

1

2 **Evolving Infection Paradox of SARS-CoV-2: Fitness Costs Virulence?**

3

4 A. S. M. Rubayet Ul Alam^{1#}, Ovinu Kibria Islam^{1#}, Md. Shazid Hasan¹, Mir Raihanul Islam²,
5 Shafi Mahmud³, Hassan M. Al-Emran⁴, Iqbal Kabir Jahid¹, Keith A. Crandall⁵, M. Anwar
6 Hossain^{6,7*}

7

8 1 Department of Microbiology, Jashore University of Science and Technology, Jashore-7408,
9 Bangladesh

10 2 BRAC James P Grant School of Public Health, BRAC University, Bangladesh

11 3 Genetic Engineering and Biotechnology, University of Rajshahi, Rajshahi-6205, Bangladesh

12 4 Department of Biomedical Engineering, Jashore University of Science and Technology,
13 Jashore-7408, Bangladesh

14 5 Computational Biology Institute and Department of Biostatistics & Bioinformatics, Milken
15 Institute School of Public Health, The George Washington University, Washington, DC,
16 USA

17 6 Jashore University of Science and Technology, Jashore-7408, Bangladesh

18 7 Department of Microbiology, University of Dhaka, Dhaka-1000, Bangladesh

19

20

21 ***Correspondence**

22 M. Anwar Hossain, Jashore University of Science and Technology, Jashore-7408,
23 Bangladesh.

24 E-mail: hossaina@du.ac.bd, Contact: +8801708818101

25

26

27 # Authors contributed equally

28

29 **ABSTRACT**

30 **Background:** SARS-CoV-2 is continuously spreading worldwide at an
31 unprecedented scale and evolved into seven clades according to GISAID where four (G, GH,
32 GR and GV) are globally prevalent in 2020. These major predominant clades of SARS-CoV-
33 2 are continuously increasing COVID-19 cases worldwide; however, after an early rise in
34 2020, the death-case ratio has been decreasing to a plateau. G clade viruses contain four co-
35 occurring mutations in their genome (C241T+C3037T+C14408T:
36 RdRp.P323L+A23403G:spike.D614G). GR, GH, and GV strains are defined by the presence
37 of these four mutations in addition to the clade-featured mutation in GGG28881-
38 28883AAC:N. RG203-204KR, G25563T:ORF3a.Q57H, and
39 C22227T:spike.A222V+C28932T-N.A220V+G29645T, respectively. The research works are
40 broadly focused on the spike protein mutations that have direct roles in receptor binding,
41 antigenicity, thus viral transmission and replication fitness. However, mutations in other
42 proteins might also have effects on viral pathogenicity and transmissibility. How the clade-
43 featured mutations are linked with viral evolution in this pandemic through gearing their
44 fitness and virulence is the main question of this study.

45 **Methodology:** We thus proposed a hypothetical model, combining a statistical and
46 structural bioinformatics approach, endeavors to explain this infection paradox by describing
47 the epistatic effects of the clade-featured co-occurring mutations on viral fitness and
48 virulence.

49 **Results and Discussion:** The G and GR/GV clade strains represent a significant
50 positive and negative association, respectively, with the death-case ratio (incidence rate ratio
51 or IRR = 1.03, $p < 0.001$ and IRR= 0.99/0.97, $p < 0.001$), whereas GH clade strains showed
52 no association with the Docking analysis showed the higher infectiousness of a spike mutant
53 through more favorable binding of G614 with the elastase-2. RdRp mutation p.P323L
54 significantly increased genome-wide mutations ($p < 0.0001$) since more expandable RdRp
55 (mutant)-NSP8 interaction may accelerate replication. Superior RNA stability and structural
56 variation at NSP3:C241T might impact upon protein or RNA interactions. Another silent
57 5'UTR:C241T mutation might affect translational efficiency and viral packaging. These G-
58 featured co-occurring mutations might increase the viral load, alter immune responses in host

59 and hence can modulate intra-host genomic plasticity. An additional viroporin
60 ORF3a:p.Q57H mutation, forming GH-clade, prevents ion permeability by cysteine (C81)-
61 histidine (H57) inter-transmembrane-domain interaction mediated tighter constriction of the
62 channel pore and possibly reduces viral release and immune response. GR strains, four G
63 clade mutations and N:p.RG203-204KR, would have stabilized RNA interaction by more
64 flexible and hypo-phosphorylated SR-rich region. GV strains seemingly gained the
65 evolutionary advantage of superspreading event through confounder factors; nevertheless,
66 N:p.A220V might affect RNA binding.

67 **Conclusion:** These hypotheses need further retrospective and prospective studies to
68 understand detailed molecular and evolutionary events featuring the fitness and virulence of
69 SARS-CoV-2.

70

71 **Key words**

72 SARS-CoV-2, COVID-19, Infection Paradox, Fitness, Virulence, Clades, Co-occurring
73 mutations

74

75

76

77

78

79

80 **Highlights**

- 81 • We speculated an association of particular SARS-CoV-2 clade with death rate.
- 82 • The polymerase mutant virus can speed up replication that corresponds to higher
83 mutations.
- 84 • The impact on viral epistasis by evolving mutations in SARS-CoV-2.
- 85 • How the virus changes its genotype and circulate with other types given the overall
86 dynamics of the epidemics?

- 87 • Human intervention seems to work well to control the viral virulence. This hygiene
88 practice will control the overall severity of the pandemic situation as recommended by the
89 WHO. Our work has given the same message but explain with the dominant co-occurring
90 mutations.

91 **1. Introduction**

92 Severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) has caused the COVID-
93 19 pandemic since the beginning of 2020 ¹. This highly contagious virus spread in 213
94 countries infecting over 11.3 million people including half a million deaths within the first six
95 months ². SARS-CoV-2 has possessed some extraordinary attributes that make it extremely
96 infectious: high replication rate, large burst size, high stability in the environment, strong
97 binding efficiency of S protein receptor-binding domain (RBD) with human angiotensin-
98 converting enzyme 2 (ACE2) receptor, and additional furin cleavage site in S protein ³⁻⁵. In
99 addition to those, it has proofreading capability ensuring high-fidelity replication ⁶. The virus
100 contains four major structural proteins: spike glycoprotein (S), envelope (E), membrane (M),
101 and nucleocapsid (N) proteins along with 16 nonstructural proteins (NSP1 to NSP16) and
102 seven accessory proteins (ORF3a, ORF6, ORF7a, ORF7b, ORF8a, ORF8b, and ORF10) ^{7,8}.

103 Mutational spectra within the SARS-CoV-2 genome ^{9,10}, spike protein ¹¹, RdRp ¹²,
104 ORF3a ¹³, and N protein ¹⁴ are reported where Islam et al. (2020) and Rahman et al. (2020)
105 studied the possible impact of mutations on the virulence of a strain ¹⁰. SARS-CoV-2 has
106 been classified into seven major clades, such as G, GH, GR, GV, S, V, and L by GISAID
107 based on the dominant mutations ¹⁵. Yin ¹⁶ reported the leader sequence mutation 241C > T is
108 co-occurring with three other mutations, 3037C > T (NSP3: C318T), 14408C > T (RdRp:
109 p.P323L), and 23403A > G (S: p.D614G). GISAID referred these co-occurring mutation
110 containing viruses as clade G, a viral strain over the wild-type ^{17,18}. These clade G or lineage
111 B.1 viruses were dominant in Europe ¹⁷ and the east coast of USA ¹⁹ during the earlier stages
112 of the pandemic that further spread in Southeast Asia ^{20,21} and Oceania ¹⁸. Remarkably, this
113 mutation variant is reported to be more transmissible ¹⁷ and was speculated to cause high
114 mortality in USA ²². The GR clade or lineage B.1.1 was classified with additional
115 trinucleotide mutations at 28881-28883 (GGG>AAC) creating two consecutive amino acid
116 changes, R203K and G204R, in N protein ¹⁸. The GR strains are now the dominant type
117 causing more than one-third of infection cases globally ¹⁵. Another derivative of G clade is
118 GH or lineage B.1.3-B.1.66, B.1.351, characterized by an additional ORF3a:p.Q57H
119 mutation frequently found in the USA and Europe ^{18,23}. The newest variant GV or lineage

120 B.1.177 is responsible for COVID-19 cases in Europe featuring an A222V mutation in the S
121 protein with other mutations of the clade G. N: A220, ORF10: V30L, and three other
122 synonymous mutations, namely T445C, C6286T, and C26801G, are also found in this clade
123 ²⁴.

124 The most frequently observed mutation is D614G of S protein ²⁵, which has direct roles in
125 receptor binding, and immunogenicity, thus viral immune-escape, transmission, and
126 replication fitness. The role of other dominant mutations remains largely underestimated.
127 However, the ORF3a:p.Q57H ²³ and RdRp:p.P323L ²⁶ have recently been investigated while
128 the effect of NSP3: C318T, 5'UTR: C241T, N:p.RG203-203KR, and N:p.A220V is still being
129 overlooked. Mutations in other proteins could also have effects on viral pathogenicity and
130 transmissibility. How the clade-featured mutations are linked with viral evolution in this
131 pandemic through gearing their fitness and virulence is the main question of this study.

132 Different mutations may work independently or through epistatic interaction, thus it is
133 difficult to spot exactly when a single mutation or co-occurring mutations become dominant
134 in populations through evolutionary fitness. Several co-occurring mutations identified in
135 SARS-CoV-2 may have complementary roles in changing its virulence and the infection
136 dynamics of COVID-19 pandemic. Yet, many questions remain: When are 'co-occurring
137 mutations' first observed? Why are particular clades, GV and GR, outcompeting their
138 speculative ancestor G clade? What are the impacts of these mutations on protein structures
139 and what are their functional roles? Do these mutant proteins interact together? Is there any
140 role of the coevolved 'silent' mutations? Do these mutations have any impact on viral fitness
141 and virulence? We attempt to answer these questions by *in silico* molecular insights of
142 SARS-CoV-2 mutants and possible interactions of coevolved proteins.

143 This study aims to identify the prevalence of dominant co-occurring mutations of G, GH,
144 GR, and GV strains circulating in all over the world, correlation of the clades with death-case
145 ratio, probable individual and/or synergistic impact of those mutants upon virulence in terms
146 of viral entry and fusion, evasion of host cell lysis, replication rate, ribonucleoprotein
147 stability, protein-protein interactions, translational capacity, and ultimately the combined
148 effect on fitness and virulence.

149 **2. Materials and Methods**

150 **2.1 Retrieval of Sequences and Mutation Analyses**

151 This study analyzed 225,526 high-coverage (<1% Ns and <0.05% unique amino acid
152 mutations) and complete (>29,000 nucleotide) genome sequences with specified collection
153 date from a total of 3,16,166 sequences submitted to GISAID until January 03, 2021. We
154 sifted the sequences generated from the non-human host out from the dataset. The Wuhan-
155 Hu-1 (Accession ID- NC_045512.2)²⁷ isolate was used as the reference genome.

156 A python script was used to partition a significant part of the dataset into two subsets
157 based on the RdRp: C14408T mutation and estimated the genome-wide variations (single
158 nucleotide changes) for each strain. For the genome-wide mutation analysis, a total of 37,179
159 sequences (RdRp wild type or ‘C’ variant: 9,815; and mutant or ‘T’ variant: 27,364) were
160 analyzed from our dataset. The frequency of mutations was tested for significance with the
161 Wilcoxon signed-rank test between RdRp ‘C’ variant and ‘T’ variant using IBM SPSS
162 statistics 25.

163 Random effect poisson regression model was performed in STATA v13.0 to identify the
164 association between death-case ratio and different clade strains (G, GH, GR, and GV); both
165 unadjusted and adjusted incidence risk ratio (IRR) were estimated where time was introduced
166 as a panel variable²⁸.

167 **2.2 Epidemiological Data Analysis and Time Plot Generation**

168 In this study, we report the prevalence of these dominant clades in 2020, both
169 individually and in combination, with disease progression and deaths allowing us to infer
170 increasing fitness of the SARS-CoV-2. A weekly-based time plot of G, GH, GR, and GV
171 clade frequencies with infection and death-cases was generated from 23 December 2019 until
172 January 3, 2021 that counted a total of 54 weeks (supplementary table S1). The total number
173 of infections and deaths by weeks were extracted from the WHO 2019-nCoV situation
174 reports. The case-death ratio was estimated by dividing the number of deaths of a particular
175 week by the number of cases identified in the earlier week based on the conservative
176 assumption of a one-week interval between diagnosis and death²⁹.

177 Regional time plots of those clades were also generated monthly (from January to
178 June) with frequencies of new infections, deaths, and death-case ratio based on the available
179 data on WHO situation reports³⁰. The studied sequences were divided into six regions;
180 Europe (n= 145,254), Americas (n= 48,014), Eastern Mediterranean (n=3,103), Southeast
181 Asia (n= 4,134), West Pacific (n= 16,974), and Africa (n= 2,740).

182 **2.3 Stability, Secondary and Three-Dimensional Structure Prediction Analyses of S,**
183 **RdRp, ORF3a, and N Proteins**

184 DynaMut ³¹ and FoldX 5.0 ^{32,33} were used to determine the stability of both wild and
185 mutant variants of N, RdRp, S, and ORF3a proteins. PredictProtein ³⁴ was utilized for
186 analyzing and predicting the possible secondary structure and solvent accessibility of both
187 wild and mutant variants of those proteins. The SWISS-MODEL homology modeling
188 webtool ³⁵ was utilized for generating the three-dimensional (3D) structures of the RdRp, S,
189 and ORF3a protein using 7c2k.1.A, 6xr8.1.A, and 6xdc.1.A PDB structure as the template,
190 respectively. Modeller v9.25 ³⁶ was also used to generate the structures against the same
191 templates. I-TASSER ³⁷ with default protein modeling mode was employed to construct the
192 N protein 3D structure of wild and mutant type since there was no template structure
193 available for the protein. The built-in structural assessment tools (Ramachandran plot,
194 MolProbity, and Quality estimate) of SWISS-MODEL were used to check the quality of
195 generated structures.

196 **2.4 Molecular Docking and Dynamics of RdRp-NSP8 and S protein-Elastase2**
197 **Complex**

198 Determination of the active sites affected by binding is a pre-requisite for docking
199 analysis. We chose 323 along with the surrounding residues (315-324) of RdRp and the
200 residues 110 to 122 of NSP8 monomer as the active sites based on the previously reported
201 structure ³⁸. The passive residues were defined automatically where all surface residues were
202 selected within the 6.5°A radius around the active residues. The molecular docking of the
203 wild and predicted mutated RdRp with the NSP8 monomer from the PDB structure 7C2K
204 was performed using the HADDOCKv2.4 to evaluate the interaction ³⁹. The binding affinity
205 of the docked RdRp-NSP8 complex was predicted using the PRODIGY ⁴⁰. The number and
206 specific interfacial contacts (IC) for each of the complexes were identified.

207 The human neutrophil elastase (hNE) or elastase-2 (PDB id: 5A0C) was chosen for
208 docking of the S protein, based on earlier reports ⁴¹. Here we employed CPORT ⁴² to find out
209 the active and passive protein-protein interface residues of hNE. The S protein active sites
210 were chosen based on the target region (594-638) interacting with the elastase-2. The passive
211 residues of S protein were defined automatically as mentioned for RdRp-NSP8 docking
212 analysis. Afterward, we individually docked wild (614D) and mutated (614G) S protein with
213 the hNE using HADDOCK 2.4. The binding affinity of the docked complexes, as well as the

214 number and specific interfacial contacts (IC), were predicted as performed after RdRP-NSP8
215 docking. We employed HDOCK server⁴³ with specifying the active binding sites residues for
216 predicting the molecular docking energy.

217 The structural stability of the protein complex and their variations were assessed
218 through YASARA Dynamics software package⁴⁴. The AMBER14 force field⁴⁵ was used for
219 these four systems, and the cubic simulation cell was created. The TIP3P (at 0.997 g/L⁻¹, 25C,
220 and 1 atm) water solvation model was used, and steepest gradient energy minimization
221 techniques by simulated annealing method was used for initial energy minimization process.
222 The hydrogen bond network system was optimized. The simulation cell was extended 20Å in
223 all cases of the protein to move freely⁴⁶. The PME or particle mesh Ewald methods was
224 applied to calculate the long-range electrostatic interaction by a cut off radius of 8Å⁴⁷. The
225 Berendsen thermostat was applied to maintain the temperature of the simulation cell. The
226 simulation system was neutralized by the addition of 0.9% NaCl, pH 7.4, and 310K
227 temperature. The time step of the simulation was set as 1.25fs⁴⁶. The simulation trajectories
228 were saved after every 100ps. Finally, the molecular dynamics simulation was conducted for
229 100ns to analyze root mean square deviation (RMSD), root mean square fluctuation (RMSF),
230 radius of gyration (Rg), solvent accessible surface area (SASA), hydrogen bond⁴⁸⁻⁵¹.

231 **2.5 Mutational Analysis of Transmembrane Domain 1 of ORF3a and SR-domain** 232 **of N protein**

233 The complete genome of 12 pangolin derived coronavirus strains, as well as 38 bat,
234 civet and human SARS-CoVs, were downloaded from GISAID and NCBI, respectively for
235 the mutational comparison between the SARS-CoV and SARS-CoV-2. We mainly targeted
236 transmembrane domain 1 (TM1), which covers 41 to 63 residues, of ORF3a to find the
237 identical mutation and scan overall variation in TM1. In the case of the N protein, CoVserver
238 of GISAID was used to study the frequent mutations other than our target ones (RG→KR:
239 203-04) to bolster our speculation on the change and subsequent effects of those surrounding
240 variations on the functions of the protein. A generalized comparison between SARS-CoV and
241 SARS-CoV-2 reference sequences was performed to identify mutation in the SR-rich region
242 that will help to postulate on N protein functions of novel coronavirus based on previous
243 related research on SARS-CoV.

244 **2.6 Analyzing RNA Folding prediction of 5'UTR and NSP3**

245 The Mfold web server ⁵² was used with default parameters to check the folding
246 pattern of RNA secondary structure in the mutated 5'UTR, synonymous leader (T445C) and
247 NSP3 regions (C3037T). The structure of complete mutant 5'UTR (variant 'T') was
248 compared with the wild type (variant 'C') secondary pattern as mentioned in the Huston, et
249 al. ⁵³. Since the wild type (variant 'C' at 318th nucleotide) RNA structure of the NSP3 was not
250 available in the literature, we generated the structure of mutant (variant 'T') to predict the
251 RNA folding in the Mfold web server. From the Mfold web server, we also estimated free
252 energy change (ΔG) for wild and mutant NSP3 RNA fold.

253 3. RESULTS AND DISCUSSION

254 We represented the global scenario of SARS-CoV-2 infection by the G, GH, GR, and GV
255 clade strains and estimated the association between the clade strains and death-case ratio.
256 Afterwards, the possible effects of the nine mutations in S, RdRp, ORF3a, N, 5'UTR, leader
257 protein, and NSP3 were discussed with associated results. Whereas researches on molecular
258 docking of the spike protein ^{54,55} and RdRp ^{56,57} in search of potential drug targets is a
259 continuous process, our study approached in a unique way to dock spike with elastase-2 and
260 RdRp with NSP8 to satisfy our purpose. The overall epistatic interactions of the mutant
261 proteins and/or RNA was then depicted (Figure 1) with appropriate explanation. Finally, we
262 endeavored to postulate using theoretical evolutionary theory how the virus might be
263 changing virulence and fitness.

264 3.1 Global Emergence of Dominant Co-occurring Mutations

265 Analysis of the SARS-CoV-2 genome sequences has indicated that 241C > T, 3037C > T,
266 14408C > T, and 23403A > G mutations were discretely identified among different viruses in
267 China on 24th January of 2020. These four mutations together in a single virus was first
268 detected in England on 3rd February 2020 (Table 1). Since then, those mutations were found
269 to cooccur with other mutations, thus formed clade G, GH, GR, and GV, and have become
270 the most dominant variants in other regions of the world (Figure 2b), for example, escalating
271 to 85% in May 2020 in Southeast Asia ²⁰. The G clade strain circulated predominantly (30%,
272 n=828) in Africa, whereas the GH, GR, and GV clades have become more recognized in the
273 Americas, Western Pacific and Europe, respectively (Table 1).

274 Our weekly based time-plot has depicted a gradual increase of G, GH, and GR viruses
275 altogether since the 10th week (24 February – 2 March 2020) recording a sudden jump to 42%
276 in that week from a mere 12% of the previous week. The global COVID-19 cases
277 exponentially increased from the 10th week with only 7,806 cases and infected almost

278 500,000 people with ca. 35,000 deaths in just seven weeks while the G, GH, and GR strains
279 reached 80% (supplementary table S1). The death rate elevated at its peak (14%) in mid of
280 March (week 12) and gradually decreased to below 2% in early September. Correspondingly,
281 the number of GR strains is increasing from the 10th week to 32nd week with a small
282 fluctuation where the G and GH strains each maintained a static ratio between ~20-30%. GV
283 clade strains initiated into the population from 30th week and became most dominant strain
284 replacing GR clade strains in just 9 weeks (Figure 2a).

285 The geographical distribution plot of the G, GH, GR, and GV clades with infection and
286 death-cases delineates that the new infections began to rise exponentially with the increase of
287 coevolved mutant variants in all regions except the West Pacific area (Figure 2b). The West
288 Pacific region, which includes East Asian countries as well, identified a very low number of
289 infection cases and deaths per million (Figure 2b). Europe and America have a high rate of
290 infections as well as a high percentage of those variants. The reason may be linked with α 1-
291 antitrypsin (AAT) deficiency, which is very rare in East Asia unlike Europe and North
292 America⁵⁸. The AAT allele deficiency facilitates entry of the 614G subtype into the host
293 cells and accelerates the spread of G, GH, GR, and GV clades. Our analysis has also found
294 that the proportion of strains containing 23403A > G mutation was 25% in East Asia and
295 >75% in Europe and America in first 6 months of this year (data not shown).

296 The death-case ratio was decreasing globally while the GR mutants were increasing until
297 August (Figure 2a). Since September, GV clade strains were increasing while the death-case
298 ratio remained low (2%). To examine the association between clade strains and death-case
299 ratio, both unadjusted and adjusted incidence risk ratio were estimated. In the adjusted model,
300 G, GR, and GV clades were found to be significantly associated with death-case ratio in both
301 models. If G clade strains increase by one percentage point then death-case ratio would be
302 expected to increase by a factor of 1.03, while holding all other variables in the model
303 constant (IRR = 1.03, 95% CI 1.01 – 1.06). However, GR and GV clades were inversely
304 related with death-case ratio. One percentage point increment in GR clade led to reduction to
305 death-case ratio by a factor of 0.99, while holding all other variables in the model constant
306 (IRR = 0.98, 95% CI 0.97 – 0.99). And then again, If GV clade strains increase by one
307 percentage point then death-case ratio would be expected to decrease by a factor of 0.97,
308 while holding all other variables in the model constant (IRR = 0.97, 95% CI 0.96 – 0.98).

309 Like other SARS-CoV-2 studies⁵⁹⁻⁶¹, this statistical analysis also suffers from some
310 limitations in dealing with genomic and calculating death-case ratio data. The death-case
311 ratio is believed to be underestimated because of the inadequate number of tests capacity and

312 asymptomatic SARS-CoV-2 cases in the general population. Moreover, fewer mutation
313 patterns are uploaded from underdeveloped or developing countries like African and Sub-
314 Saharan countries which might lead to spatial biasedness in the analysis. Therefore, the
315 global epidemiological scenario of different clades was explored to mitigate this problem.

316 Regional monthly data depicted a similar increase of GR and GV strains while the death
317 ratio was decreasing in studied regions with some rare exceptions. In Europe, GR strains
318 were predominant from April to August, while GV strains became predominant from
319 September. In the Americas a high abundance (over 50%) of GH clade strains was found
320 from March to December. In Eastern Mediterranean, Africa and Southeast Asian region there
321 was an increased rise of GR strains in May, while in Western Pacific region GR strains
322 became predominant in June. The increase of G strains and the death-case ratio at the same
323 time was observed in most regions, except Eastern Mediterranean and Southeast Asia, where
324 a limited number of sequence data were produced. (Figure 2b).

325 Researchers around the globe are now trying to explore the factors associated with
326 variable mortality rates due to the infection by SARS-CoV2 in different regions of the world.
327 A recent report identified a Neanderthal-derived haplotype in chromosome 3 which is
328 prevalent in South Asia as a substantial contributor to COVID-19 risk for hospitalized
329 patients along with other risk factors ⁶². Also, a haplotype inherited from Neandertals on
330 chromosome 12 was found at ~25 to 30% in Eurasia and the Americas, can reduce relative
331 risk of becoming severely ill with COVID-19 by ~22% ⁶³. But other reports say that
332 mortality associated to COVID-19 in Indian and south Asian subcontinent is lower than in the
333 west ⁶⁴. In-hospital mortality was also found to be considerably lower in Asia than Europe
334 and America ⁶⁵. Several factors such as the virulence of the pathogen, host factors like innate
335 immunity, genetic diversity in immune responses, epigenetic factors, genetic polymorphisms
336 of ACE2 receptors, micro RNAs and universal BCG vaccination, along with environmental
337 factors may have contributed in low mortality in Indian region ⁶⁶. In European countries,
338 mortality rates were observed to vary from region to region ⁶⁷. In Africa, where the
339 Neanderthal risk haplotype is almost absent, youthful population, and weather conditions
340 might have also contributed to low morbidity and mortality of COVID-19 ⁶⁸. No association
341 between the SARS-CoV-2 variants and mortality rates was observed in the Eastern
342 Mediterranean Region ⁶⁹. However, another report showed that, GR, GH and L clade viruses
343 are predominant in countries with higher deaths and GR clade showed higher prevalence
344 among severe/deceased patient ⁷⁰. Besides these host, pathogen and environment associated

345 factors, human and social intervention like political decision making, scientific advice, and
346 health system and public health capacity can also contribute to variable COVID-19 mortality
347 rates in different regions of the world ⁷¹.

348 **3.2 S Protein D614G Mutation Favors Elastase-2 Binding**

349 Korber et al. (2020) linked the S:p.D614G mutation with viral spread or transmission
350 capacity rather than virus infectivity or virulence [17] whereas Becerra Flores et al. (2020)
351 approached to demonstrate a potential link of D614G with virus pathogenicity [22], using
352 same data and method alike Korber et al. (2020). A preview published by Grubaugh et al.
353 (2020) stated those major points as unclear without any concrete interpretations, and
354 questioned on the correlation of D614G mutation with virus spread or transmissibility.
355 Dearlove et al. (2020) acknowledged that this 614G variant may constitute an exception to
356 their conclusion on the arising of other mutations of the virus due to genetic drift ⁷². Rausch
357 et al. (2020) ⁷³ also stressed upon the possible effect of D614G on viral fitness through
358 revisiting the points raised by Dearlove et al. (2020). Recently, Plante et al. (2020) ⁷⁴ and Hou
359 et al. (2020) ⁷⁵ reported that D614G mutation alters virus fitness by enhancing viral load in
360 COVID-19 patients and may increase transmission in human, which was speculated from the
361 increasing transmission of mutant-type viruses in hamsters.

362 S: p.D614G mutation may facilitate the exposure of the cleavage domains of S1-S2
363 and S2' to proteases, elastase-2, Furin, or TMPRSS2, as discussed by Easwarkhanth, et al. ⁷⁶.
364 Another explanation of the more efficient cellular entry for D614G variant is the possibly
365 breakage of hydrogen bonds that modulates the interactions between S protein protomers and
366 ACE2 receptor binding ^{17,75,77}. Moreover, G614 mutation may increase S protein stability and
367 participate in N-linked glycosylation at N616 ⁴¹. Several recent experiments have suggested
368 that mutated (G₆₁₄) protein contains a novel serine protease cleavage site at 615-616 that is
369 cleaved by host elastase-2, a potent neutrophil elastase playing important roles in
370 inflammatory diseases of human, more efficiently than wild protein, D₆₁₄
371 ^{25,41,58,78}. Bhattacharyya et al. (2020) hypothesized that the neutrophil elastase level at the site
372 of infection will facilitate the host cell entry for G₆₁₄. We explained here, with the aid of
373 structural bioinformatics, how the mutated protein gets the advantage over wild type with a
374 single mutation.

375 This study found interesting structural features of the S protein while comparing and
376 superimposing the wild (D₆₁₄) over mutant (G₆₁₄). The secondary structure prediction and

377 surface accessibility analyses showed that there was a slight mismatch at the S1-S2 junction
378 (⁶⁸¹PRRAR↓S⁶⁸⁶) where serine at 686 (S⁶⁸⁶) was found covered in G₆₁₄ and exposed to the
379 surface in D₆₁₄. However, S⁶⁸⁶ in both G₆₁₄ and D₆₁₄ were exposed to an open-loop region to
380 have possible contact with the proteases (supplementary Figure S1). Further investigation on
381 the aligned 3D structures showed no conformational change at the Furin or TMPRSS2
382 cleavage site (Figure 3c). We observed no structural variation in the surrounding residues of
383 the protease-targeting S1-S2 site (Figure 3c), which eliminated the assumption of Phan⁷⁹.
384 The predictive 3D models and structural assessment of D₆₁₄ and G₆₁₄ variants also confirmed
385 that the cleavage site at 815-16 of S2 subunit (⁸¹²PSKR↓S⁸¹⁶) or S2'^{3,80} had no structural and
386 surface topological variation (Figure 3d-e). Rather, the superimposed 3D structures suggested
387 a conformational change in the immediate downstream region (⁶¹⁸TEVPVAIHADQLTPT⁶³²)
388 of the 614th position of G₆₁₄ that was not observed in D₆₁₄ variants (Figure 3a-b).

389 The elastase-2 restrictedly cut valine at 615, due to its valine-dependent constriction
390 of catalytic groove⁸¹. The present sequence setting surrounding of G₆₁₄ (P6-
391 ⁶¹⁰VLYQGV↓NCTEV⁶²⁰-P'5) showed a higher enzymatic activity over the D₆₁₄, which
392 cannot be completely aligned with previous works on the sequence-based substrate
393 specificity of elastase-2⁸². However, the first misaligned residue of the superimposed G₆₁₄,
394 located at the P'4 position (T⁶¹⁸), may also be important for binding with the elastase-2, and
395 further down the threonine (T) at 618, the residues may affect the bonding with the respective
396 amino acids of the protease. This changed conformation at the downstream binding site of
397 G₆₁₄ may help overcome unfavorable adjacent sequence motifs in the mutated S protein as the
398 elastase-2 substrate. Therefore, the simultaneous and/or sequential processing of the mutated
399 S protein by TMPRSS2/Furin/Cathepsin and elastase-2 facilitates a more efficient SARS-
400 CoV-2 entry into the host cells and cell-cell fusion^{41,58,78}.

401 This study further observed the possible association of the S protein with elastase-2
402 and found an increased binding affinity in case of G₆₁₄ (Table 3). Hence, the active sites of
403 the mutated protein interacted efficiently with more amino acids of elastase-2 (Table 4),
404 possibly providing a better catalytic activity as shown by Hu, et al.⁴¹. The mutation may have
405 changed the structural configuration of the elastase-2 cleavage site in a way that the enzyme
406 is facing less challenge to get near to the cutting site of the altered protein (Figure 3a-b and
407 4). The efficient cleaving of this enzyme, although located in an upstream position of the S1-
408 S2 junction, may assist in releasing S1 from S2 and change the conformation in a way to get
409 later cleaved at the S2' site, and then help in the fusion process^{83,84}. Mutated spike protein

410 and elastase-2 complex was more flexible than the wild one, and the interactions with
411 enzyme was also different as shown in RMSD deviation between the complexes (Figure 5).

412 Besides, this G⁶¹⁴ amino acid replacement may have a destabilization effect on the
413 overall protein structure (Table 4 and Figure 3a-b). The deformed flexible region at or near
414 G⁶¹⁴ is the proof of that destabilizing change (Figure 3f and supplementary figure S3). Zhang
415 et al.⁸⁵ explained less S1 shedding through more stable hydrogen bonding between Q⁶¹³
416 (glutamine) and T⁸⁵⁹ (threonine) of protomer due to greater backbone flexibility provided by
417 G⁶¹⁴. On the other hand, increasing the number of RBD up conformation or increasing the
418 chance of 1-RBD-up conformation due to breakage of both intra-and inter-protomer
419 interactions of the spike trimer and symmetric conformation will give a better chance to bind
420 with ACE2 receptor and can also increase antibody-mediated neutralization⁸⁶. The S1 will
421 release from S2 more effectively in G614 protein due to introduction of glycine that will
422 break hydrogen bond in between the D614 and the neighboring protomer T859 amino acid
423^{86,87}. Our analyses have provided the *in silico* proof of this point by showing that the mutated
424 protein was less stable than the wild type protein by missing a hydrophobic interaction with
425 Phe⁵⁹² (Figure 3g-h). This new result accorded well with the later consensus as explained in
426 Weismann et al. (2020) that G614 will increase its overall flexibility. Therefore, the overall
427 structural change may assist the mutated S protein by providing elastase-2 a better binding
428 space and attachment opportunity onto the cleavage site, thus providing a more stable
429 interaction with the S2 domain that increases the credibility of an efficient infection (Figure
430 1).

431 **3.3 Increased Flexibility of RdRp-NSP8 Complex: Compromise Proof-Reading** 432 **Efficiency with Replication Speed**

433 The binding free energy (ΔG) of the RdRp-NSP8 complexes have been predicted as -
434 10.6 and -10.5 Kcalmol⁻¹, respectively, in wild (P₃₂₃) and mutant (L₃₂₃) type that suggests less
435 strong interaction for mutant protein (Table 4). The number of contacts made at the interface
436 (IC) per property and interacting amino acids increased between L₃₂₃ and NSP8 (Table 3 and
437 4). This increasing contacts due to higher hydrogen bonds (Figure 5d) made the complex of
438 mutated RdRp and NSP8 less flexible. Our analyses have identified that proline (P³²³) or
439 leucine (L³²³) of RdRp can interact with the aspartic acid (D¹¹²), cysteine (C¹¹⁴), valine (V¹¹⁵),
440 and proline (P¹¹⁶) of NSP8 (Figure 6). RdRp binds with NSP8 in its interface domain
441 (residues alanine:A²⁵⁰ to arginine:R³⁶⁵), forming positively charged 'sliding poles' for RNA

442 exit and enhance the replication speed probably by extending the RNA-binding surface on
443 NSP8^{88,89}. Molecular dynamics of the mutated RdRp-NSP8 complex supported this by
444 showing a more expanded surface area in interacting sites (Figure 5b) and maintained
445 integrity throughout the simulation (Figure 5c). Another study reported that the NSP8 binding
446 sites on RdRp and the RNA exit tunnel were comparatively neutral and conserved⁹⁰. Besides,
447 a zinc ion in the conserved metal-binding motif (H²⁹⁵, C³⁰¹, C³⁰⁶, and C³¹⁰), close to the 323
448 residue, is responsible for maintaining the integrity of the RdRp architecture^{90,91}. We did not
449 find any interaction of NSP8 with the zinc-binding residues of RdRp proteins (Table 3 and
450 Figure 6). Therefore, the P323L mutation within this conserved site of the RdRp interface
451 domain may only affect the RdRp-NSP8 interaction without changing metal binding affinity.

452 Recent reports on structural analyses revealed that the RdRp:p.P323L led to the
453 stabilization of the protein structure^{92,93}. The results from the six state-of-the-art tools used in
454 this study have suggested that mutant (L₃₂₃) protein cannot be concluded as 'stable' from the
455 ambiguous $\Delta\Delta G$ values (Table 4), rather the interaction with the adjacent amino acid⁹⁴
456 mainly defined the stability. Although the hydrophobic amino acid (L³²³) embedded in the
457 L₃₂₃ protein, our secondary structure analysis indicated the presence of only loop structure in
458 323 position for both wild and mutant types (Figure S1). However, Chand et al. (2020)
459 reported the turn structures of 323 and 324 shifted into five sheets at positions 321, 322, 323,
460 324 and 327 to bury leucine. The superimposed 3D structures solved this contradiction about
461 secondary structure by showing that there was no deviation in loop/turn structure of mutant
462 protein (supplementary Figure S4).

463 Overall, the mutation at 323 position, to some extent, stabilized the L₃₂₃ structure,
464 made the protein more rigid, binds less strongly with the NSP8, and thus expanded the
465 interacting region with NSP8. These variations may together increase the replication speed by
466 exiting the processed RNA genome from the RdRp groove structure more swiftly (Figure 1).
467 The increasing replication speed might be due to the perturbation of interaction between
468 RdRp and NSP8^{88,94}, or less possibly, the complex tripartite interactions (RdRp, NSP8, and
469 NSP14) responsible for the speculated decrease of proof-reading efficiency⁶. Thus, RdRp
470 mutants might increase the mutation rate by a trade-off between high replication speed and
471 low fidelity of the mutant polymerase⁹⁵. Another possibility could be the lower proof-reading
472 efficiency of NSP14 that was not linked to the replication speed⁶. Analysis of our study
473 sequences revealed that the frequency of mutation (median=8) in L₃₂₃ mutants (n=27,364) is
474 significantly higher (p<0.0001) than the frequency (median=6) of wild-type (P₃₂₃) strains

475 (n=9,815). This increased mutation rate may play a vital role in genetic drifts and provide
476 next generations a better adaptation to adverse environments.

477 **3.4 Q57H Substitution in ORF3a Viroporin: the Roles of Decreased Ion Permeability**

478 This study has found that the replacement of glutamine (Q⁵⁷) with positively charged
479 histidine (H⁵⁷) at 57 position of ORF3a transmembrane region 1 (TM1) does not change
480 secondary transmembrane helical configuration (supplementary Figure S1), and aligned 3D
481 structures have also shown no variation of TM1 in the monomeric state (Figure 7a). The
482 mutant (H₅₇) protein has a non-significant increase in structural stabilization and a minimal
483 decrease in molecular flexibility (Table 4 and supplementary Figure S5). This is because of
484 the weak ionic interaction of C_α with the sulfur atom of cysteine (C⁸¹) in TM2 and hydrogen
485 bond of terminal N_ε of lysine (61K) with one of the endocyclic nitrogens of H⁵⁷ (Figure 7b).

486 Selection analysis predicted the accumulation of those mutations as a result of
487 pervasive positive selection with an increasing trend ⁹⁶. One possibility of this positive
488 selection of histidine over glutamine might be the role of H⁵⁷ (TM1) in increasing
489 constrictions wherein the diagonal C⁸¹ (TM2) may assist through the ionic interaction (Figure
490 7b). Notably, the Q⁵⁷ in wild-type protein forms the major hydrophilic constriction within the
491 ORF3a channel pore ⁹⁷. Thus, further favorable increasing constrictions within the H₅₇ protein
492 channel pore and the replacement of charge-neutral Q with a positive-charged H in the
493 selectivity filter may reduce the passing of positive ions such as Ca²⁺, Na⁺, and K⁺ by either
494 electrostatic repulsion or blocking ⁹⁸⁻¹⁰¹. This speculation for ORF3a mutant protein was
495 supported by another study that showed the reduction of ion permeability of Na²⁺ and Ca²⁺
496 through the H₅₇, however, that decrease was not found statistically significant (p>0.05) ⁹⁷.

497 The decrease intracellular concentration of cytoplasmic Ca²⁺ ions potentially reduces
498 caspase-dependent apoptosis of the host cell ¹⁰², mainly supporting viral spread without
499 affecting replication ²³ as shown in Figure 1. Moreover, the ORF3a can drive necrotic cell
500 death ¹⁰³ wherein the permeated ions ¹⁰⁴ and the insertion of ORF3a viroporin into lysosome
501 ¹⁰⁵ play vital roles. The H₅₇ mutant may thus decrease pathogenicity and symptoms during the
502 early stages of the infection, i.e., reducing ‘cytokine storm’ in the host ¹⁰⁶. Besides, ORF3a
503 was proved to affect inflammasome activation, viral release, and cell death, as shown by
504 Castaño-Rodríguez, et al. ¹⁰⁷ that the deletion of ORF3a reduced viral load and morbidity in
505 animal models.

506 Even though similar proteins of ORF3a have been identified in the sarbecovirus
507 lineage infecting bats, pangolins, and humans¹⁰⁸, only one pangolin derived strain from 2017
508 in Guangxi, China contains H⁵⁷ residue as shown by our mutation analyses (supplementary
509 Figure S6) and also reported by Kern, et al.⁹⁷. However, bat or civet did not contain this
510 mutation in TM1 and the flanking vicinity was not identical; whereas, the TM1
511 (⁴¹LPGWLIVGVALLAVFQSASKII⁶³) does not have amino acid replacements for the
512 strains of SARS-CoV-2 and pangolin-coronaviruses (Figure S5a-b). The presence of this
513 mutation in pangolin could be an accidental case or might explain its impact on modulating
514 host-specific immune response, which needs functional experimental verification. A possible
515 explanation behind that presence might be the more accustomed nature of the virus towards
516 reverse transmission by being less virulent, i.e., from human to other animals, as observed in
517 recent reports^{109,110}.

518 **3.5 N Protein Mutation: Augmenting Nucleocapsid Stability and Exerting** 519 **Miscellaneous Effects**

520 Our study has observed that the combined mutation (N: p.RG203-204KR) causes no
521 conformational change in secondary and 3D structures (Figure S1 and Figure 8, respectively)
522 of the conserved SR-rich region of the LKR (supplementary Figure S7), but there is a minor
523 alteration in the degree of buried or exposed site (Figure 3). This result contradicted the
524 prediction of¹¹¹ about the change in the length and arrangements of the alpha-helix in the SR-
525 rich region. The superimposed 3D structures showed structural deviation, rather at
526 ²³¹ESKMSGKGQQQGGQTVT²⁴⁷ of the LKR (Figure 9), corresponding to the high
527 destabilization of the KR₂₀₃₋₂₀₄ protein (Table 3). On the other hand, A220V mutation in the
528 N protein of the GV clade showed a slightly more stable formation of the mutated N protein
529 with no change in the chemical properties (Table 4), that might affect RNA binding affinity
530¹¹².

531 Impedance to form particular SR-motif due to RG→KR mutation might disrupt the
532 phosphorylation catalyzed by host glycogen synthase kinase-3¹¹³. Similar hypo-
533 phosphorylation events could arise due to the conversion of serine to nonpolar or neutral
534 amino acids (L^{188/194/197}, I¹⁹³, and N²⁰²), as represented in supplementary Table S2 and¹¹⁴.
535 Consequently, the low phosphorylation level after entering into the cell should unwind the
536 viral ribonucleoprotein (RNP) in a slower but more organized fashion that might have an
537 impact upon translation and immune-modulation¹¹⁵⁻¹¹⁷. In KR₂₀₃₋₂₀₄, replacement of a glycine
538 with lysine that may increase the nucleocapsid (N protein- RNA complex) stability by

539 forming stronger electrostatic and ionic interactions due to increased positive charge^{118,119}.
540 Besides, the more disordered orientation of the associated LKR¹¹⁵ and highly destabilizing
541 property of KR₂₀₃₋₂₀₄ may assist in the packaging of a stable RNP^{112,120}. These interactions
542 and impact upon mutations are depicted in Figure 1.

543 N protein also utilizes the dynamic nature of the intrinsically disordered linker region
544 (LKR) that controls its affinity towards M protein, self-monomer, 5'UTR, and cellular
545 proteins¹²¹⁻¹²³. The phosphorylation at the LKR site may play an essential role to regulate
546 these interactions¹¹⁸. It was speculated that KR²⁰³⁻²⁰⁴ attained more selective advantage⁹⁶
547 over the other mutations of N protein (Table S2), probably because of stronger RNA binding
548 and synchronized hypo-phosphorylation.

549 **3.6 Silent Mutations may not be Silent**

550 The C241T of 5'UTR (untranslated region), a single nucleotide change, or 'silent'
551 mutation, located at the UUCGU pentaloop part of the stem-loop region (SLR5B) has a
552 potential role in viral packaging¹²⁴. This pentaloop of 5'UTR remains unchanged and
553 maintains a particular structure^{122,125}. The RNA secondary structural analysis in our study
554 predicted that there is no change in the 241T structure (supplementary Figure S8a). However,
555 the silent mutation in the loop region upstream to the ORF1a start codon (266-268 position)
556 may be involved in differential RNA binding affinity to the ribosome and translational factors
557¹²⁶.

558 In the case of multi-domain NSP3 (papain-like protease), we have observed superior
559 stability of the RNA after gaining of the synonymous mutation 3037C<T (C318T) where
560 wild and mutant RNA structure has -151.63 and -153.03 Kcal/mol, respectively
561 (Supplementary Figure S8b-c). A more stable secondary structure of (+)-ssRNA as observed
562 in the mutant NSP3 protein corresponds to the slower translational elongation. It contributes
563 to a range of abnormalities resulting in low translation efficiency, which affects
564 posttranslational modifications as a part of protein regulation¹²⁷. Because silent mutations
565 have an impact on the ribosome occupancy time depending on the structural modification and
566 mRNA stability which guide co-translational folding kinetics of a protein¹²⁷. This silent
567 mutation is located within the flexible loop of the NSP3 ubiquitin-like domain 1 (Ubl1). In
568 SARS-CoV, Ubl1 was reported to bind with single-stranded RNA containing AUA patterns
569 and interact with the nucleocapsid (N) protein^{128,129}. Besides, Ubl1 was likely to bind with
570 several signature repeats in 5'-UTR in SARS-CoV-2 genome¹³⁰. Figure 1 represents the
571 overall possible scenario due to these silent mutations.

572 Change in T445C in leader protein may not cause any change in expression or others
573 since the structure and energy are same -172.34 kcal/mol. The change C6286T is at between
574 in the NSP3 region whereas the C26801G is in between the Nucleic acid-binding (NAB)
575 domain and betacoronavirus specific marker (β SM) domain. The change C26801G is in the
576 TM3 of virion membrane.

577 **3.7 Epistatic Effects of Co-occurring Mutations: Increasing Viral Fitness Costing** 578 **Virulence?**

579 The co-occurring mutations, as defined by the presence of simultaneous multi-site
580 variations in the same or different proteins or in the genome [115], have provided new
581 insights into the dynamic epistatic network by employing differential molecular interactions.
582 The synergistic effects of the mutations were reported to control viral fitness and virulence,
583 as observed in Influenza and Ebola virus outbreak¹³¹⁻¹³⁶. For instance, the detrimental effect
584 of R384G on influenza A fitness was overcome by the co-occurring mutation E375G^{132,135},
585 and co-occurring mutations at the antigenic sites of influenza hemagglutinin can also drive
586 viral evolution^{134,136}. The correlation of the co-occurring GP-L mutations affect Ebola virus
587 virulence and thus case fatality rate¹³³. The viral fitness determined by the efficiency of
588 viruses to spread throughout the population and infect more new victims^{137,138}. A general
589 definition of virulence as proposed by Geoghegan et al. (2018) is ‘the harm caused by
590 pathogen infection, particularly in terms of host morbidity and mortality’¹³⁹. Evolution
591 generally drives a pathogen to gain more fitness while reducing the virulence¹³⁹ during the
592 natural course of mutations¹⁴⁰; nevertheless, there is contradiction between short-term¹⁴⁰ and
593 long-term evolution¹³⁹ of virus. However, the virulence can also increase during the course
594 of the outbreak as observed in HIV in 2002¹⁴¹. Here we used death-case ratio and the
595 number of infection cases as equivalent to virulence and fitness, respectively.

596 Between two important G clade-featured co-occurring mutations, the p:D614G of the
597 S protein and p:P323L of RdRp, we observed no interlinked functional relationship (Figure
598 1). The former assists mainly by more rapid entering into the host cells with an efficient
599 elastase-2 activity and higher aggressiveness of G₆₁₄ mutant is related to elastase-2 or human
600 neutrophil elastase (hNE) concentration during inflammation¹⁴². Besides, higher levels of
601 functional S protein observed in G₆₁₄ strains can increase the chance of host-to-host
602 transmission⁸⁵. The later one may boost up the replication by a faster RNA processing
603 (exiting) that can open up the avenue to generate strains with significantly ($p < 0.0001$) higher
604 number of mutations. This increased mutation rate in L₃₂₃ mutants can surpass the constant

605 proof-reading fidelity¹⁴³, with an average of ~8 (range 0-45) mutations per strain (n=37,179),
606 and evolve a greater number of variants in a population (Figure 1). Nevertheless, SARS-
607 CoV-2 has still low genetic diversity, and these mutations mostly are the product of genetic
608 drift with no or unknown effect on virus virulence and fitness [53, 54], and without
609 possessing a deleterious effect on viral fitness, a mutation fixed early on in the epidemic is
610 unlikely to be lost.

611 The S protein and RdRp mutations, albeit seemingly unrelated, can cumulatively
612 escalate the infectiousness of the virus as a result of higher viral load and shorter burst time.
613 This rapid within-host replication might be directly correlated with the virulence, in turn,
614 morbidity and mortality rate¹⁴⁴. This study also showed that co-occurring mutations, by
615 acting together, may benefit viral populations by incrementing the ability to produce a diverse
616 viral population (Figure 1), that might be able to adapt more quickly in adverse climatic
617 condition, evade the immune response, and survive within different selective pressure^{145,146}.
618 The S protein mutation is also related to human allelic variation⁵⁸, and may get an advantage
619 in particular populations.

620 NSP3 is a scaffolding protein for the replication-transcription complex, and the
621 possible change in its structure may affect the overall dynamics of viral replication^{128,129}.
622 P323L mutation of RdRp may change binding affinity to the Ubl1 region of NSP3¹⁴⁷ (Figure
623 1). The p.C241T mutation, on its own, may affect the transcription and viral packaging,
624 although we have not found any possible association of it with the S, RdRp, and NSP3
625 mutant proteins.

626 The mutant N protein may have an impact on viral replication and transcription, like
627 other coronaviruses¹²³, through the binding with NSP3 protein that is linked to RdRp
628 centered replication complex. This N protein can affect the membrane stability by yet
629 uncharacterized interaction with the M protein, which should ultimately produce more stable
630 virion particles¹⁴⁸⁻¹⁵⁰. A stronger N protein-RNA complex provokes slower intracellular
631 immune response¹¹⁶, and at the same time, can remain highly contagious and aggressive
632 because of the concurrent presence of G clade-featured S protein and RdRp mutations (Figure
633 1). In the case of mutant ORF3a protein, we have not found any report that correlates with
634 other co-occurring mutations. H₅₇ mutant, possibly linked to the mild or asymptomatic cases,
635 may allow the silent transmission and increase the chance of viral spread by lowering the

636 activation of inflammatory response (Figure 1), such as reduced viral particle release and
637 cytokine storm^{23,151,152}.

638 The fittest viral strains will dominate in a population considering other selective
639 parameters associated with the virulence^{139,144}. Public health interventions were able to
640 create a selective pressure to make a strain less virulent and highly competitive^{153,154}, which
641 might increase the chance of viral transmission¹⁵⁵; however, there was also argument against
642 this established theory¹⁵⁶. The SARS-CoV-2 with multiple clades and variants needs to be
643 more efficient to maintain the delicate evolutionary trade-off between fitness and virulence.
644 A recent publication on COVID-19 demonstrated a reduced viral load in the upper respiratory
645 system for a particular genotype¹⁵⁷, which may in turn influence viral transmission and
646 spread. We proposed here a hypothesis (model) on SARS-CoV-2 infection as to how the
647 fitness of the major clade strains might cost virulence.

648 The presence of mostly G strains in the early pandemic period (starting weeks of
649 March) might be linked to a higher mortality rate (Figure 2a). G strains were continuously
650 present in any region at a particular level (~20% among these clades), which could be the
651 reason for maintaining a balanced high mortality rate. This stage of the pandemic might be
652 linked to the initial lack of awareness, hygienic practice, and social distancing that led to a
653 large number of transmissions. The high viral load and quick immune response will spread
654 the G strains more efficiently in a crowded area as shown in the case of Africa (Figure 1)
655 where the lockdown and other social interventions were not very strict¹⁵⁸. The G strains thus
656 spread unprecedentedly from Europe to the other parts of the world. Being more
657 transmissible alone can also cause problem by putting a greater strain on hospitals and
658 leading to a sharper spike in deaths.

659 The GH strains were mainly restricted to the USA and partly Eastern Mediterranean
660 (Fig. 1), and mostly spread by cryptic¹⁵⁹ and pre/a-symptomatic transmission¹⁶⁰. These
661 variants might trade-off virulence by a slower release of virions, and in exchange, benefited
662 from the induction of low immune response in asymptomatic hosts. They were more fit at
663 that time, in theory, when the people were dealing with the pandemic in panic. The GH type
664 might be able to maximize their transmission by residing within the host unknowingly and
665 spread at ease. We can not nullify that there could be no significant link of the mutation with
666 virulence or death-case ratio since the effect of ORF3a was not found significant in wet lab
667 experiment and in our statistical analysis as well.

668 The GR strains possibly attained an advantage over G and GH by a more orchestrated,
669 delicately balanced synergistic effects on replication and transmission fitness. These epistatic
670 effects could have increased the fitness by hiding the virus from immunity and increasing
671 stability in the environment, i.e., more transmissible through air and surface. The
672 asymptomatic patient infected with GR strains, similar to GH, would have a weaker immune
673 response and shed the virus for a longer period ¹⁵². Hamed et al. (2020) ⁷⁰ reported a result
674 against this hypothesis with a statistical approach using GISAID data; however only 4.75%
675 (27,608/ 581,120 on February 20, 2021) sequence data contains patient status, that will
676 definitely misinterpret true scenario.

677 GV strains featuring an A222V mutation in the S protein with other mutations of the
678 clade G was reported to have a less pronounced effect due to its structural position and have
679 no role in antibody binding, thus founder effect might be the cause of its spreading within
680 Europe and the age group bias towards young adults ²⁴. There was hence probably no effect
681 of the mutations on the viral transmission, severity, and escaping antibody ^{24,161}, rather
682 superspreading events after lifting up of travel restriction in Europe and lack of effective
683 containment could be the main factors for the rapid spread and establishment of this clade.
684 There was, however, speculation about the effect of this mutation on immunogenicity since
685 computational analysis predicted its location within T-cell target region ^{162,163}. A plausible
686 link of the mutation with viral re-infection in some patients was of prime concern and
687 demands further study ¹⁶⁴. How the more stabilized linker region of the mutated N protein
688 due to A220V mutation might give an edge for the GV strains in terms of association with
689 lower death-case ratio but seemingly higher infection number by interacting with RNA is a
690 question for further study. Together, the mutations of GV strains can affect in an unknown
691 way in transmission and spread that cannot be captured alone in terms of spike A222V
692 mutation.

693 Lineage B.1.1.7, a variant under GR clade, consists of 17 mutations (14
694 nonsynonymous and 3 deletions) in spike, N, ORF1ab, and ORF8 proteins, as well as, 6
695 silent mutations ^{165,166}. This lineage shows a 50% more transmissibility, realistic possibility in
696 increased risk of death, and is rapidly dominating over other variants in more than 75
697 countries. The mutations on the spike protein might play most crucial role both in increasing
698 transmission fitness and a slightly reduced neutralization to antibody ^{167,168}. Other mutations
699 may still play roles in viral replication and pathogenesis through different epistatic
700 interactions. B.1.351 lineage of clade GH also showed a faster spread and higher load of virus

701 in swabs, as well as, higher binding affinity to ACE2 receptor and escape from neutralizing
702 antibody due to three mutations (K417N, E484K and N501Y) at key sites of the RBD ¹⁶⁹.
703 This lineage has mutations on the E, N, and ORF1a proteins, functions of which are yet to be
704 determined ¹⁷⁰. Similar important mutational package in RBD with known biological
705 importance to B.1.351 lineage, except for K417T and several others in the ORF1a, N, ORF8,
706 and ORF3a, was found in P.1 (descendent of B.1.1.28) lineage of GR clade ^{171,172}. B.1.258
707 lineage under G clade contained N439K that has higher binding affinity with new contact
708 point to ACE2 and can evade antibody-mediated immunity^{173,174}. The mutations in proteins
709 other than spike of these new variants could have effects that needs further investigation.
710 Recently, the variants are supplanting their original clade strains throughout the world due to
711 RBD mutation effects, which makes correlating the clade with case fatality rate difficult and
712 evolutionary theory complex. Due to these variants, recent data of January suggests an
713 increase of GR strains where GH and GV are similar in different regions and worldwide in a
714 weekly basis (data not shown). However, the death-case ratio is following the previous trends
715 of around 2% in the weeks after 3rd January.

716 The GH and GR strains have seemingly arisen during the initial phase of the pandemic
717 with a limited frequency, and at that time, have not attained the necessary viral fitness. Both
718 strains have evolved from Europe and Australia after spreading quietly alongside the
719 vigorous G strain. After the public awareness and other social interventions, these strains,
720 especially GR have probably increased in fitness. Although Hu and Riley ¹⁵⁹ have recently
721 reported the co-circulation of both G and GH strains correlating with mortality in different
722 states, how the association among multiple clades would modulate the evolutionary dynamics
723 is now a burning question. With the new data of GV strains as shown in our study, the
724 previous analysis by Hodcroft et al. (2020) ²⁴ can be updated to recent time that may change
725 result of the association. Our analysis has considered the evolutionary trade-off between
726 virulence and sustainable fitness of these strains based on the global emergence of the strains,
727 their death-case ratio, structural stability, and predicted molecular mechanisms. The relation
728 of the host-pathogen interactions, host allelic variations, and host immunity needs to correlate
729 with virulence by in vivo studies.

730 **4. CONCLUSION**

731 The course of COVID-19 pandemic was continuing, although the death rate was
732 gradually decreasing in 2020. Our study hypothesized that this paradoxical scenario was

733 related to the effect of the dominant co-occurring mutations. The significant association of
734 death-case ratio with GR clade mutant could be linked with the signature mutation in
735 nucleocapsid protein. GH clade mutant might contribute to cryptic transmission by ORF3a-
736 mediated low immune response in asymptomatic hosts. GV strains have spread quickly by
737 superspreading events, but we speculated a role of more stable linker region of N protein on
738 low virulence. G clade mutants, however, was speculated to assist in frequent transmissions
739 through severe symptoms lead to the isolation of infected persons restricting further
740 spreading and significantly increased the risk of death by a factor of 1.03. Therefore, it can be
741 speculated that the fitness of SAR-CoV-2 virus may increase in terms of survivability and
742 transmission, while the next phase of the COVID-19 pandemic may continue with less
743 virulent mutants. However, further investigations are required to strengthen this hypothesis.
744 Effective management of the infected patients and social distancing to prevent viral
745 transmission should be the first priority that may prevent the emergence of more virulent
746 strains.

747 **Data Availability**

748 All the sequence data were taken from the GISAID (<https://www.gisaid.org/>) and
749 RCSB PDB (<https://www.rcsb.org/>) as mentioned in the methodology section. We provide all
750 the necessary information such as accession numbers, date-based data source for helping
751 readers and reviewers to check the authenticity of the work.

752 **Acknowledgments**

753 We would like to acknowledge the team at GISAID for creating SARS-CoV-2 global
754 database. The funding of the research was provided by Jashore University of Science and
755 Technology. We appreciate to the Microbial Genetics and Bioinformatics Laboratory of
756 University of Dhaka for the support of the high-performance computer access. We thanked
757 M. Shaminur Rahman and Spencer Mark Mondol for their technical assistance in protein
758 structure analyses.

759 **Biographical Note:**

760 A. S. M. Rubayet Ul Alam is a lecturer in the Department of Microbiology, Jashore
761 University of Science and Technology. His research mainly focuses on molecular biology,
762 data analysis, and bioinformatics.

763 Ovinu Kibria Islam is an assistant professor in the Department of Microbiology,
764 Jashore University of Science and Technology. His research activity is focused on molecular
765 biology and big data analysis.

766 Md. Shazid Hasan is a lecturer in the Department of Microbiology, Jashore University
767 of Science and Technology. His research activity is focused on molecular biology and
768 bioinformatics.

769 Mir Raihanul Islam is the Senior Research Associate of BRAC James P Grant School
770 of Public Health, BRAC University who is expert in statistical analysis and epidemiological
771 study and is currently working on many active researches related to public health.

772 Shafi Mahmud is working as a thesis student in the department of Genetic
773 Engineering and Biotechnology, University of Rajshahi, Rajshai-6205, Bangladesh.

774 Hassan M. Al-Emran is an assistant professor in the Department of Biomedical
775 Engineering, Jashore University of Science and Technology, Jashore-7408, Bangladesh. His
776 research mainly focuses on clinical microbiology and infectious diseases.

777 Dr. Iqbal Kabir Jahid is a professor and chairman in the Department of Microbiology,
778 Jashore University of Science and Technology. His research mainly focuses on molecular
779 biology.

780 Keith A. Crandall is a professor in the Department of Biostatistics & Bioinformatics
781 and the director of Computational Biology Institute, Director of Genomics Core, Milken
782 Institute School of Public Health, The George Washington University, Washington, DC,
783 USA.

784 M. Anwar Hossain is the director of Genome Center, and vice chancellor of Jashore
785 University of Science and Technology. He is an expert in molecular biology, virology and
786 vaccine development.

787 REFERENCES

- 788 1 Lu, R. *et al.* Genomic characterisation and epidemiology of 2019 novel coronavirus:
789 implications for virus origins and receptor binding. *The Lancet* **395**, 565-574 (2020).
790 2 WHO. Coronavirus disease (COVID-19) Situation Report – 168 *WHO situation*
791 *Reports* (2020).
792 3 Hoffmann, M., Kleine-Weber, H. & Pöhlmann, S. A multibasic cleavage site in the
793 spike protein of SARS-CoV-2 is essential for infection of human lung cells.
794 *Molecular Cell* (2020).

- 795 4 Hoque, M. N., Chaudhury, A., Akanda, M. A. M., Hossain, M. A. & Islam, M. T.
796 Genomic diversity and evolution, diagnosis, prevention, and therapeutics of the
797 pandemic COVID-19 disease. *PeerJ* **8**, e9689 (2020).
- 798 5 Petersen, E. *et al.* Comparing SARS-CoV-2 with SARS-CoV and influenza
799 pandemics. *The Lancet infectious diseases* (2020).
- 800 6 Romano, M., Ruggiero, A., Squeglia, F., Maga, G. & Berisio, R. A Structural View of
801 SARS-CoV-2 RNA Replication Machinery: RNA Synthesis, Proofreading and Final
802 Capping. *Cells* **9**, 1267 (2020).
- 803 7 Rahman, M. S. *et al.* Epitope-based chimeric peptide vaccine design against S, M and
804 E proteins of SARS-CoV-2 etiologic agent of global pandemic COVID-19: an in
805 silico approach. *PeerJ* **8**, e9572 (2020).
- 806 8 Yoshimoto, F. K. The Proteins of Severe Acute Respiratory Syndrome Coronavirus-2
807 (SARS CoV-2 or n-COV19), the Cause of COVID-19. *The protein journal* **39**, 198-
808 216, doi:10.1007/s10930-020-09901-4 (2020).
- 809 9 Callaway, E. The coronavirus is mutating-does it matter? *Nature* **585**, 174-177
810 (2020).
- 811 10 Islam, M. R. *et al.* Genome-wide analysis of SARS-CoV-2 virus strains circulating
812 worldwide implicates heterogeneity. *Scientific reports* **10**, 1-9 (2020).
- 813 11 Rahman, M. S. *et al.* Comprehensive annotations of the mutational spectra of
814 SARS-CoV-2 spike protein: a fast and accurate pipeline. *Transboundary and*
815 *emerging diseases* (2020).
- 816 12 Eskier, D., Karakulah G, Suner A, Oktay Y. RdRp mutations are associated with
817 SARS-CoV-2 genome evolution. *PeerJ* **8**, 9587,
818 doi:<https://doi.org/10.7717/peerj.9587> (2020).
- 819 13 Hassan, S. S., Choudhury, P. P., Basu, P. & Jana, S. S. Molecular conservation and
820 Differential mutation on ORF3a gene in Indian SARS-CoV2 genomes. *Genomics*
821 (2020).
- 822 14 Shaminur Rahman, M. *et al.* Evolutionary dynamics of SARS-CoV-2 nucleocapsid
823 (N) protein and its consequences. *Journal of medical virology*.
- 824 15 Shu, Y. & McCauley, J. GISAID: Global initiative on sharing all influenza data—from
825 vision to reality. *Eurosurveillance* **22**, 30494 (2017).
- 826 16 Yin, C. Genotyping coronavirus SARS-CoV-2: methods and implications. *Genomics*
827 **112**, 3588-3596, doi:10.1016/j.ygeno.2020.04.016 (2020).
- 828 17 Korber, B. *et al.* Tracking changes in SARS-CoV-2 Spike: evidence that D614G
829 increases infectivity of the COVID-19 virus. *Cell* **182**, 812-827. e819 (2020).
- 830 18 Mercatelli, D. & Giorgi, F. M. Geographic and Genomic Distribution of SARS-CoV-
831 2 Mutations. *Frontiers in Microbiology* **11**, doi:10.3389/fmicb.2020.01800 (2020).
- 832 19 Brufsky, A. Distinct Viral Clades of SARS-CoV-2: Implications for Modeling of
833 Viral Spread. *Journal of medical virology* (2020).
- 834 20 Islam, O. *et al.* Emergence of European and North American mutant variants of
835 SARS-CoV-2 in Southeast Asia. (2020).
- 836 21 Ul Alam, A. R., Rafiul Islam, M., Shaminur Rahman, M., Islam, O. K. & Anwar
837 Hossain, M. Understanding the possible origin and genotyping of first Bangladeshi
838 SARS-CoV-2 strain. *Journal of Medical Virology* (2020).
- 839 22 Becerra-Flores, M. & Cardozo, T. SARS-CoV-2 viral spike G614 mutation
840 exhibits higher case fatality rate. *International Journal of Clinical Practice* (2020).
- 841 23 Issa, E., Merhi, G., Panossian, B., Salloum, T. & Tokajian, S. SARS-CoV-2 and
842 ORF3a: Nonsynonymous Mutations, Functional Domains, and Viral Pathogenesis.
843 *Msystems* **5** (2020).

- 844 24 Hodcroft, E. B. *et al.* Emergence and spread of a SARS-CoV-2 variant through
845 Europe in the summer of 2020. *medRxiv* (2020).
- 846 25 Grubaugh, N. D., Hanage, W. P. & Rasmussen, A. L. Making sense of mutation: what
847 D614G means for the COVID-19 pandemic remains unclear. *Cell* (2020).
- 848 26 Pachetti, M. *et al.* Emerging SARS-CoV-2 mutation hot spots include a novel RNA-
849 dependent-RNA polymerase variant. *Journal of Translational Medicine* **18**, 1-9
850 (2020).
- 851 27 of the International, C. S. G. The species Severe acute respiratory syndrome-related
852 coronavirus: classifying 2019-nCoV and naming it SARS-CoV-2. *Nature*
853 *Microbiology* **5**, 536 (2020).
- 854 28 Dalgaard, P. in *Introductory statistics with R* 259-274 (Springer, 2008).
- 855 29 Wang, W., Tang, J. & Wei, F. Updated understanding of the outbreak of 2019 novel
856 coronavirus (2019-nCoV) in Wuhan, China. *Journal of medical virology* **92**, 441-447
857 (2020).
- 858 30 WHO. Coronavirus disease (COVID-19) Weekly Epidemiological Update and
859 Weekly Operational Update. *WHO situation Reports* (2020).
- 860 31 Rodrigues, C. H., Pires, D. E. & Ascher, D. B. DynaMut: predicting the impact of
861 mutations on protein conformation, flexibility and stability. *Nucleic acids research*
862 **46**, W350-W355 (2018).
- 863 32 Schymkowitz, J. *et al.* The FoldX web server: an online force field. *Nucleic acids*
864 *research* **33**, W382-W388 (2005).
- 865 33 Delgado, J., Radusky, L. G., Cianferoni, D. & Serrano, L. FoldX 5.0: working with
866 RNA, small molecules and a new graphical interface. *Bioinformatics* **35**, 4168-4169
867 (2019).
- 868 34 Rost, B., Yachdav, G. & Liu, J. The predictprotein server. *Nucleic acids research* **32**,
869 W321-W326 (2004).
- 870 35 Waterhouse, A. *et al.* SWISS-MODEL: homology modelling of protein structures and
871 complexes. *Nucleic acids research* **46**, W296-W303 (2018).
- 872 36 Webb, B. & Sali, A. Comparative protein structure modeling using MODELLER.
873 *Current protocols in bioinformatics* **54**, 5.6. 1-5.6. 37 (2016).
- 874 37 Roy, A., Kucukural, A. & Zhang, Y. I-TASSER: a unified platform for automated
875 protein structure and function prediction. *Nature protocols* **5**, 725-738 (2010).
- 876 38 Wang, Q. *et al.* Structural basis for RNA replication by the SARS-CoV-2 polymerase.
877 *Cell* **182**, 417-428. e413 (2020).
- 878 39 Van Zundert, G. *et al.* The HADDOCK2. 2 web server: user-friendly integrative
879 modeling of biomolecular complexes. *Journal of molecular biology* **428**, 720-725
880 (2016).
- 881 40 Xue, L. C., Rodrigues, J. P., Kastritis, P. L., Bonvin, A. M. & Vangone, A.
882 PRODIGY: a web server for predicting the binding affinity of protein-protein
883 complexes. *Bioinformatics* **32**, 3676-3678, doi:10.1093/bioinformatics/btw514
884 (2016).
- 885 41 Hu, J. *et al.* The D614G mutation of SARS-CoV-2 spike protein enhances viral
886 infectivity. *bioRxiv* (2020).
- 887 42 de Vries, S. J. & Bonvin, A. M. CPORT: a consensus interface predictor and its
888 performance in prediction-driven docking with HADDOCK. *PloS one* **6**, e17695
889 (2011).
- 890 43 Yan, Y., Tao, H., He, J. & Huang, S.-Y. The HDock server for integrated protein-
891 protein docking. *Nature protocols* **15**, 1829-1852 (2020).
- 892 44 Land, H. & Humble, M. S. in *Protein Engineering* 43-67 (Springer, 2018).

- 893 45 Dickson, C. J. *et al.* Lipid14: the amber lipid force field. *Journal of chemical theory*
894 *and computation* **10**, 865-879 (2014).
- 895 46 Krieger, E. & Vriend, G. New ways to boost molecular dynamics simulations.
896 *Journal of computational chemistry* **36**, 996-1007 (2015).
- 897 47 Krieger, E., Nielsen, J. E., Spronk, C. A. & Vriend, G. Fast empirical pKa prediction
898 by Ewald summation. *Journal of molecular graphics and modelling* **25**, 481-486
899 (2006).
- 900 48 Chowdhury, K. H. *et al.* Drug Repurposing Approach against Novel Coronavirus
901 Disease (COVID-19) through Virtual Screening Targeting SARS-CoV-2 Main
902 Protease. *Biology* **10**, 2 (2021).
- 903 49 Pramanik, S. K. *et al.* Fermentation optimization of cellulase production from
904 sugarcane bagasse by *Bacillus pseudomycolides* and molecular modeling study of
905 cellulase. *Current Research in Microbial Sciences* **2**, 100013 (2021).
- 906 50 Swargiary, A., Mahmud, S. & Saleh, M. A. Screening of phytochemicals as potent
907 inhibitor of 3-chymotrypsin and papain-like proteases of SARS-CoV2: an in silico
908 approach to combat COVID-19. *Journal of Biomolecular Structure and Dynamics*, 1-
909 15 (2020).
- 910 51 Yan, H. & Yu, C. Repair of full-thickness cartilage defects with cells of different
911 origin in a rabbit model. *Arthroscopy: The Journal of Arthroscopic & Related Surgery*
912 **23**, 178-187 (2007).
- 913 52 Zuker, M. Mfold web server for nucleic acid folding and hybridization prediction.
914 *Nucleic acids research* **31**, 3406-3415 (2003).
- 915 53 Huston, N. C., Wan, H., Tavares, R. d. C. A., Wilen, C. & Pyle, A. M.
916 Comprehensive in-vivo secondary structure of the SARS-CoV-2 genome reveals
917 novel regulatory motifs and mechanisms. *BioRxiv* (2020).
- 918 54 Rane, J. S., Chatterjee, A., Kumar, A. & Ray, S. (ChemRxiv, 2020).
- 919 55 Souza, P. F., Lopes, F. E., Amaral, J. L., Freitas, C. D. & Oliveira, J. T. A molecular
920 docking study revealed that synthetic peptides induced conformational changes in the
921 structure of SARS-CoV-2 spike glycoprotein, disrupting the interaction with human
922 ACE2 receptor. *International journal of biological macromolecules* **164**, 66-76
923 (2020).
- 924 56 Elfiky, A. A. Ribavirin, Remdesivir, Sofosbuvir, Galidesivir, and Tenofovir against
925 SARS-CoV-2 RNA dependent RNA polymerase (RdRp): A molecular docking study.
926 *Life sciences*, 117592 (2020).
- 927 57 Aftab, S. O. *et al.* Analysis of SARS-CoV-2 RNA-dependent RNA polymerase as a
928 potential therapeutic drug target using a computational approach. *Journal of*
929 *translational medicine* **18**, 1-15 (2020).
- 930 58 Bhattacharyya, C. *et al.* SARS-CoV-2 mutation 614G creates an elastase cleavage site
931 enhancing its spread in high AAT-deficient regions. *Infection, Genetics and*
932 *Evolution*, 104760 (2021).
- 933 59 Backhaus, A. Common Pitfalls in the Interpretation of COVID-19 Data and Statistics.
934 *Intereconomics* **55**, 162-166 (2020).
- 935 60 Sáez, C., Romero, N., Conejero, J. A. & García-Gómez, J. M. Potential limitations in
936 COVID-19 machine learning due to data source variability: A case study in the
937 nCov2019 dataset. *Journal of the American Medical Informatics Association* **28**, 360-
938 364 (2021).
- 939 61 L.R. Lopes Junior, S. A., M. Qazvini COVID-19: Statistics, Ratios and understanding
940 limitations of data available. *ICAT: Data Limitations sub-group* (2020).
- 941 62 Zeberg, H. & Pääbo, S. The major genetic risk factor for severe COVID-19 is
942 inherited from Neanderthals. *Nature* **587**, 610-612 (2020).

- 943 63 Zeberg, H. & Pääbo, S. A genomic region associated with protection against severe
944 COVID-19 is inherited from Neandertals. *Proceedings of the National Academy of*
945 *Sciences* **118** (2021).
- 946 64 Jain, V. K., Iyengar, K., Vaish, A. & Vaishya, R. Differential mortality in COVID-19
947 patients from India and western countries. *Diabetes & Metabolic Syndrome: Clinical*
948 *Research & Reviews* **14**, 1037-1041 (2020).
- 949 65 Goel, S. *et al.* Clinical characteristics and in-hospital mortality for COVID-19 across
950 the Globe. *Cardiology and Therapy* **9**, 553-559 (2020).
- 951 66 Samaddar, A., Gadepalli, R., Nag, V. L. & Misra, S. The enigma of low COVID-19
952 fatality rate in India. *Frontiers in Genetics* **11**, 854 (2020).
- 953 67 Villani, L., McKee, M., Giraldi, L., Ricciardi, W. & Boccia, S. Comparison of deaths
954 rates for COVID-19 across Europe. (2020).
- 955 68 Njenga, M. K. *et al.* Why is there low morbidity and mortality of COVID-19 in
956 Africa? *The American journal of tropical medicine and hygiene* **103**, 564-569 (2020).
- 957 69 Omais, S. O., Kharroubi, S. O. & Zaraket, H. No association between the SARS-CoV-
958 2 variants and mortality rates in the Eastern Mediterranean Region. *medRxiv* (2021).
- 959 70 Hamed, S. M., Elkhatib, W. F., Khairallah, A. S. & Noreddin, A. M. Global dynamics
960 of SARS-CoV-2 clades and their relation to COVID-19 epidemiology. (2020).
- 961 71 McKee, M., Gugushvili, A., Koltai, J. & Stuckler, D. Are populist leaders creating the
962 conditions for the spread of COVID-19?; Comment on “A scoping review of populist
963 radical right parties’ influence on welfare policy and its implications for population
964 health in Europe”. *International journal of health policy and management* (2020).
- 965 72 Dearlove, B. *et al.* A SARS-CoV-2 vaccine candidate would likely match all currently
966 circulating variants. *Proceedings of the National Academy of Sciences* **117**, 23652-
967 23662 (2020).
- 968 73 Rausch, J. W., Capoferri, A. A., Katusime, M. G., Patro, S. C. & Kearney, M. F. Low
969 genetic diversity may be an Achilles heel of SARS-CoV-2. *Proceedings of the*
970 *National Academy of Sciences* **117**, 24614-24616 (2020).
- 971 74 Plante, J. A. *et al.* Spike mutation D614G alters SARS-CoV-2 fitness. *Nature*, 1-9
972 (2020).
- 973 75 Hou, Y. J. *et al.* SARS-CoV-2 D614G variant exhibits efficient replication *ex vivo*
974 and transmission *in vivo*. *Science* (2020).
- 975 76 Eaaswarkhanth, M., Al Madhoun, A. & Al-Mulla, F. Could the D614 G substitution
976 in the SARS-CoV-2 spike (S) protein be associated with higher COVID-19 mortality?
977 *International Journal of Infectious Diseases* (2020).
- 978 77 Omotuyi, I. O. *et al.* Atomistic simulation reveals structural mechanisms underlying
979 D614G spike glycoprotein-enhanced fitness in SARS-CoV-2. *Journal of*
980 *computational chemistry* **41**, 2158-2161 (2020).
- 981 78 Korber, B. *et al.* Tracking changes in SARS-CoV-2 Spike: evidence that D614G
982 increases infectivity of the COVID-19 virus. *Cell* (2020).
- 983 79 Phan, T. Genetic diversity and evolution of SARS-CoV-2. *Infection, genetics and*
984 *evolution* **81**, 104260 (2020).
- 985 80 Belouzard, S., Chu, V. C. & Whittaker, G. R. Activation of the SARS coronavirus
986 spike protein via sequential proteolytic cleavage at two distinct sites. *Proceedings of*
987 *the National Academy of Sciences* **106**, 5871-5876 (2009).
- 988 81 Perona, J. J. & Craik, C. S. Structural basis of substrate specificity in the serine
989 proteases. *Protein Science* **4**, 337-360 (1995).
- 990 82 Fu, Z., Thorpe, M., Akula, S., Chahal, G. & Hellman, L. T. Extended cleavage
991 specificity of human neutrophil elastase, human proteinase 3, and their distant

- 992 ortholog clawed frog PR3—three elastases with similar primary but different
993 extended specificities and stability. *Frontiers in Immunology* **9**, 2387 (2018).
- 994 83 Li, F. Structure, function, and evolution of coronavirus spike proteins. *Annual review*
995 *of virology* **3**, 237-261 (2016).
- 996 84 Walls, A. C. *et al.* Tectonic conformational changes of a coronavirus spike
997 glycoprotein promote membrane fusion. *Proceedings of the National Academy of*
998 *Sciences* **114**, 11157-11162 (2017).
- 999 85 Zhang, L. *et al.* SARS-CoV-2 spike-protein D614G mutation increases virion spike
1000 density and infectivity. *Nature communications* **11**, 1-9 (2020).
- 1001 86 Weissman, D. *et al.* D614G spike mutation increases SARS CoV-2 susceptibility to
1002 neutralization. *Cell host & microbe* **29**, 23-31. e24 (2021).
- 1003 87 Wrapp, D. *et al.* Cryo-EM structure of the 2019-nCoV spike in the prefusion
1004 conformation. *Science* **367**, 1260-1263 (2020).
- 1005 88 Hillen, H. S. *et al.* Structure of replicating SARS-CoV-2 polymerase. *Nature*, 1-6
1006 (2020).
- 1007 89 Yin, W. *et al.* Structural basis for inhibition of the RNA-dependent RNA polymerase
1008 from SARS-CoV-2 by remdesivir. *Science* (2020).
- 1009 90 Kirchdoerfer, R. N. & Ward, A. B. Structure of the SARS-CoV nsp12 polymerase
1010 bound to nsp7 and nsp8 co-factors. *Nature communications* **10**, 1-9 (2019).
- 1011 91 Gao, Y. *et al.* Structure of the RNA-dependent RNA polymerase from COVID-19
1012 virus. *Science* **368**, 779-782 (2020).
- 1013 92 Chand, G., Banerjee A, Azad GK. Identification of novel mutations in RNA-
1014 dependent RNA polymerases of SARS-CoV-2 and their implications on its protein
1015 structure. *PeerJ* **8**, e9492, doi:<https://doi.org/10.7717/peerj.9492> (2020).
- 1016 93 Begum, F., Mukherjee, D., Das, S., Thagriki, D., Tripathi, P. P., Banerjee, A. K., &
1017 Ray, U. Specific mutations in SARS-CoV2 RNA dependent RNA polymerase and
1018 helicase alter protein structure, dynamics and thus function: Effect on viral RNA
1019 replication. *bioRxiv.org* (2020).
- 1020 94 Chand, G. B., Banerjee, A. & Azad, G. K. Identification of novel mutations in RNA-
1021 dependent RNA polymerases of SARS-CoV-2 and their implications on its protein
1022 structure. *bioRxiv* (2020).
- 1023 95 Eskier, D., Karakulah, G., Suner, A. & Oktay, Y. RdRp mutations are associated with
1024 SARS-CoV-2 genome evolution. *bioRxiv* (2020).
- 1025 96 Pond, S. Genomic diversity and divergence of SARS-CoV-2/COVID-19 from
1026 GISAID. (2020).
- 1027 97 Kern, D. M. *et al.* Cryo-EM structure of the SARS-CoV-2 3a ion channel in lipid
1028 nanodiscs. *BioRxiv* (2020).
- 1029 98 Malasics, A. *et al.* Protein structure and ionic selectivity in calcium channels:
1030 Selectivity filter size, not shape, matters. *Biochimica et Biophysica Acta (BBA)-*
1031 *Biomembranes* **1788**, 2471-2480 (2009).
- 1032 99 Naranjo, D., Moldenhauer, H., Pincuntureo, M. & Díaz-Franulic, I. Pore size matters
1033 for potassium channel conductance. *Journal of General Physiology* **148**, 277-291
1034 (2016).
- 1035 100 Stephens, R. F., Guan, W., Zhorov, B. S. & Spafford, J. D. Selectivity filters and
1036 cysteine-rich extracellular loops in voltage-gated sodium, calcium, and NALCN
1037 channels. *Frontiers in Physiology* **6**, 153 (2015).
- 1038 101 Suárez-Delgado, E. & Islas, L. D. Ion Channels: A novel origin for calcium
1039 selectivity. *Elife* **9**, e55216 (2020).

- 1040 102 Kondratskyi, A., Kondratska, K., Skryma, R. & Prevarskaya, N. Ion channels in the
1041 regulation of apoptosis. *Biochimica et Biophysica Acta (BBA)-Biomembranes* **1848**,
1042 2532-2546 (2015).
- 1043 103 Yue, Y. *et al.* SARS-Coronavirus Open Reading Frame-3a drives multimodal necrotic
1044 cell death. *Cell death & disease* **9**, 1-15 (2018).
- 1045 104 Pinton, P., Giorgi, C., Siviero, R. & Zecchini, E. Rizzuto R. *Bcl-2 and Ca* **2**, 1409-
1046 1418 (2006).
- 1047 105 Nieva, J. L., Madan, V. & Carrasco, L. Viroporins: structure and biological functions.
1048 *Nature Reviews Microbiology* **10**, 563-574 (2012).
- 1049 106 Ren, Y. *et al.* The ORF3a protein of SARS-CoV-2 induces apoptosis in cells. *Cellular*
1050 *& molecular immunology*, 1-3 (2020).
- 1051 107 Castaño-Rodríguez, C. *et al.* Role of severe acute respiratory syndrome coronavirus
1052 viroporins E, 3a, and 8a in replication and pathogenesis. *MBio* **9** (2018).
- 1053 108 Boni, M. F. & Lemey, P. Evolutionary origins of the SARS-CoV-2 sarbecovirus
1054 lineage responsible for the COVID-19 pandemic. *Nature microbiology*,
1055 doi:10.1038/s41564-020-0771-4 (2020).
- 1056 109 Halfmann, P. J. *et al.* Transmission of SARS-CoV-2 in domestic cats. *New England*
1057 *Journal of Medicine* (2020).
- 1058 110 Shi, J. *et al.* Susceptibility of ferrets, cats, dogs, and other domesticated animals to
1059 SARS-coronavirus 2. *Science* **368**, 1016-1020 (2020).
- 1060 111 Ayub, M. I. Reporting Two SARS-CoV-2 Strains Based on A Unique Trinucleotide-
1061 Bloc Mutation and Their Potential Pathogenic Difference. *Preprints*
1062 doi:10.20944/preprints202004.0337.v1 (2020).
- 1063 112 Chang, C.-k., Hou, M.-H., Chang, C.-F., Hsiao, C.-D. & Huang, T.-h. The SARS
1064 coronavirus nucleocapsid protein-forms and functions. *Antiviral research* **103**, 39-50
1065 (2014).
- 1066 113 Tylor, S. *et al.* The SR-rich motif in SARS-CoV nucleocapsid protein is important for
1067 virus replication. *Canadian journal of microbiology* **55**, 254-260 (2009).
- 1068 114 Rahman, M. S. *et al.* Evolutionary dynamics of SARS-CoV-2 nucleocapsid protein
1069 and its consequences. *Journal of medical virology* (2020).
- 1070 115 Järvelin, A. I., Noerenberg, M., Davis, I. & Castello, A. The new (dis) order in RNA
1071 regulation. *Cell Communication and Signaling* **14**, 9 (2016).
- 1072 116 Kikkert, M. Innate immune evasion by human respiratory RNA viruses. *Journal of*
1073 *innate immunity* **12**, 4-20 (2020).
- 1074 117 Kopecky-Bromberg, S. A., Martínez-Sobrido, L., Frieman, M., Baric, R. A. & Palese,
1075 P. Severe acute respiratory syndrome coronavirus open reading frame (ORF) 3b, ORF
1076 6, and nucleocapsid proteins function as interferon antagonists. *Journal of virology*
1077 **81**, 548-557 (2007).
- 1078 118 McBride, R., Van Zyl, M. & Fielding, B. C. The coronavirus nucleocapsid is a
1079 multifunctional protein. *Viruses* **6**, 2991-3018 (2014).
- 1080 119 Sokalingam, S., Raghunathan, G., Soundrarajan, N. & Lee, S.-G. A study on the
1081 effect of surface lysine to arginine mutagenesis on protein stability and structure using
1082 green fluorescent protein. *PloS one* **7**, e40410 (2012).
- 1083 120 Haynes, C. & Iakoucheva, L. M. Serine/arginine-rich splicing factors belong to a class
1084 of intrinsically disordered proteins. *Nucleic acids research* **34**, 305-312 (2006).
- 1085 121 Carlson, C. R., Asfaha, J. B., Ghent, C. M., Howard, C. J., Hartooni, N., Safari, M., ...
1086 & Morgan, D. O. Phosphoregulation of Phase Separation by the SARS-CoV-2 N
1087 Protein Suggests a Biophysical Basis for its Dual Functions. *Molecular Cell* **11**, 025,
1088 doi:<https://doi.org/10.1016/j.molcel.2020.11.025>. (2020).

- 1089 122 Schuster, N. A. Using the nucleocapsid protein to investigate the relationship between
1090 SARS-CoV-2 and closely related bat and pangolin coronaviruses. *BioRxiv* (2020).
- 1091 123 de Haan, C. A. & Rottier, P. J. Molecular interactions in the assembly of
1092 coronaviruses. *Advances in virus research* **64**, 165-230 (2005).
- 1093 124 Rangan, R. *et al.* RNA genome conservation and secondary structure in SARS-CoV-2
1094 and SARS-related viruses: a first look. *Rna* **26**, 937-959 (2020).
- 1095 125 Huston, N. C. *et al.* Comprehensive in vivo secondary structure of the SARS-CoV-2
1096 genome reveals novel regulatory motifs and mechanisms. *Molecular cell* **81**, 584-598.
1097 e585 (2021).
- 1098 126 Kristofich, J. *et al.* Synonymous mutations make dramatic contributions to fitness
1099 when growth is limited by a weak-link enzyme. *PLoS genetics* **14**, e1007615 (2018).
- 1100 127 Mitra, S., Ray, S. K. & Banerjee, R. Synonymous codons influencing gene expression
1101 in organisms. *Research and Reports in Biochemistry* **6**, 57 (2016).
- 1102 128 Hurst, K. R., Ye, R., Goebel, S. J., Jayaraman, P. & Masters, P. S. An interaction
1103 between the nucleocapsid protein and a component of the replicase-transcriptase
1104 complex is crucial for the infectivity of coronavirus genomic RNA. *Journal of*
1105 *virology* **84**, 10276-10288 (2010).
- 1106 129 Hurst, K. R., Koetzner, C. A. & Masters, P. S. Characterization of a critical
1107 interaction between the coronavirus nucleocapsid protein and nonstructural protein 3
1108 of the viral replicase-transcriptase complex. *Journal of virology* **87**, 9159-9172
1109 (2013).
- 1110 130 Serrano, P. *et al.* Nuclear magnetic resonance structure of the N-terminal domain of
1111 nonstructural protein 3 from the severe acute respiratory syndrome coronavirus.
1112 *Journal of virology* **81**, 12049-12060 (2007).
- 1113 131 Du, X. *et al.* Networks of genomic co-occurrence capture characteristics of human
1114 influenza A (H3N2) evolution. *Genome research* **18**, 178-187 (2008).
- 1115 132 Rimmelzwaan, G. *et al.* Full restoration of viral fitness by multiple compensatory co-
1116 mutations in the nucleoprotein of influenza A virus cytotoxic T-lymphocyte escape
1117 mutants. *Journal of general virology* **86**, 1801-1805 (2005).
- 1118 133 Deng, L. *et al.* Network of co-mutations in Ebola virus genome predicts the disease
1119 lethality. *Cell research* **25**, 753-756 (2015).
- 1120 134 Chen, H., Zhou, X., Zheng, J. & Kwoh, C.-K. Rules of co-occurring mutations
1121 characterize the antigenic evolution of human influenza A/H3N2, A/H1N1 and B
1122 viruses. *BMC medical genomics* **9**, 69 (2016).
- 1123 135 Rimmelzwaan, G. F., Kreijtz, J. H., Bodewes, R., Fouchier, R. A. & Osterhaus, A. D.
1124 Influenza virus CTL epitopes, remarkably conserved and remarkably variable.
1125 *Vaccine* **27**, 6363-6365 (2009).
- 1126 136 Lyons, D. M. & Luring, A. S. Mutation and epistasis in influenza virus evolution.
1127 *Viruses* **10**, 407 (2018).
- 1128 137 Crow, J. F. *Basic concepts in population, quantitative, and evolutionary genetics.*
1129 (WH Freeman and Company, 1986).
- 1130 138 Di Giallonardo, F. & Holmes, E. C. Viral biocontrol: grand experiments in disease
1131 emergence and evolution. *Trends in microbiology* **23**, 83-90 (2015).
- 1132 139 Geoghegan, J. L. & Holmes, E. C. The phylogenomics of evolving virus virulence.
1133 *Nature Reviews Genetics* **19**, 756-769 (2018).
- 1134 140 Berngruber, T. W., Froissart, R., Choisy, M. & Gandon, S. Evolution of virulence in
1135 emerging epidemics. *PLoS Pathog* **9**, e1003209 (2013).
- 1136 141 Pantazis, N. *et al.* Temporal trends in prognostic markers of HIV-1 virulence and
1137 transmissibility: an observational cohort study. *The Lancet HIV* **1**, e119-e126 (2014).
- 1138 142 Mohamed, M. M., El-Shimy, I. A. & Hadi, M. A. (BioMed Central, 2020).

- 1139 143 Smith, E. C., Sexton, N. R. & Denison, M. R. Thinking outside the triangle:
1140 replication fidelity of the largest RNA viruses. *Annual Review of Virology* **1**, 111-132
1141 (2014).
- 1142 144 Skjesol, A. *et al.* IPNV with high and low virulence: host immune responses and viral
1143 mutations during infection. *Virology journal* **8**, 1-14 (2011).
- 1144 145 Pfeiffer, J. K. & Kirkegaard, K. Increased fidelity reduces poliovirus fitness and
1145 virulence under selective pressure in mice. *PLoS Pathog* **1**, e11 (2005).
- 1146 146 Vignuzzi, M., Stone, J. K., Arnold, J. J., Cameron, C. E. & Andino, R. Quasispecies
1147 diversity determines pathogenesis through cooperative interactions in a viral
1148 population. *Nature* **439**, 344-348 (2006).
- 1149 147 Lei, J., Kusov, Y. & Hilgenfeld, R. Nsp3 of coronaviruses: Structures and functions of
1150 a large multi-domain protein. *Antiviral research* **149**, 58-74 (2018).
- 1151 148 Schoeman, D. & Fielding, B. C. Coronavirus envelope protein: current knowledge.
1152 *Virology journal* **16**, 1-22 (2019).
- 1153 149 Escors, D., Ortego, J., Laude, H. & Enjuanes, L. The membrane M protein carboxy
1154 terminus binds to transmissible gastroenteritis coronavirus core and contributes to
1155 core stability. *Journal of virology* **75**, 1312-1324 (2001).
- 1156 150 J Alsaadi, E. A. & Jones, I. M. Membrane binding proteins of coronaviruses. *Future*
1157 *Virology* **14**, 275-286 (2019).
- 1158 151 Bai, Y. *et al.* Presumed asymptomatic carrier transmission of COVID-19. *Jama* **323**,
1159 1406-1407 (2020).
- 1160 152 Quan-Xin, L. *et al.* Clinical and immunological assessment of asymptomatic SARS-
1161 CoV-2 infections. *Nature medicine* (2020).
- 1162 153 Dennis, C. (Nature Publishing Group, 2001).
- 1163 154 Dieckmann, U., Metz, J. A. & Sabelis, M. W. *Adaptive dynamics of infectious*
1164 *diseases: in pursuit of virulence management*. Vol. 2 (Cambridge University Press,
1165 2005).
- 1166 155 Day, T., Gandon, S., Lion, S. & Otto, S. P. On the evolutionary epidemiology of
1167 SARS-CoV-2. *Current Biology* (2020).
- 1168 156 Ebert, D. & Bull, J. J. Challenging the trade-off model for the evolution of virulence:
1169 is virulence management feasible? *Trends in microbiology* **11**, 15-20 (2003).
- 1170 157 Lorenzo-Redondo, R. *et al.* A clade of SARS-CoV-2 viruses associated with lower
1171 viral loads in patient upper airways. *EBioMedicine* **62**, 103112 (2020).
- 1172 158 Abbas, K. *et al.* Routine childhood immunisation during the COVID-19 pandemic in
1173 Africa: a benefit–risk analysis of health benefits versus excess risk of SARS-CoV-2
1174 infection. *The Lancet Global Health* (2020).
- 1175 159 Hu, Y. & Riley, L. W. Dissemination and co-circulation of SARS-CoV2 subclades
1176 exhibiting enhanced transmission associated with increased mortality in Western
1177 Europe and the United States. *medRxiv* (2020).
- 1178 160 He, X. *et al.* Temporal dynamics in viral shedding and transmissibility of COVID-19.
1179 *Nature medicine* **26**, 672-675 (2020).
- 1180 161 McCallum, M. *et al.* N-terminal domain antigenic mapping reveals a site of
1181 vulnerability for SARS-CoV-2. *bioRxiv* (2021).
- 1182 162 Zhang, B.-z. *et al.* Mining of epitopes on spike protein of SARS-CoV-2 from
1183 COVID-19 patients. *Cell research* **30**, 702-704 (2020).
- 1184 163 Mateus, J. *et al.* Selective and cross-reactive SARS-CoV-2 T cell epitopes in
1185 unexposed humans. *Science* **370**, 89-94 (2020).
- 1186 164 To, K. K.-W. *et al.* Coronavirus disease 2019 (COVID-19) re-infection by a
1187 phylogenetically distinct severe acute respiratory syndrome coronavirus 2 strain
1188 confirmed by whole genome sequencing. *Clinical Infectious Diseases* (2020).

- 1189 165 Volz, E. *et al.* Transmission of SARS-CoV-2 Lineage B. 1.1. 7 in England: Insights
1190 from linking epidemiological and genetic data. *medRxiv*, 2020.2012. 2030.20249034
1191 (2021).
- 1192 166 Muik, A. *et al.* Neutralization of SARS-CoV-2 lineage B. 1.1. 7 pseudovirus by
1193 BNT162b2 vaccine-elicited human sera. *Science* (2021).
- 1194 167 McCarthy, K. R. *et al.* Recurrent deletions in the SARS-CoV-2 spike glycoprotein
1195 drive antibody escape. *Science* (2021).
- 1196 168 Kemp, S. *et al.* Neutralising antibodies drive Spike mediated SARS-CoV-2 evasion
1197 (medRxiv). *bioRxiv* (2020).
- 1198 169 Tegally, H. *et al.* Emergence and rapid spread of a new severe acute respiratory
1199 syndrome-related coronavirus 2 (SARS-CoV-2) lineage with multiple spike mutations
1200 in South Africa. *medRxiv* (2020).
- 1201 170 grinch. B.1.351_South African lineage defined by new variant 501Y.V2 - A more
1202 detailed description of the lineage. *global report investigating novel coronavirus*
1203 *haplotypes* (2021).
- 1204 171 grinch. P.1_Brazilian lineage with variants of biological significance E484K, N501Y
1205 and K417T. *global report investigating novel coronavirus haplotypes* (2021).
- 1206 172 Faria, N. R. *et al.* Genomic characterisation of an emergent SARS-CoV-2 lineage in
1207 Manaus: preliminary findings. *Virological* (2021).
- 1208 173 Zhou, W. *et al.* N439K variant in spike protein may alter the infection efficiency and
1209 antigenicity of SARS-CoV-2 based on molecular dynamics simulation. *bioRxiv*
1210 (2020).
- 1211 174 Thomson, E. C. *et al.* Circulating SARS-CoV-2 spike N439K variants maintain
1212 fitness while evading antibody-mediated immunity. *Cell* (2021).
- 1213
- 1214
- 1215
- 1216
- 1217
- 1218

1219 **Figure 1. Schematic diagram of SARS-CoV-2 replication in cell showcasing the related**
1220 **to S, N, ORF3a, RdRp, NSP3 and 5'-UTR based epistatic interactions.**

1221 The replication cycle starts with the ACE2 receptor binding of the spike glycoprotein (S) as
1222 cornered at the top-left and finishes with the exocytosis at the top-right. The viruses which do
1223 not carry G-, GH- and/or GR-featured mutations in the S, N, ORF3a, RdRp, NSP3 and 5'-
1224 UTR are denoted as the wild type where mutants contain those. Throughout the diagram, the
1225 red and green color icons such as proteins, genome, and virion represent the wild and mutant
1226 type, respectively. For a generalized virion, we used the blue color. Although this theme is

1227 not show the co-infection of both types, which might occur in rare occasions, we showed the
1228 comparative epistatic effects side-by-side fashion during the whole replication cycle that will
1229 make it easy to grasp. Related figure(s) for each protein are shown in the enclosed box. To
1230 mean the uncertainty or unknown effects of any mutant proteins/RNA structure, we used the
1231 ‘question’ mark in a pathway and on explanatory box. RdRp- RNA dependent RNA
1232 polymerase; NSP14- proof-reading enzyme of SARS-CoV-2; ER- endoplasmic reticulum;
1233 ERGIC- endoplasmic reticulum (ER) Golgi intermediate compartment.

1234 **Figure 2. Global epidemiological scenario of G, GH and GR clades.**

1235 (a) Weekly time-plot showing the percentage of G, GHGR and GV clade viruses with
1236 worldwide SARS-CoV-2 infections, death cases and death rates from 24 December, 2019
1237 (collection date of first sequenced virus) to 19th October 2020 where the weeks are shown in
1238 X axis. As low as 9 (for 2nd week) to 9587 sequences (for 14th Week) were analyzed to
1239 determine the ratio of G, GH and GR clade (shown in right Y axis) in the sequenced viruses
1240 based on the availability of high coverage genomes submitted in GISAID until 28th of
1241 October, 2020. For analyzing the death rate for each week, the number of detected infections
1242 and death cases (left Y axis) were collected from WHO situation reports.

1243 (b) Monthly region wise time-plot showing the percentage of G, GH and GR clades with
1244 infections, deaths and death rates from January to October of 2020 in 6 WHO regions. 68,125
1245 sequences were analyzed for Europe, 30,395 for Americas, 14778 for West Pacific, 3050 for
1246 Southeast Asia, 1170 for Eastern Mediterranean and 1523 for Africa. Infections and death
1247 cases were collected from WHO situation reports. The data labeling and values in the X and
1248 Y axis are same to figure 2a.

1249 **Figure 3. Different structural and stability comparison of wild and mutant spike**

1250 **protein.** Structural superposition of wild and mutant spike proteins (a-b); conformation in the
1251 S1-S2 (c) and S2'sites (d-e); representation of vibrational entropy energy change on the
1252 mutant type structure (f); and interatomic interaction prediction of both wild (g) and mutant
1253 (h) types. For Figure a-e, the gray and yellow color represent the wild and mutant protein,
1254 respectively. (a) The downstream (617-636) of D614G in wild (green) and mutant (red) S
1255 protein was focused. Overlapping of the wild (D₆₁₄) and mutant (G₆₁₄) S protein showed
1256 conformational change in the 3D structures. (b) However, the conformational change are in
1257 the loop region (618-632) of the proteins thus may potentially play role in interacting with
1258 other proteins or enzymes, such as elastase-2 as we focused in this work. (c) No change was

1259 found in the S1-S2 cleavage site (685-686), depicted in blue color, of the wild and mutant
1260 protein. **(d)** Surface and **(e)** cartoon (2°) structure of the superimposed wild and mutant
1261 proteins where the S2' (pink) is situated in surface region and do not show any change in
1262 accessibility in the residual loop region. **(f)** The mutant (G₆₁₄) protein showed higher
1263 flexibility in the G⁶¹⁴ (sticks) and its surroundings (red). The intra-molecular interaction
1264 determined the overall stability of the **(g)** wild and **(h)** mutant structure where C_β of D⁶¹⁴
1265 (aspartic acid at 614; green stick modelled) had two hydrophobic interaction with the benzene
1266 rings. This intramolecular contacts stabilize the S protein of wild type and missing of this
1267 bond destabilize the mutant (G₆₁₄) protein. The mutant protein has glycine at 614 which has
1268 less chance of interacting with other neighboring amino acids due to its shorter and nonpolar
1269 R-group. The color code representing the bond type is presented in each **(g)** and **(h)**.

1270 **Figure 4. The molecular docking of wild and mutant with elastase-2.** Both the (upper
1271 figure) wild (D₆₁₄) and (lower figure) mutant (G₆₁₄) version of S protein was shown in golden
1272 color whereas the elastase-2 docked to D₆₁₄ and G₆₁₄ in blue and green color, respectively.
1273 The enlarged views of the docked site were shown in separate boxes. **(a)** The possible docked
1274 residues (stick model) on the wild S protein (warm pink) and elastase-2 (green) are 618(Thr)-
1275 619(Glu)-620(Val) and 198(Cys)-199(Phe)-225:227 (Gly, Gly, Cys), respectively. The
1276 aspartic acid at 614 is 17.3°A far away from the valine (101), apparently the nearest amino
1277 acid of elastase-2 to the cleavage site (615-616). **(b)** The possible interacting residues (stick
1278 model) on the mutant S protein (blue) and elastase-2 (warm pink) are 614(Gly)-618(Thr)-
1279 619(Glu)-620(Val) and 101(Val)-103(Leu)-181(Arg)-222:227(Phe, Val, Arg, Gly, Gly, Cys),
1280 and 236 (Ala) respectively. In this case, the glycine at 614 is only 5.4°A far away from the
1281 valine (101), the nearest amino acid of elastase-2 to the cleavage site (615-616).

1282 **Figure 5.** (a) Both the wild and mutated spike protein had lower RMSD profile till 60ns, then
1283 it rised and maintained steady state. Although the spike protein had higher degree of
1284 deviation in RMSD profile than RdRp but they did not exceed 3.0Å. The RMSD from
1285 demonstrated that mutant and wild RdRp protein complex has initial rise of RMSD profile
1286 due to flexibility. Therefore, both RdRp complexes stabilized after 30ns and maintained
1287 steady peak. The wild type RdRp complex had little bit higher RMSD peak than mutant
1288 RdRp which indicates the more flexible nature of the wild type. (b) The spike protein
1289 complex had similar SASA profile and did not change its surface volume and maintained
1290 similar trend during the whole simulation time. The higher deviation of SASA indicates that
1291 mutant and wild type RdRp had straight line but mutant structure had higher SASA profile

1292 which indicates the protein complex had enlarged its surface area. Therefore, mutation in
1293 RdRp protein leads to more expansion of the surface area than wild types as their average
1294 SASA value had significant difference. (c) Mutated spike exhibits little more Rg profile than
1295 the wild type which correlates with the comparative labile nature of the mutant. The higher
1296 level of Rg value defines the loose packaging system and mobile nature of the protein
1297 systems. The mutant RdRp had lower level of fluctuations and maintains its integrity in
1298 whole simulation time. The wild type RdRp complexes had higher deviations and more
1299 mobility than the mutant complex. (d) Any aberration in hydrogen bond number can lead to
1300 a higher flexibility. Therefore, the mutant and wild spike protein exhibit same flexibility in
1301 terms of H-bonding. The mutant RdRp protein had more hydrogen bonding than the wild
1302 types, but they did not differentiate too much and relatively straight line was observed for the
1303 protein.

1304 **Figure 6. The molecular interaction of mutant RdRp with NSP8.** The mutant (L₃₂₃) RdRp
1305 (pale green) and NSP8 (light blue) are interacting as shown in center of the lower figure and
1306 an enlarged view of the docked site is presented above within a box. The leucine at 323
1307 interacted with the Asp (112), Cys (114), Val (115), and Pro (116). The wild (P₃₂₃) RdRp has
1308 identical docking interactions with NSP8 (table 4), thus is not presented as separate figure
1309 here.

1310 **Figure 7. The effect on transmembrane channel pore of ORF3a viroporin due to**
1311 **p.Q57H mutation.** (a) The wild (Q₅₇) and mutant (H₅₇) ORF3a protein are presented in light
1312 gray and green color, respectively. The structural superposition displays no overall
1313 conformation change, however the histidine at 57 position of mutant ORF3a (deep blue) has
1314 slightly rotated from glutamine at same position of the wild protein (bright red). This change
1315 in rotamer state at 57 residue may influence (b) the overall stability of H₅₇ (upper part) over
1316 Q₅₇ (lower part) because of ionic interaction of histidine (green; stick model) of
1317 transmembrane domain 1 (TM1) with cysteine at 81 (yellow stick) of TM2. The color code
1318 defined different bond types is shown in inset.

1319 **Figure 8. Structural superposition of wild and mutant N protein.** The light grey color
1320 represents both wild (RG₂₀₃₋₂₀₄) and mutant (KR₂₀₃₋₂₀₄) N protein. The linker region (LKR:
1321 183-247 amino acids) of wild (RG₂₀₃₋₂₀₄) and mutant (KR₂₀₃₋₂₀₄) are in pale yellow and warm
1322 pink color, respectively. (a) The aligned structures showed a highly destabilizing (Table 3)
1323 conformational change from 231 to 247 amino acids within LKR. Other regions of the N

1324 protein, especially the SR-rich region (184-204 amino acids) where the mutations occur, do
 1325 not change. **(b)** A more emphasized look into the SR-rich and mutated sites (RG203-204KR)
 1326 of wild and mutant N protein represent slight deviation in the Ser (197) and Thr (198) while
 1327 only glycine (green) to arginine (blue) substitution at 204 position shows changing at rotamer
 1328 state. The enlarged view is shown in the bottom part.

1329 **Table 1. First occurrence, global frequencies and prevalent zones for co-occurring**
 1330 **mutation were shown.**

Mutation/Clade	First occurrence (date, country, accession)	Global Frequencies (frequency, percentage)	Prevalent zone (zone, frequency, percentage)
241C>T, 3037C>T, 14408C>T, 23403A>G (G clade)	03-02-20, England EPI_ISL_464302	47943 (16.57%)	Africa (828, 30.22%)
241C>T, 3037C>T, 14408C>T, 23403A>G, 25563G>T (GH clade)	05-02-20, Australia EPI_ISL_480608	61323 (21.19%)	Americas (26647, 55.50%)
241C>T, 3037C>T, 14408C>T, 23403A>G, 28881-3GGG>AAC (GR Clade)	16-02-20, England EPI_ISL_466615	95217 (32.91%)	Western Pacific (10779, 63.50%)
241C>T, 3037C>T, 14408C>T, 23403A>G, 22227C>T (GV Clade)	15-05-20, Mexico EPI_ISL_516622	65195 (22.53%)	Europe (52864, 36.39%)

1331

1332 **Table 2. Incidence risk ratio (IRR) for different clade strains against death-case ratio.**

Variables	Unadjusted			Adjusted		
	IRR	P value	95% CI	IRR	P value	95% CI
G clade	1.03	0.020	1.01 - 1.06	1.03	0.012	1.01 - 1.06
GH clade	1.01	0.171	0.99 - 1.03	1.00	0.789	0.99 - 1.02
GR clade	0.99	0.044	0.98 - 0.99	0.98	<0.001	0.97 - 0.99
GV clade	0.98	<0.001	0.97 – 0.98	0.97	<0.001	0.96 - 0.98

1333

1334 **Table 3. The scores of HADDOCK, PRODIGY (ΔG and K_d (M) at 37.0 °C) for RdRp/NSP8**
 1335 **and Spike-Elastase docked complex.**

Variables	types	RdRp/NSP8	Spike-Elastase
HADDOCK score	Wild	-82.2 +/- 7.8	-43.0 +/- 8.9
	Mutant	-118.3 +/- 2.5	-61.9 +/- 4.5
ΔG (kcal mol⁻¹)	Wild	-10.6	-13.3
	Mutant	-10.5	-13.7
K_d (M) at 37.0 °C	Wild	3.5E ⁻⁰⁸	4.5E ⁻¹⁰
	Mutant	3.9E ⁻⁰⁸	2.3E ⁻¹⁰
Number of interfacial contacts (ICs) per property	Wild	charged-charged (5); charged-polar (9); charged-apolar (15); polar-polar (2); polar-apolar (16); and apolar-apolar (21)	charged-charged (17); charged-polar (22); charged-apolar (32); polar-polar (5); polar-apolar (31); and apolar-apolar (23)
	Mutant	charged-charged(5); charged-polar (16); charged-apolar (19); polar-polar (3); polar-apolar(15); apolar-apolar (23)	charged-charged (13); charged-polar (18); charged-apolar (27); polar-polar (4); polar-apolar (28) and apolar-apolar (36)
Associated amino acids of Elastase-2 with possible docking interactions (for spike) or NSP8 (for RdRp)	Wild	P323: Asp(112), Cys(114), Val(115) and Pro (116)	605 (Ser) and 607 (Gln): 36 (Arg); 618 (Thr): 199 (Phe); 619 (Glu): 199 (Phe), Cys (227); 620 (Val): 198 (Cys), 225:227(Gly, Gly, Cys)
	Mutant	P323: Asp(112), Cys(114), Val(115) and Pro (116)	614 (Gly): 101 (Val); 618 (Thr): 181 (Arg), 223-226(Val, Arg, Gly, Gly); 619 (Glu): 103 (Leu), 181(Arg), 222-225(Phe, Val, Arg, Gly), 236 (Ala); 620 (Val): 223-227 (Val, Arg, Gly, Gly, Cys)

1336 **Table 4. Assess the effect of mutations on structural dynamics of NSP-12/ RDRP, Spike,**

1337 **NS3 and N Protein of SARS CoV-2 using DynaMut. The value of $\Delta\Delta G$ <0 indicates that**

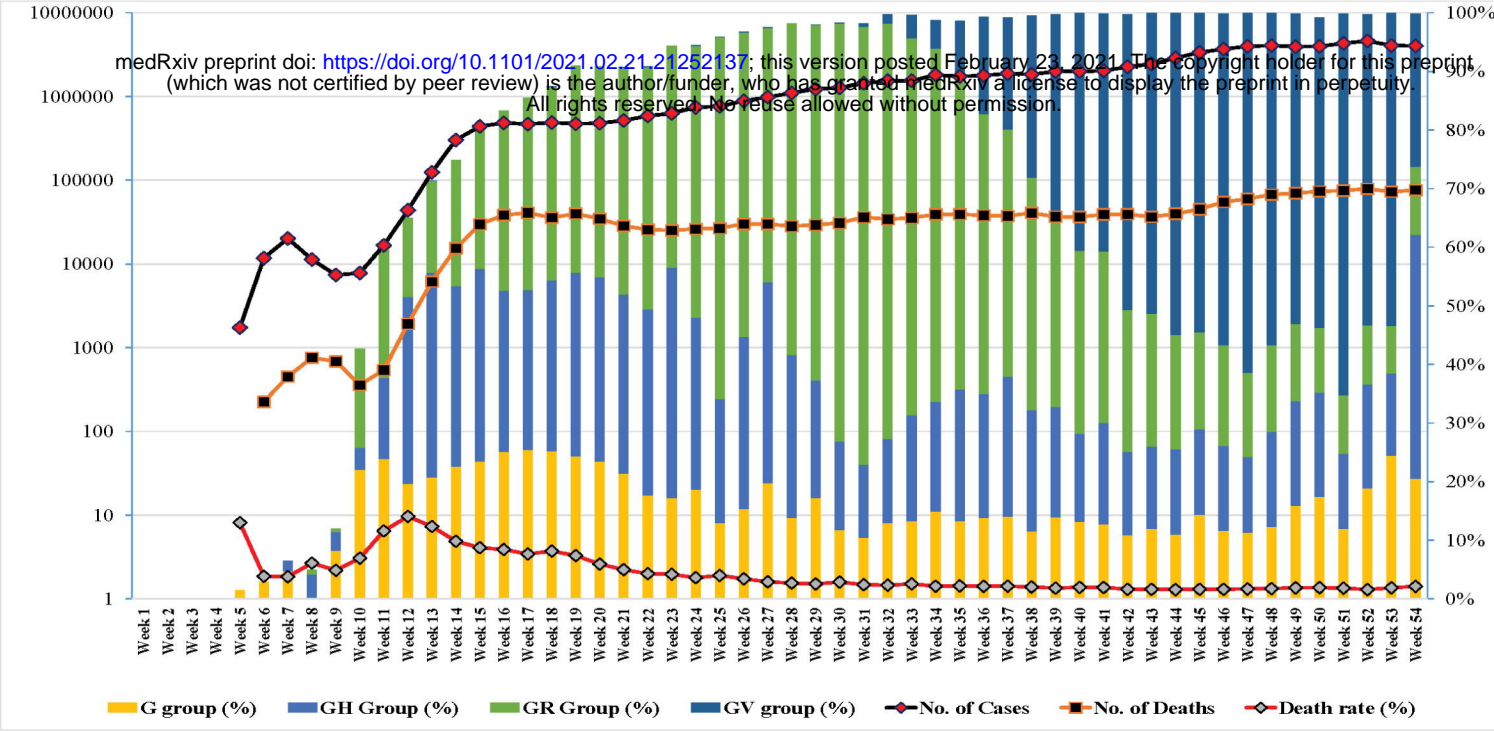
1338 **the mutation causes destabilization and $\Delta\Delta G$ > 0 represents protein stabilization. For**

1339 $\Delta\Delta S_{vib}ENCoM$, positive and negative value denotes the increase and decrease of
 1340 molecular flexibility, respectively.

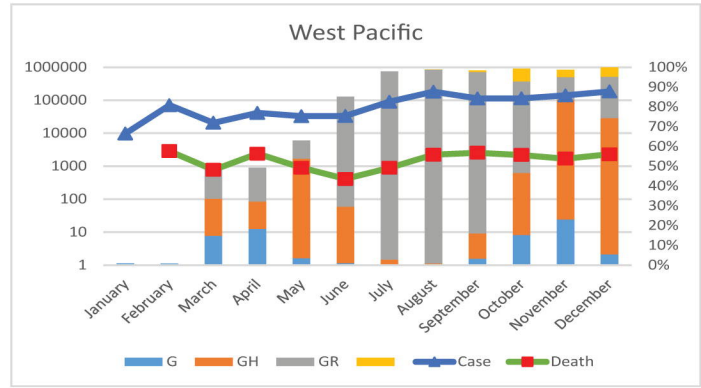
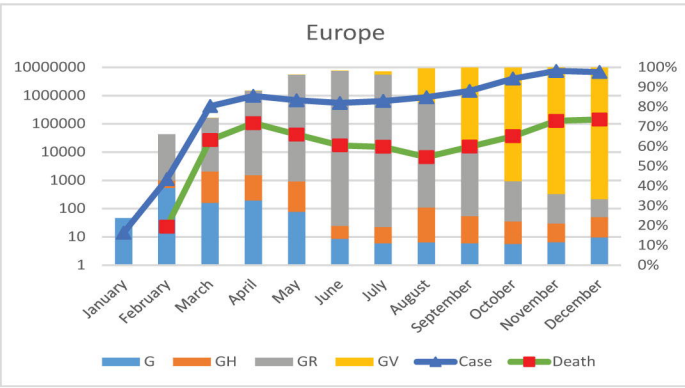
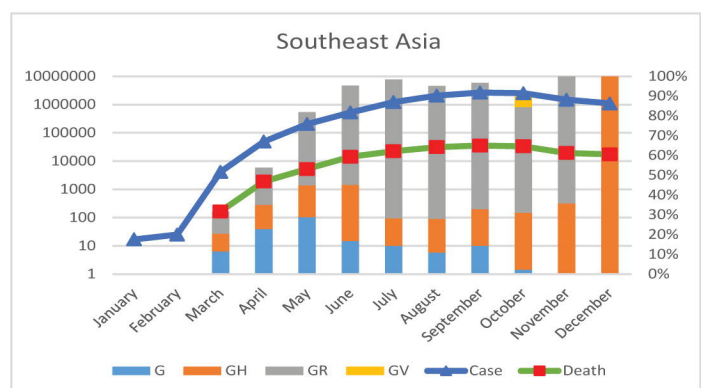
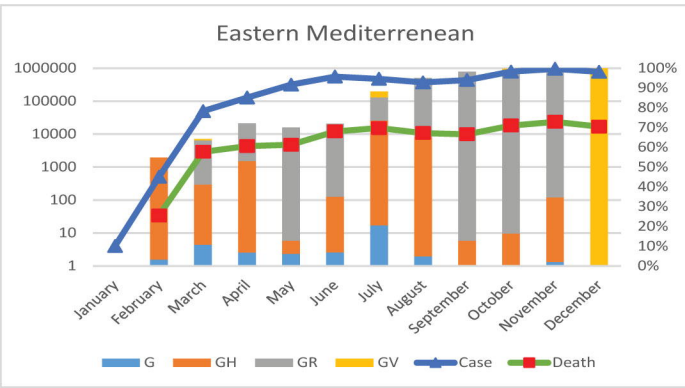
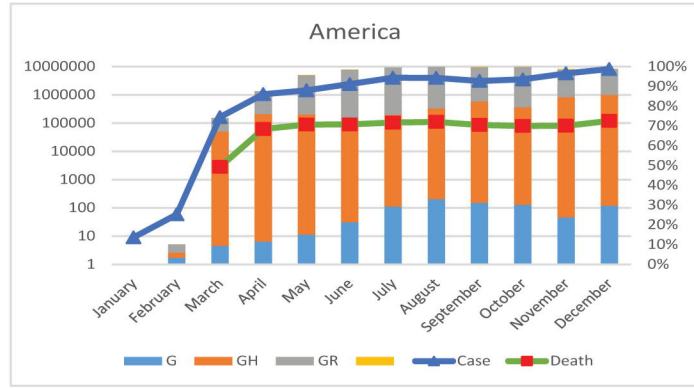
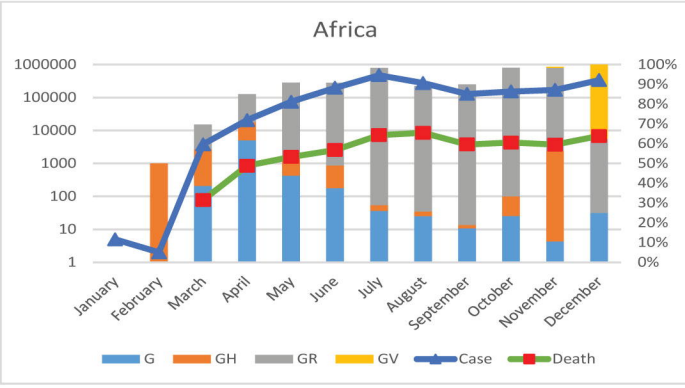
Protein Name	Mutation with position	$\Delta\Delta G$ DynaMut kcal/mol	$\Delta\Delta G$ ENCoM kcal/mol	$\Delta\Delta G$ mCSM kcal/mol	$\Delta\Delta G$ SDM kcal/mol	$\Delta\Delta G$ DUET kcal/mol	$\Delta\Delta G$ FoldX (kcal/mol)	Results*	$\Delta\Delta S_{vib}ENCoM$ kcal.mol ⁻¹ .K ⁻¹
RdRp	P323L	1.054	-0.441	-0.264	0.700	0.118	-0.733	Stabilizing	-0.551
Spike	D614G	-0.769	+0.408	-0.492	2.530	0.195	+0.289	Destabilizing	0.510
ORF3a	Q57H	0.275	-0.128	0.788	0.520	-0.464	-1.438	Stabilizing	-0.160
N Protein	RG203-04KR	-	-	-	-	-	-3.42262	Highly Destabilizing	-
N protein	A220V	0.109	0.458	-0.586	-1.460	-0.567	+1.6	Stabilizing	-0.572

1341
 1342 ***The final result of the stability for each protein was determined based on the intra-**
 1343 **molecular interactome analysis.**

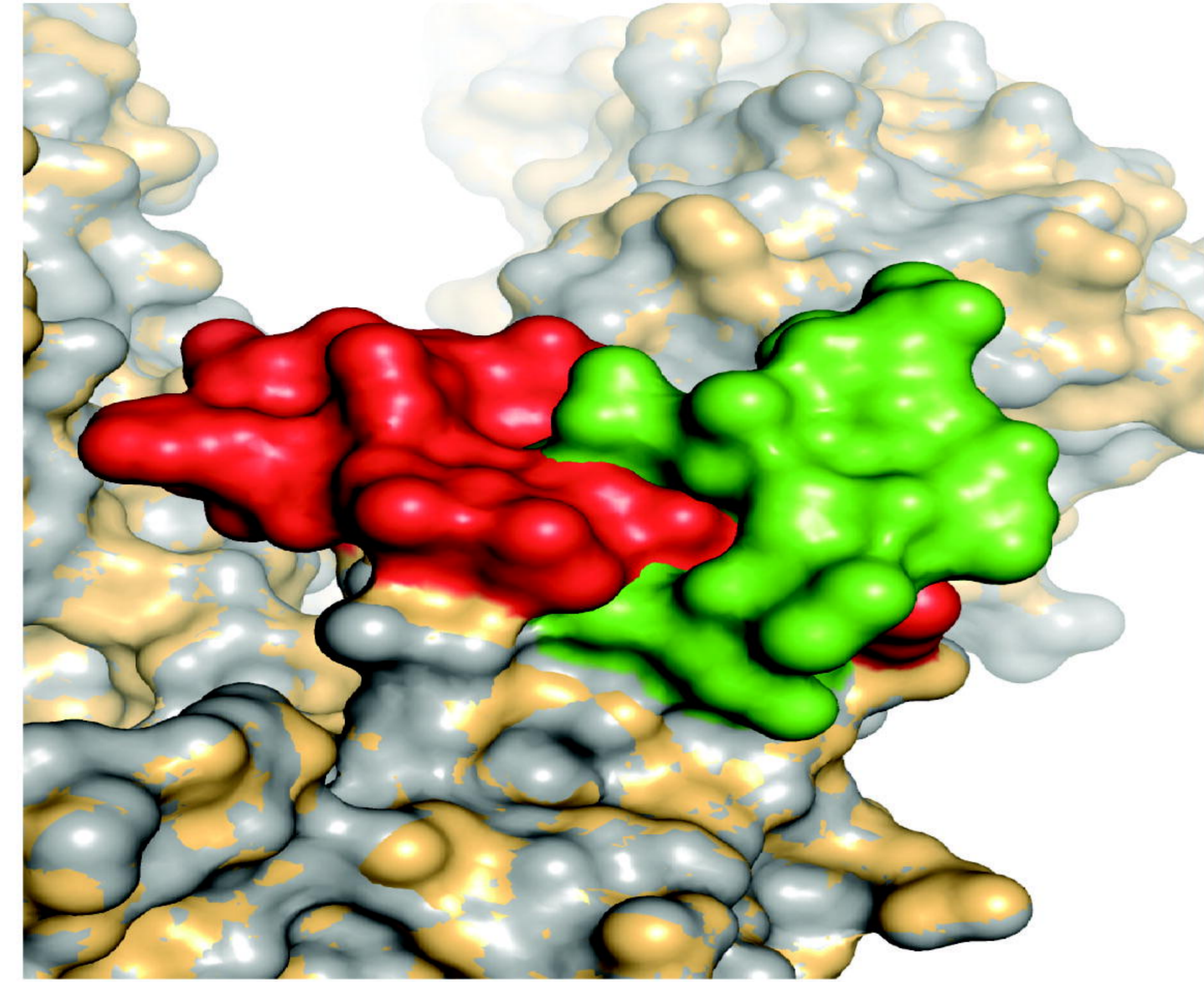
medRxiv preprint doi: <https://doi.org/10.1101/2021.02.21.21252137>; this version posted February 23, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted medRxiv a license to display the preprint in perpetuity. All rights reserved. No reuse allowed without permission.



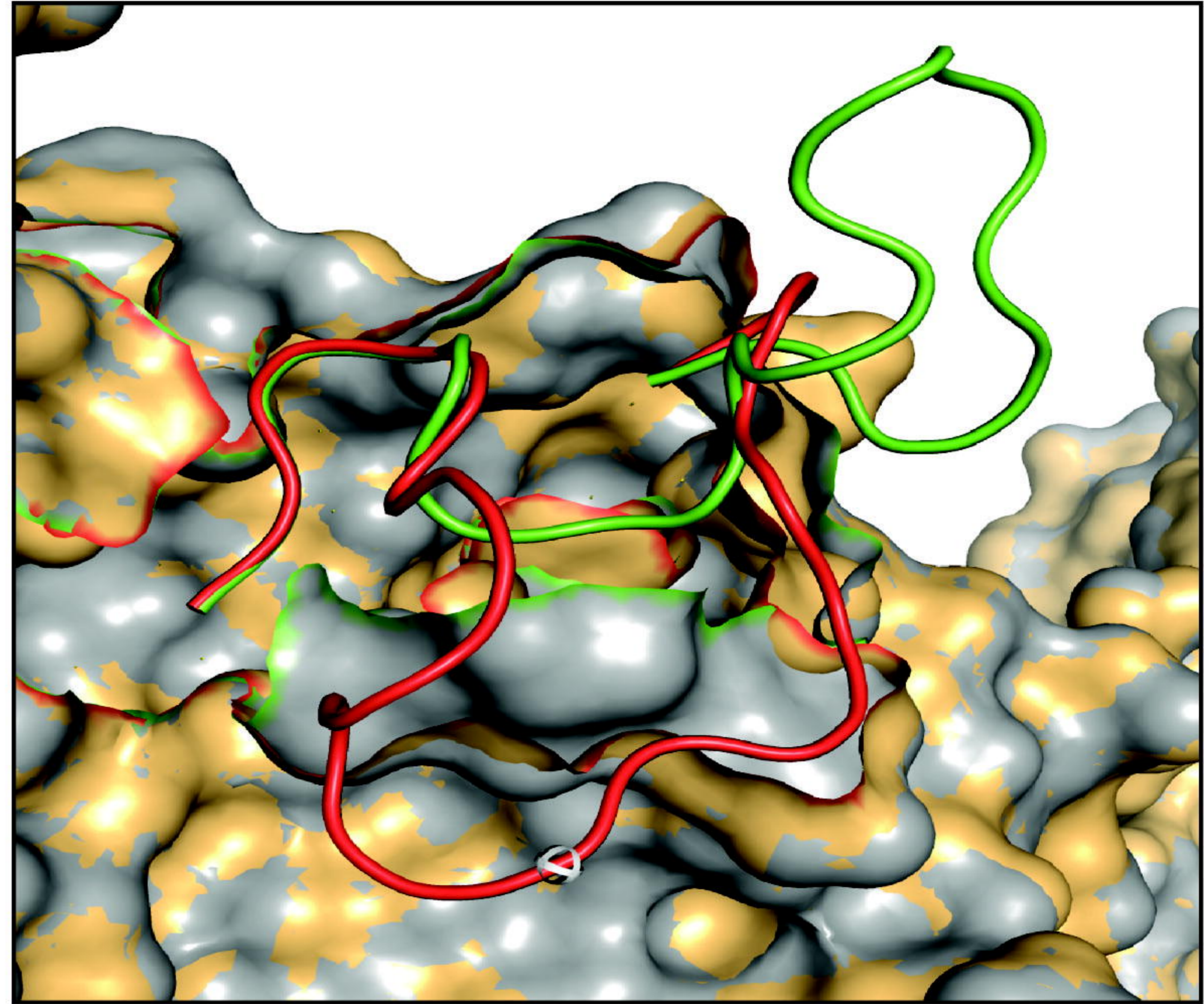
a



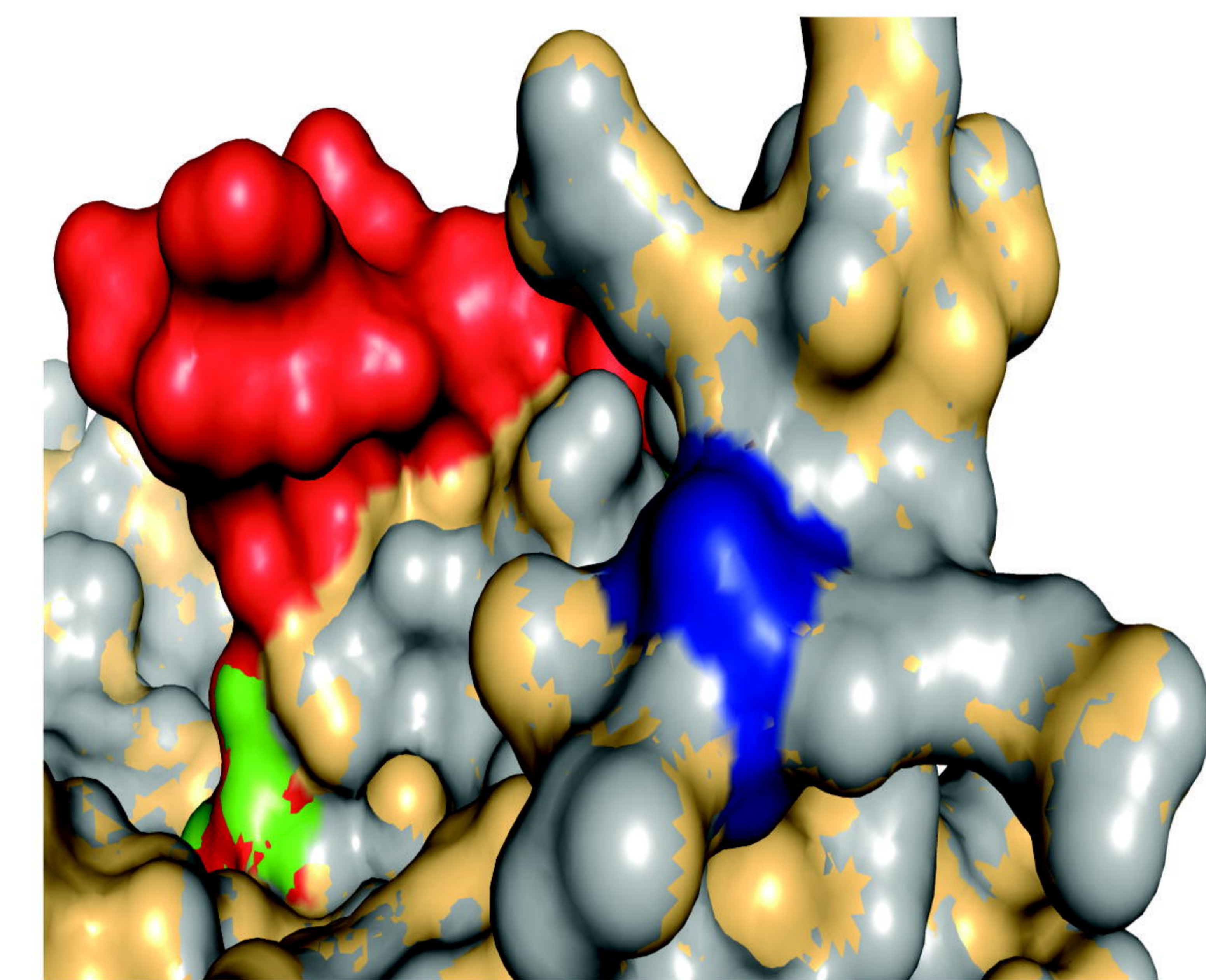
b



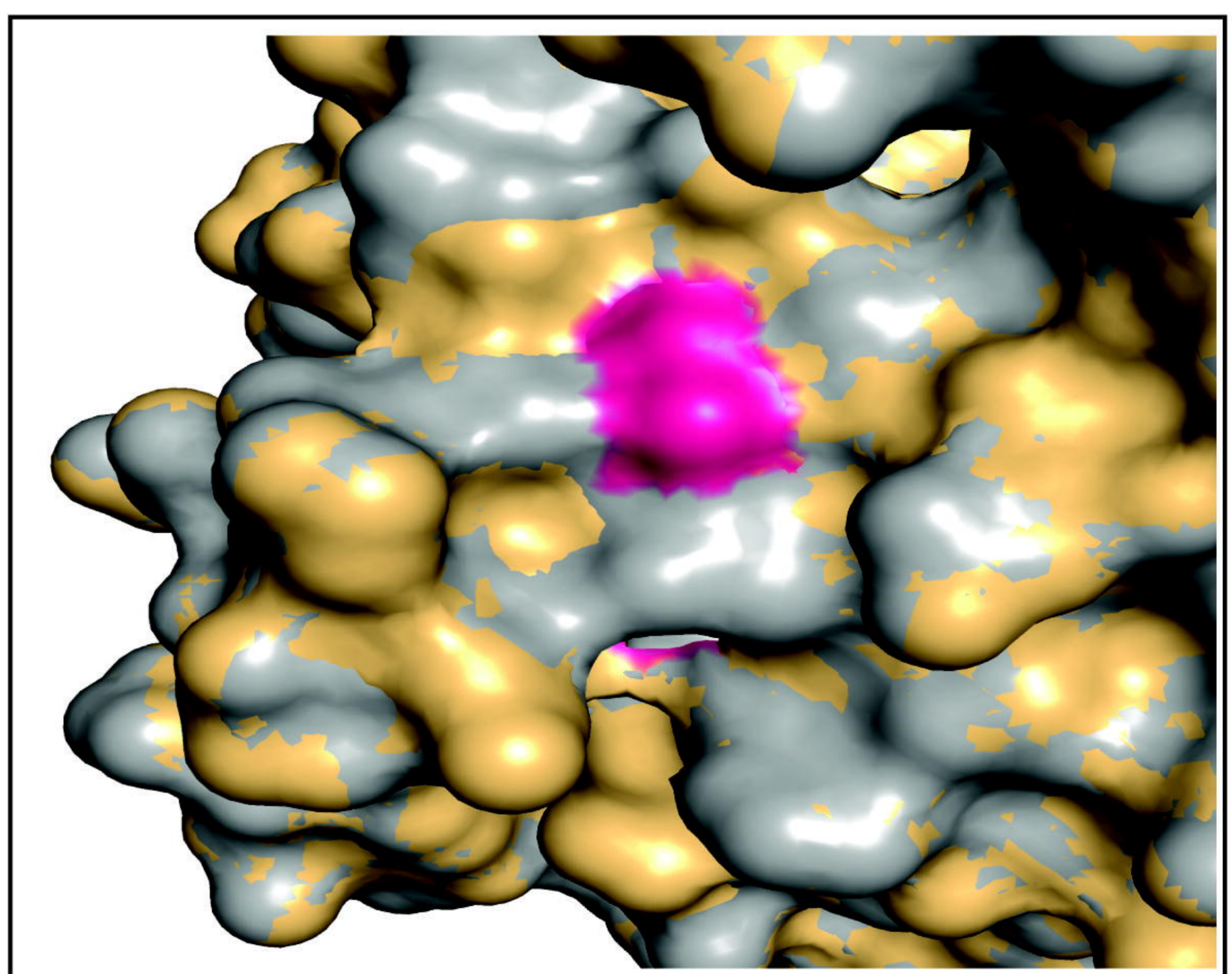
a



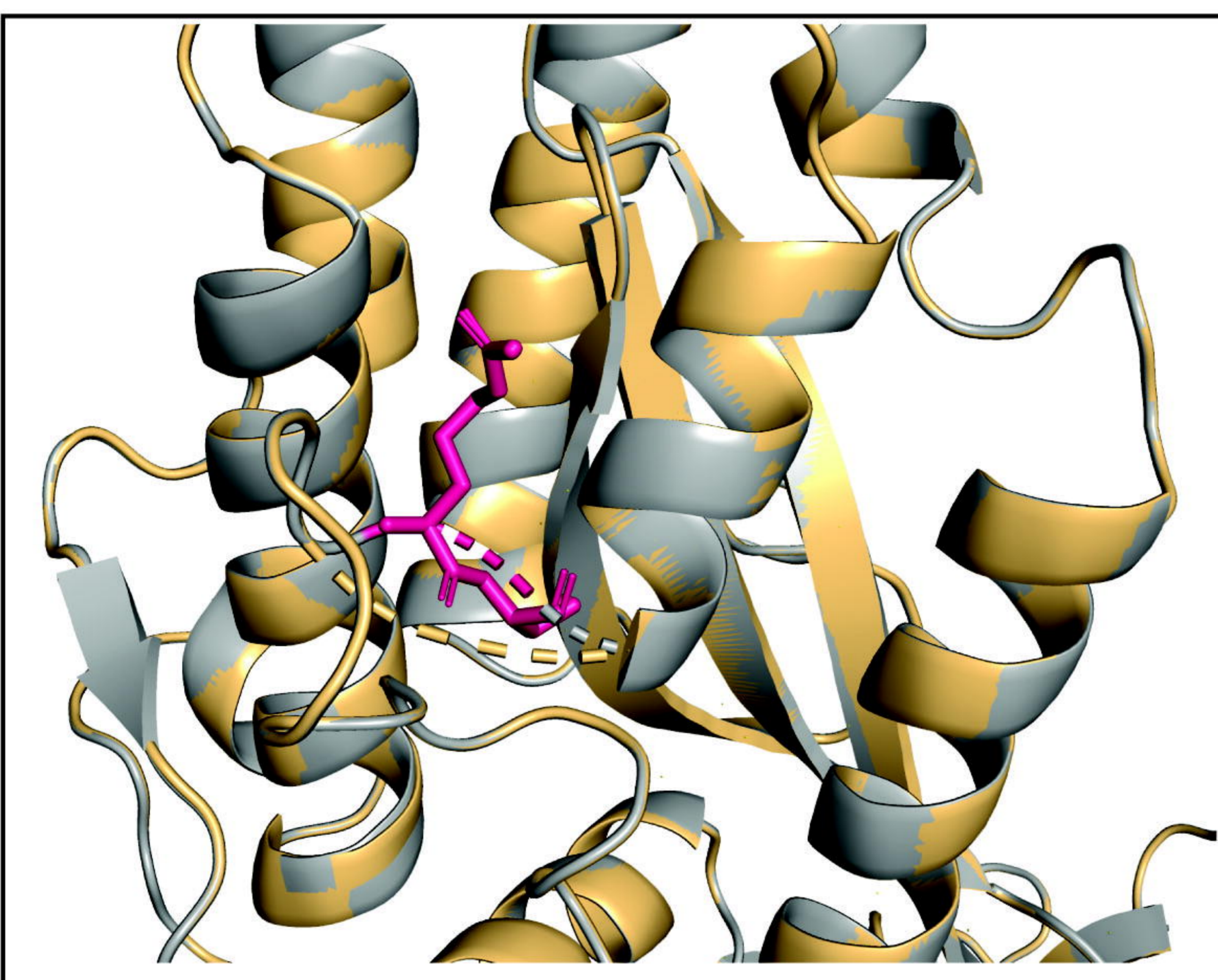
b



c



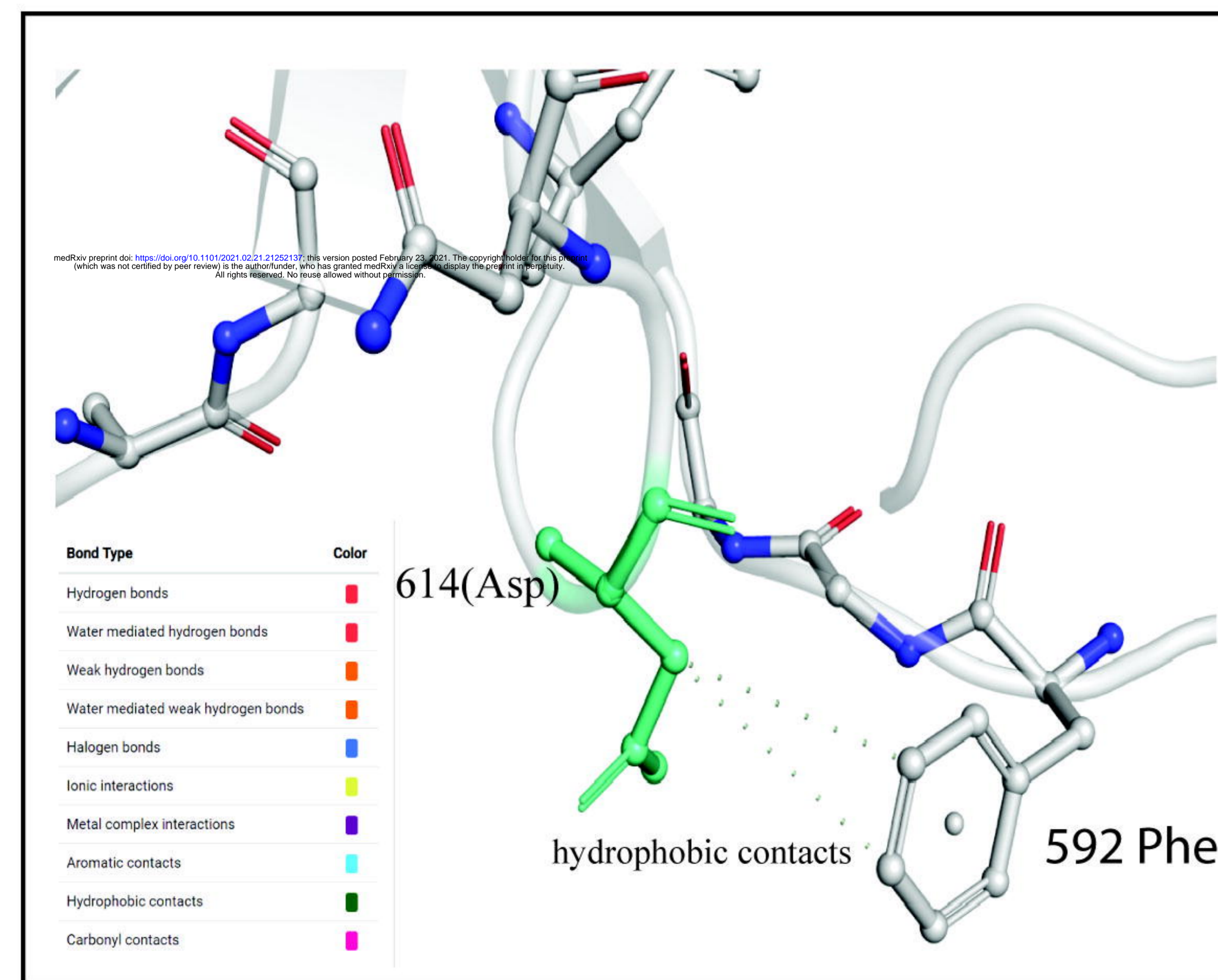
d



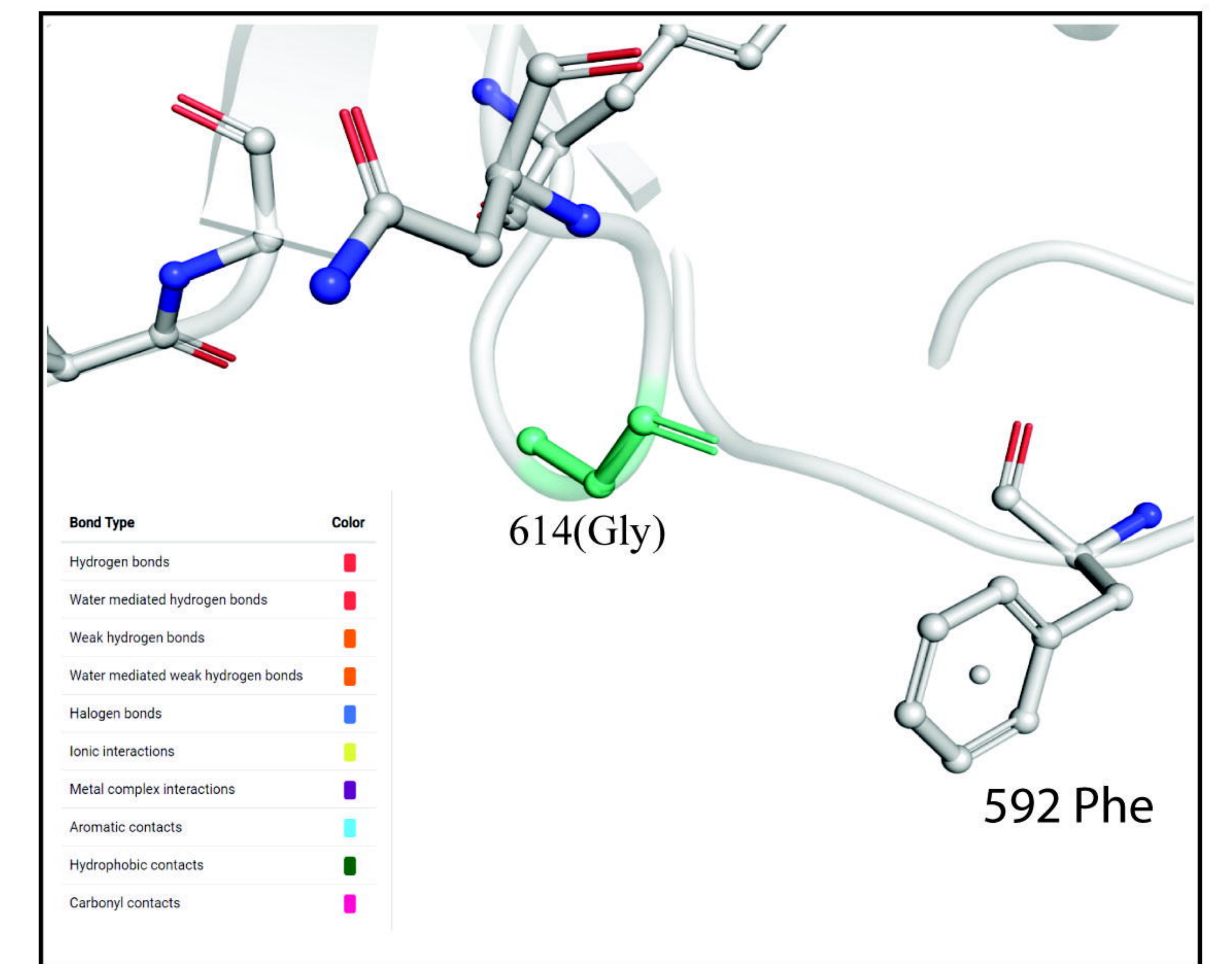
e



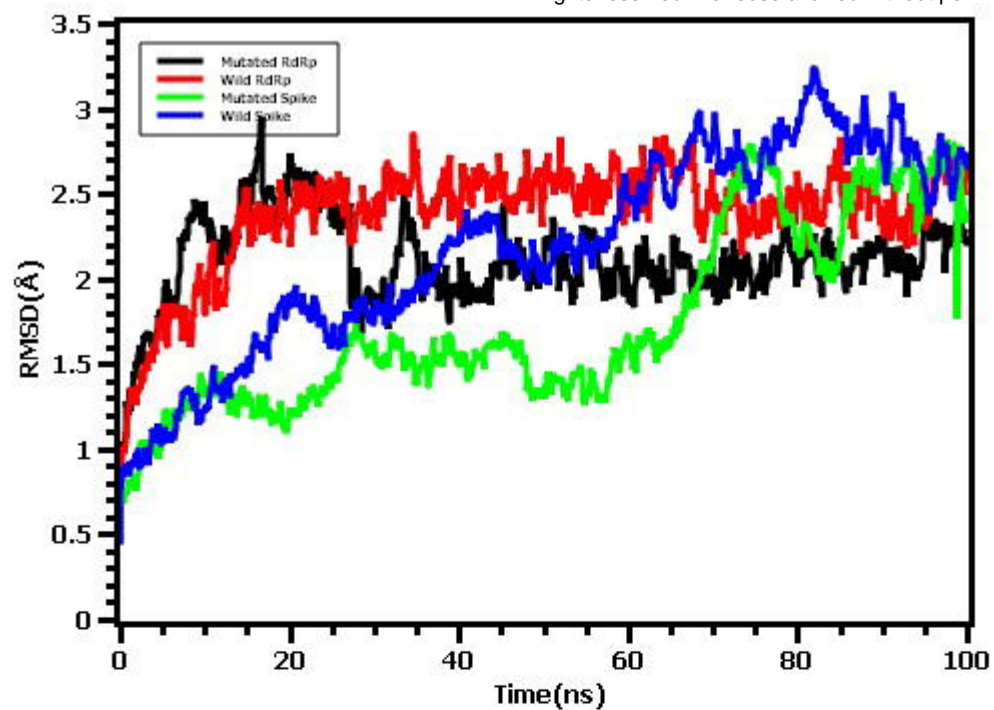
f



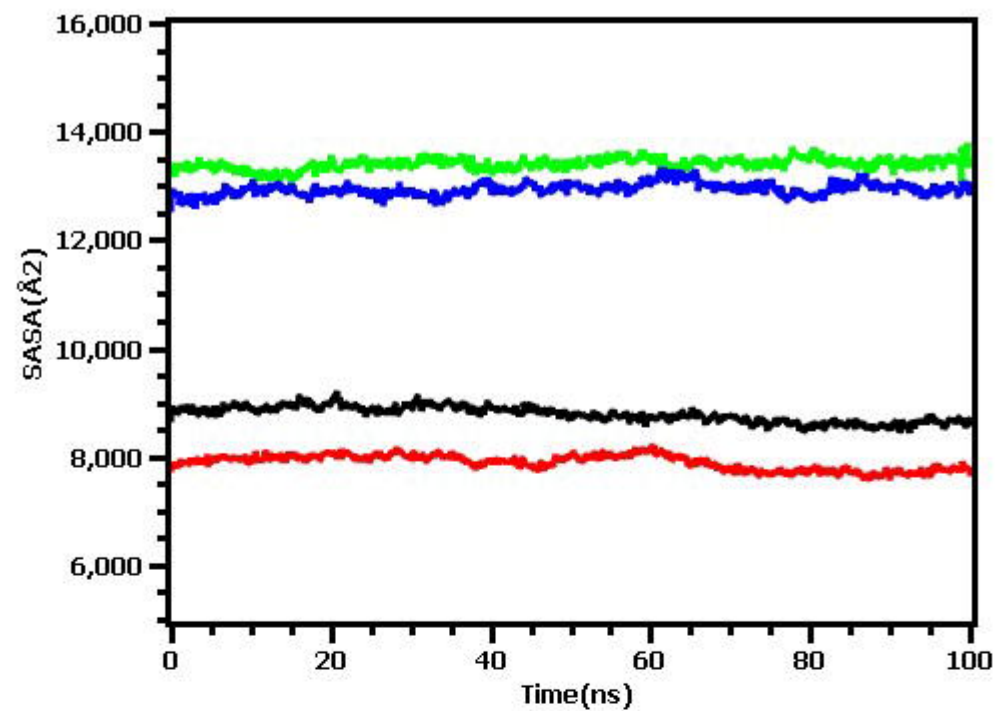
g



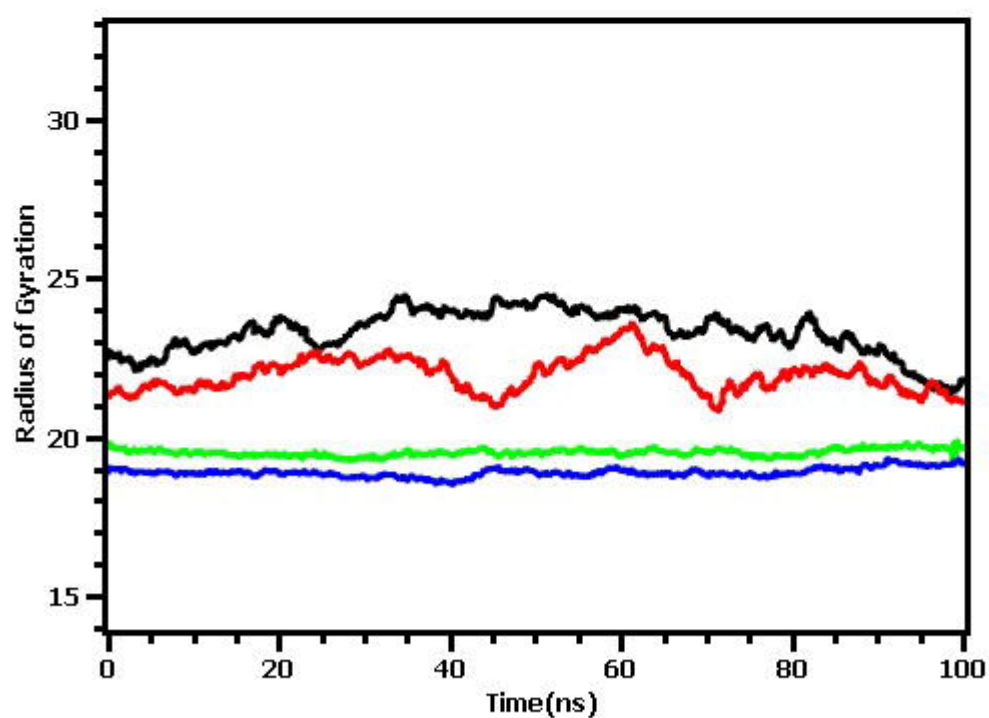
h



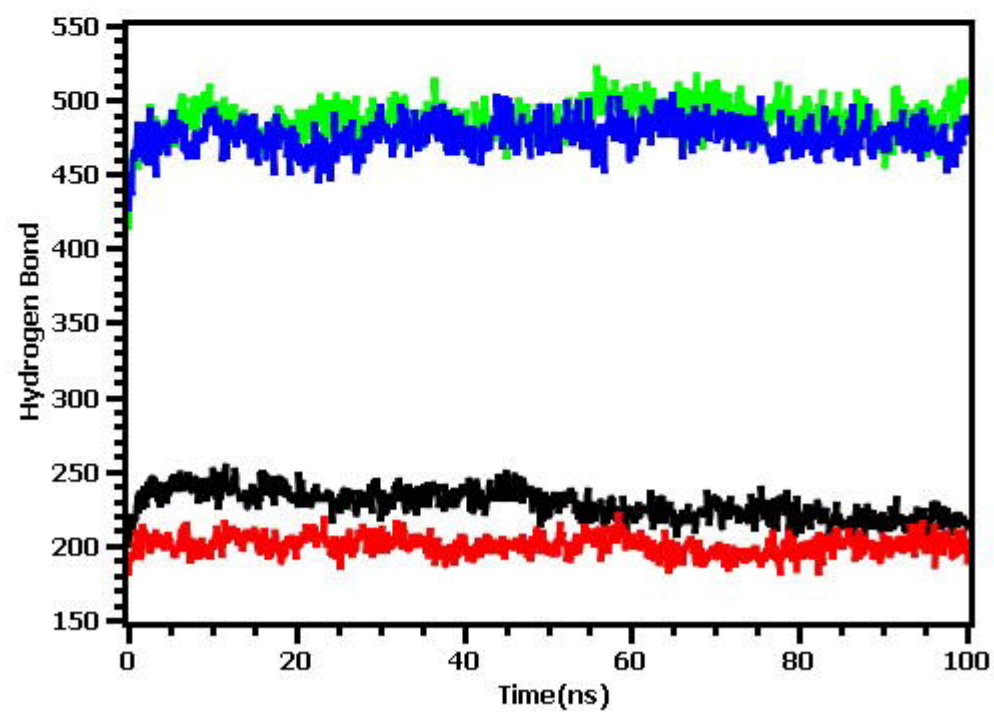
(a)



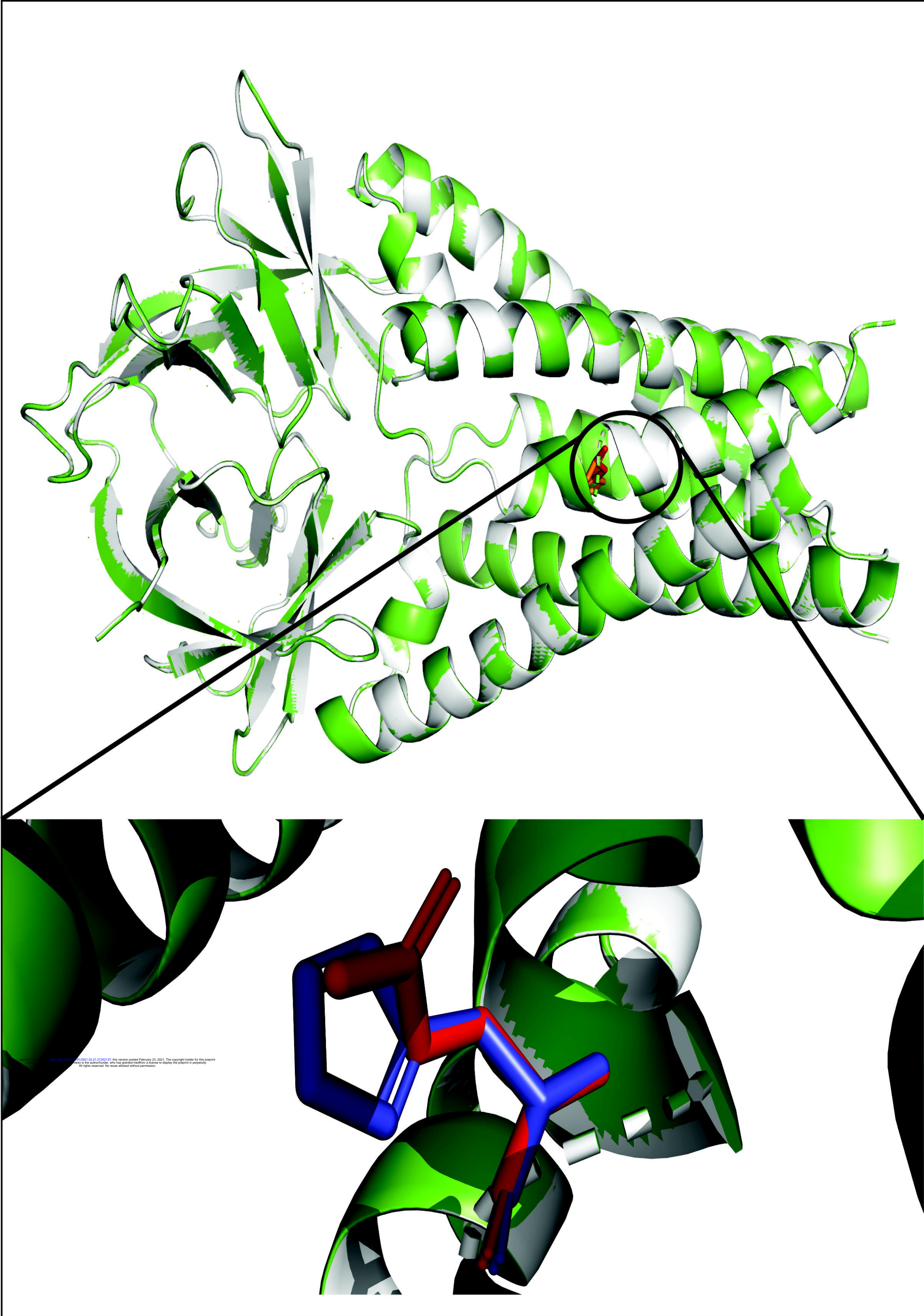
(b)



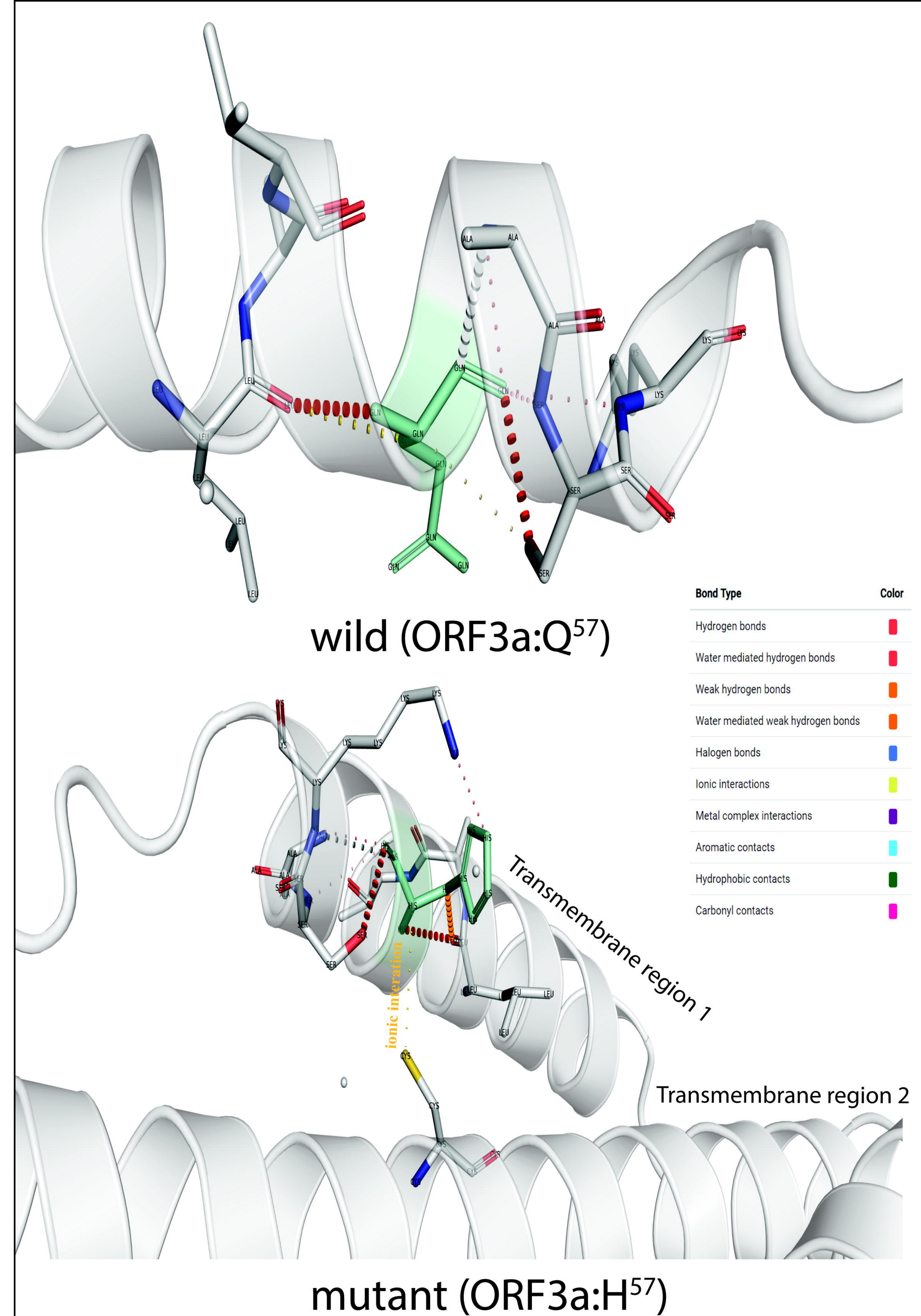
(c)



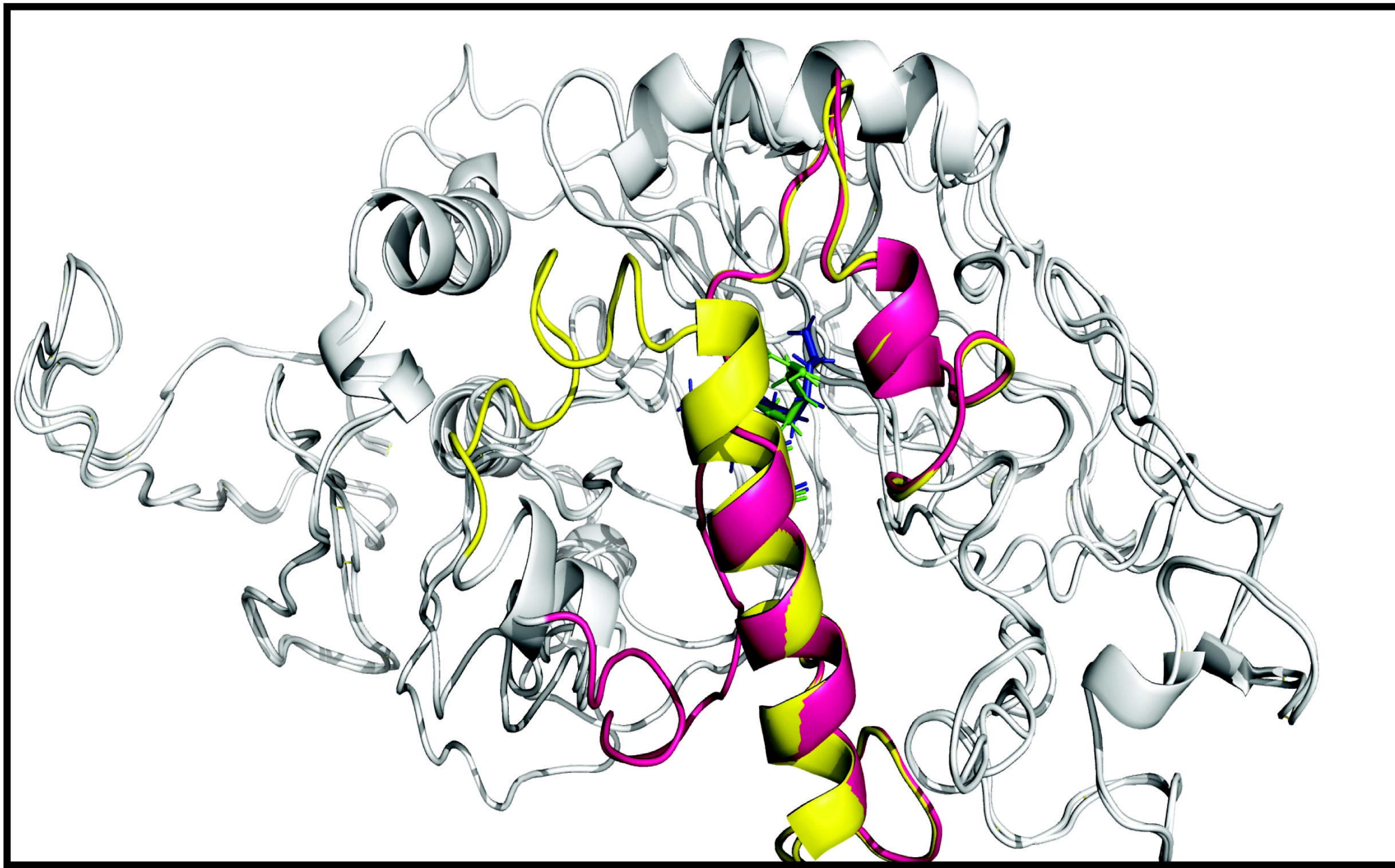
(d)



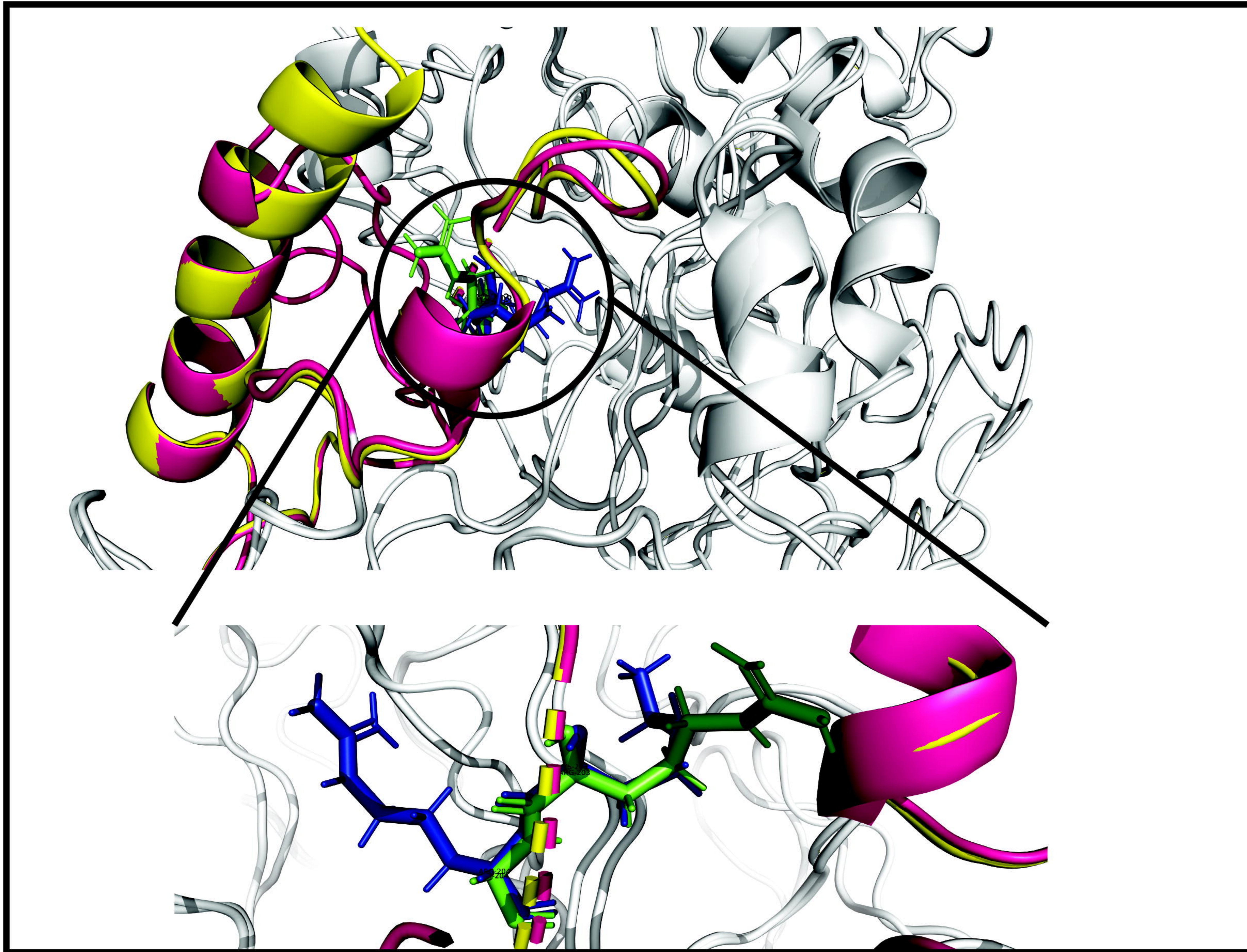
a



b



a



b