucleocapsid protein – Forms and functions. Antiviral Res. 2014 Mar;103:39–50.

498

499

500 **Supplementary Information**

501



502

503

504 **Supplemental Figure 1. Maximum-likelihood phylogenetic tree comparing**

505 **Uganda lineage A.23 and A.23.1 strains to global lineage A.23 and A.23.1**

506 **genomes.** A maximum-likelihood (ML) phylogenetic tree comparing Ugandan A.23

507 and A.23.1 (N=191) with the global A.23 and A.23.1 (N=336). The tree was rooted by

508 the A.23 lineage and strains were coloured according to the countries where they were

509 identified. Branch length was drawn to the scale of number of nucleotide substitutions

510 per site and only bootstrap values at the major nodes were shown. The tree was

511 visualised in Figtree (38).

512

513 **Supplementary Table 1. Lineage distribution in Uganda**

Lineage	Count in Uganda ¹	
A	8	
A.20	1	
A.22	1	
A.23	72	
A.23.1	119	
A.25	20	A: 221 (69%)
B	4	
B.1	35	
B.1.1.1	13	
B.1.1.161	1	
B.1.1.27	6	
B.1.1.29	1	
B.1.214	4	
B.1.222	3	
B.1.363	1	
B.1.380	2	
B.1.393	24	
B.1.398	1	
B.31	3	
B.4	3	B: 101(31%)
Total	322	

514

515 1. All available genomes from Uganda were classified into Pango lineages using pangolin
516 (19).

517

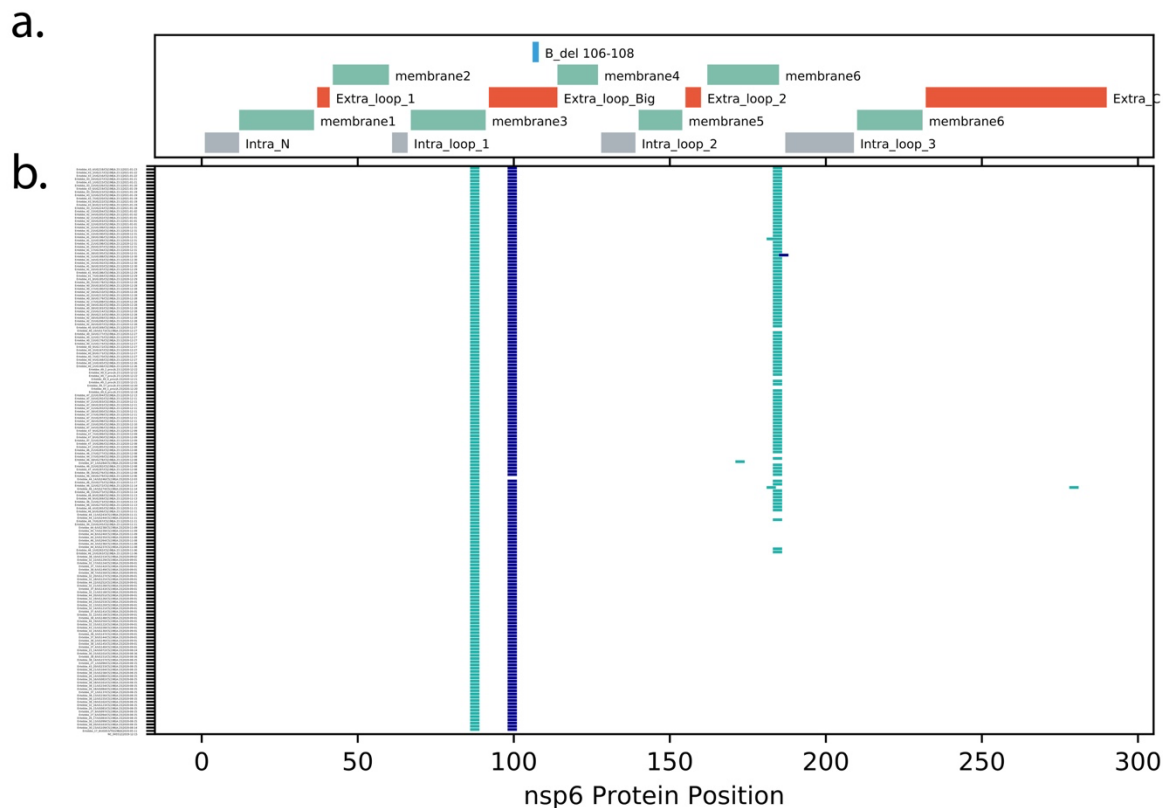
518 **Supplementary Information**

519 **Supplementary Table 2.** Summary of replacements in the A.23.1 and 3 VOC lineages¹.

Lineage	nsp2	nsp3	nsp5	nsp6	nsp12	spike	ORF3a	ORF4 E	ORF8	ORF9 N
A.23.1		nsp3: E95K		nsp6: M86I		S: R102I			ORF 8: L84S	N: S202N
A.23.1				nsp6: L98F		S: F157L			ORF 8: E92K	N: Q418H
A.23.1				nsp6: M183I		S: V367F				
A.23.1						S: Q613H				
A.23.1						S: P681R				
B.1.1.7		nsp3: T183I		nsp6: del_11288_11296	nsp12: P323L	S: N501Y			ORF 8: Q27*	N: D3L
B.1.1.7		nsp3: A690D				S: A570D			ORF 8: R52I	N: G204R
B.1.1.7		nsp3: I1412T				S: D614G			ORF 8: Y73C	N: R203K
B.1.1.7						S: P681H				N: S235F
B.1.1.7						S: T716I				
B.1.1.7						S: S982A				
B.1.1.7						S: D1118H				
B.1.1.7						S: del_21765_21770				
B.1.1.7						S: del_21991_21993				
B.1.351	nsp2: T85I	nsp3: K837N	nsp5: K90R	nsp6: del_11288_11296	nsp12: P323L	S: D80A	ORF 3a: Q57H	E: P71L		N: T205I
B.1.351						S: D215G				
B.1.351						S: K417N				
B.1.351						S: E484K				
B.1.351						S: N501Y				
B.1.351						S: D614G				
B.1.351						S: A701V				
B.1.351						S: non_cod_del_22281_22289				
B.1.281.1		nsp3: S970L		nsp6: del_11288_11296	nsp12: P323L	S: L18F	ORF 3a: S253P		ORF 8: E92K	N: P80R
B.1.281.1		nsp3: K977Q				S: T20N				N: R203K
B.1.281.1						S: P26S				N: G204R
B.1.281.1						S: D138Y				
B.1.281.1						S: R190S				
B.1.281.1						S: K417T				
B.1.281.1						S: E484K				
B.1.281.1						S: N501Y				
B.1.281.1						S: D614G				
B.1.281.1						S: H655Y				
B.1.281.1						S: T1027I				
B.1.281.1						S: V1176F				

520
521
522
523
524

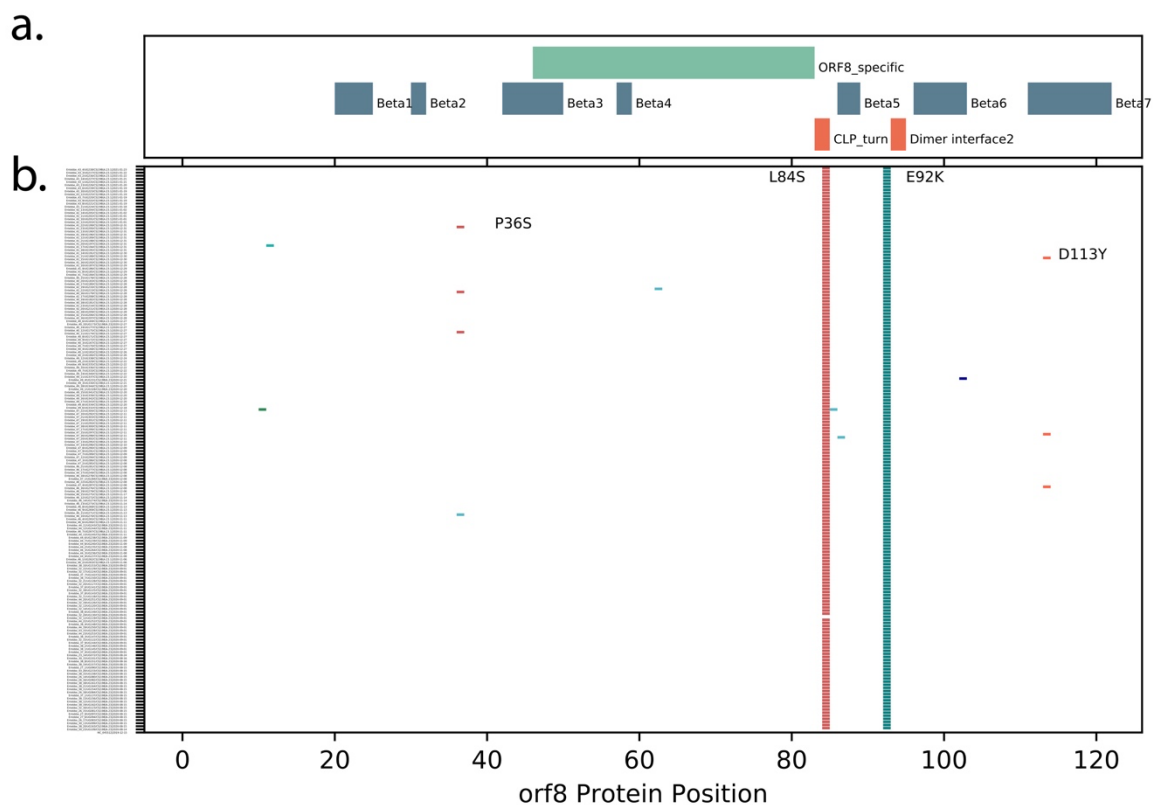
1. Representative sets of each lineage (either all available gap-free full genomes (B.1.351 and B.1.281.1) or all genomes gap-free full genomes annotated as B.1.1.7 were analyzed using CoV-GLUE (39) to identify frequent amino acid replacements associated with each lineage.



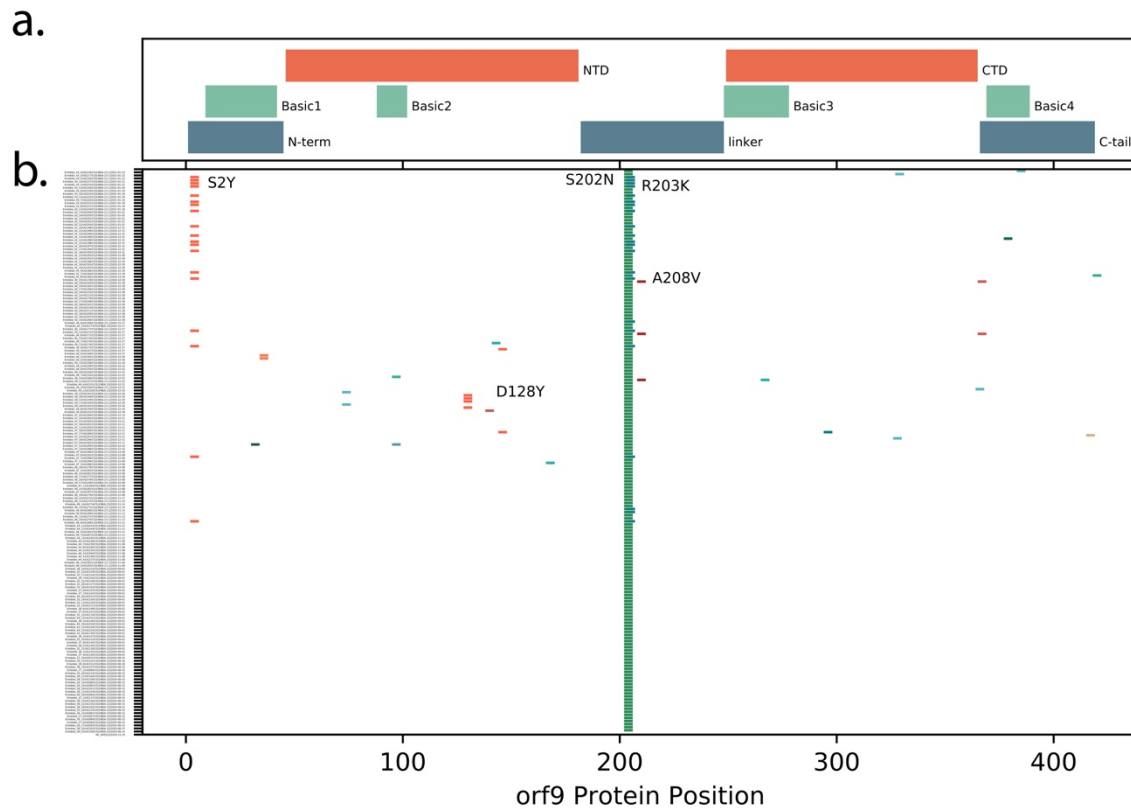
525
526
527
528

Supplementary Figure 2. Changes in A.23/A.23.1 nsp6 protein. The encoded nsp6 protein from all Ugandan A.23 and A.23.1 genomes gather, aligned and compared to the nsp6 protein

529 from GenBank NC_045512.2. **Panel a:** The locations of important nsp6 protein features are indicated based on the analysis of nsp6 from Benvenuto et al. (29). Intra_N: intravesicular
530 amino-terminal region, Extra_loop_1: extravesicular loop1, Intra_loop_1: intravesicular loop
531 1, B_del 106-108: the region of nsp6 deleted in the lineage B VOC genomes, Extra_loop_Big:
532 large extravesicular loop, Intra_loop_2: intravesicular loop 2, Extra_loop_2: extravesicular
533 loop 2, Intra_loop_3: intravesicular loop 3, Extra_C: carboxy-terminal extra-vesicular portion.
534 All features with "membrane" indicate membrane-spanning regions of nsp6. **Panel b:** Each
535 line represents the encoded nsp6 protein sequence from a single genome, ordered by date of
536 samples collection (bottom earliest, top most recent). Markers indicating the positions of amino
537 acid (aa) differences from the reference strain, changes observed in multiple genomes are
538 annotated with the annotation (original aa position new aa).
539
540



541 **Supplementary Figure 3. Changes in A.23/A.23.1 ORF8 protein.** The encoded ORF8
542 protein from all Ugandan A.23 and A.23.1 genomes gather, aligned and compared to the
543 ORF8 protein from GenBank NC_045512.2. **Panel a:** The locations of important ORF8 protein
544 features are indicated based on the analysis of ORF8 from Flower et al. (40). Features with
545 "Beta" indicate beta-sheets, ORF8_specific is a region unique to SARS-CoV-2 ORF8,
546 CLP_turn: indicates a cysteine, Leucine, Proline motif essential for a fold in the mature protein,
547 Dimer interface2 indicates the region of the protein the forms the interface between two
548 monomers. **Panel b:** Each line represents the encoded ORF8 protein sequence from a single
549 genome, ordered by date of samples collection (bottom earliest, top most recent). Markers
550 indicating the positions of amino acid (aa) differences from the reference strain, changes
551 observed in multiple genomes are annotated with the annotation (original aa position new aa).
552



553
554
555
556
557
558
559
560
561
562
563
564
565
566

Supplementary Figure 4. Changes in A.23/A.23.1 ORF9 protein. The encoded ORF9 protein from all Ugandan A.23 and A.23.1 genomes gather, aligned and compared to the ORF9 protein from GenBank NC_045512.2. **Panel a:** The locations of important ORF9 protein features are indicated based on the analysis of ORF9 from Chang et al.(41). N-term: amino-terminal extension, NTD: amino-terminal domain, linker: linker region between the NTD and CTD, CTD: carboxy-terminal domain, C-tail: carboxy-terminal extension, Regions with "Basic" indicate the 4 regions enriched in positively charged amino acids. **Panel b:** Each line represents the encoded ORF9 protein sequence from a single genome, ordered by date of samples collection (bottom earliest, top most recent). Markers indicating the positions of amino acid (aa) differences from the reference strain, changes observed in multiple genomes are annotated with the annotation (original aa position new aa).