

1 **Title:** Novel clinical subphenotypes in COVID-19: derivation, validation, prediction, temporal  
2 patterns, and interaction with social determinants of health

3  
4 **Authors:** Chang Su, PhD<sup>1</sup>; Yongkang Zhang, PhD<sup>1</sup>; James H Flory<sup>2</sup>, MD; Mark G. Weiner, MD<sup>1</sup>;  
5 Rainu Kaushal\*, MD<sup>1,3</sup>; Edward J. Schenck\*, MD<sup>3,4</sup>; Fei Wang\*, PhD<sup>1</sup>

6  
7 **Affiliations:**

8 <sup>1</sup>Department of Population Health Sciences, Weill Cornell Medicine, New York, NY;

9 <sup>2</sup>Memorial Sloan-Kettering Cancer Center, New York, NY;

10 <sup>3</sup>New York-Presbyterian Hospital, Weill Cornell Medicine, New York, NY;

11 <sup>4</sup>Division of Pulmonary & Critical Care Medicine, Joan and Sanford I. Weill Department of  
12 Medicine, Weill Cornell Medicine, New York, NY

13 \*Corresponding authors

14  
15  
16  
17  
18 **Abstract**

19 The coronavirus disease 2019 (COVID-19) is heterogeneous and our understanding of the  
20 biological mechanisms of host response to the novel viral infection remains limited. Identification  
21 of meaningful clinical subphenotypes may benefit pathophysiological study, clinical practice, and  
22 clinical trials. Here, our aim was to derive and validate COVID-19 subphenotypes using machine  
23 learning and routinely collected clinical data, assess temporal patterns of these subphenotypes  
24 during the pandemic course, and examine their interaction with social determinants of health  
25 (SDoH). We retrospectively analyzed 14418 COVID-19 patients in five major medical centers in  
26 New York City (NYC), between March 1 and June 12, 2020. Using clustering analysis, four  
27 biologically distinct subphenotypes were derived in the development cohort (N = 8199).  
28 Importantly, the identified subphenotypes were highly predictive of clinical outcomes (especially  
29 60-day mortality). Sensitivity analyses in the development cohort, and re-derivation and  
30 prediction in the internal (N = 3519) and external (N = 3519) validation cohorts confirmed the  
31 reproducibility and usability of the subphenotypes. Further analyses showed varying  
32 subphenotype prevalence across the peak of the outbreak in NYC. We also found that SDoH  
33 specifically influenced mortality outcome in Subphenotype IV, which is associated with older age,  
34 worse clinical manifestation, and high comorbidity burden. Our findings may lead to a better  
35 understanding of how COVID-19 causes disease in different populations and potentially benefit  
36 clinical trial development. The temporal patterns and SDoH implications of the subphenotypes  
37 may add new insights to health policy to reduce social disparity in the pandemic.

38  
39

## 40 [Main text]

### 41 Introduction

42 The outbreak of coronavirus disease 2019 (COVID-19), caused by the severe acute respiratory  
43 syndrome coronavirus 2 (SARS-CoV-2) infection, has led to a pandemic that imposed  
44 tremendous pressure on healthcare systems globally<sup>1</sup>. As the pandemic continues and the  
45 second wave has emerged in the US and many other countries, research is still needed to  
46 understand how SARS-CoV-2 causes the wide spectrum of COVID-19 disease. Previous  
47 studies have uncovered substantial variation in the host response to SARS-CoV-2 and the  
48 variable clinical manifestations of this disease, including respiratory failure, kidney injury, and  
49 cardiovascular dysfunction<sup>2-8</sup>. Pivotal studies of corticosteroids<sup>9</sup> and anticoagulation<sup>10,11</sup>  
50 demonstrate differential responses in distinct subpopulations based on severity of disease. The  
51 pathophysiology of differential organ dysfunction in COVID-19 remains unclear across varied  
52 patient populations. Prior to the COVID-19 pandemic, identification of biologically distinct, data  
53 driven subphenotypes<sup>12,13</sup> has helped to disentangle complex syndromic disease such as  
54 sepsis<sup>14,15</sup>, ARDS<sup>16</sup>, heart failure<sup>17,18</sup>, diabetes<sup>19</sup>, and Alzheimer's disease<sup>20</sup>.

55  
56 Identifying robust subphenotypes in COVID-19 patients could lead to improved understanding of  
57 biological mechanisms of host response to SARS-CoV-2 infection and may identify  
58 subpopulations that could be prioritized for clinical trial enrollment<sup>13,21</sup>. Previous efforts<sup>22-25</sup> have  
59 been made in this area but remain limited probably due to cohort size, data availability, and  
60 lacking evaluation of robustness and usability of the identified subphenotypes. In addition, the  
61 hospitalized case fatality rate of COVID-19 has varied over the course of the pandemic<sup>26,27</sup> and  
62 according to social determinants of health (SDoH)<sup>28-30</sup>. Exploration of temporal patterns and

63 SDoH characteristics in conjunction with subphenotypes may derive new insights to improve  
64 public health.

65  
66 In this analysis, our goal was to derive and validate COVID-19 subphenotypes amongst a  
67 population of patients who presented to the emergency department (ED) or were hospitalized in  
68 multiple health systems in New York City (NYC). Specifically, we used routinely collected clinical  
69 data to first derive subphenotypes using the agglomerative hierarchical clustering model. Then,  
70 multiple strategies in data pre-processing, data filtering, and data-driven models (both  
71 unsupervised clustering model and supervised predictive model) were used to confirm  
72 reproducibility and usability of the identified subphenotypes. After that, statistical analyses were  
73 conducted to evaluate the characteristics and clinical outcomes of the subphenotypes. Further  
74 analyses were performed to examine temporal patterns of the subphenotypes and impacts of  
75 SDoH status on subphenotype-level outcomes. The overall workflow of our study is illustrated in  
76 [Figure 1](#).

77

78

## 79 **Results**

### 80 **Patients**

81 A total of 14418 patients with confirmed COVID-19 between March 1st and June 12<sup>th</sup> 2020,  
82 treated in ED (N=2354, 16.3%) or inpatient (N=12064, 83.7%) settings, were included for  
83 analysis from the five major medical centers in New York City (NYC), including New York  
84 University Langone Medical Center (NYU-LMC), New York Presbyterian - Weill Cornell Medical  
85 Center (NYP-WCMC), Mount Sinai Health System (MSHS), Montefiore Medical Center (MMC),

86 and New York Presbyterian - Columbia University Medical Center (NYP-CUMC). Details of  
87 inclusion and exclusion criteria are presented in [eFigure 1](#) in Supplement. We identified 2853  
88 (19.8%) deaths within 60-day after COVID-19 confirmation in total, including 2801 (19.4%) in-  
89 hospital deaths and 52 (4%) deaths after discharge from COVID related hospitalization or ED  
90 visits. Considering population diversity (especially race) of the five medical centers (see [eTable](#)  
91 [1](#) in Supplement), we combined four centers and randomly divided them into the development  
92 cohort (70%) and internal validation cohort (30%); patients of the remaining center were used as  
93 the external validation cohort (see [Figure 1](#) and [eFigure. 1](#) in Supplement).

94  
95 The development cohort contained a total of 8199 patients with a median age of 65.35  
96 (interquartile range [IQR] [50.57, 75.17]) years old, consisting of 3787 (46.2%) females, 2036  
97 (24.8%) white patients, and 2155 (26.3%) black patients. The internal validation cohort  
98 contained a total of 3519 patients with similar patient characteristics when compared with the  
99 development cohort, with a median age of 63.51 (IQR [50.95, 75,17]) years old, consisting of  
100 1585 (45.0%) females, 838 (23.8%) white patients, and 915 (26%) black patients. The external  
101 validation cohort contained a total of 2700 patients. It had a median age of 65.85 (IQR [51.08,  
102 77.38]) years old and consisted of 1305 (48.3%) females, 675 (25.0%) white patients, and 545  
103 (20.2%) black patients. Across the three cohorts, the overall 60-day mortality rates after ED or  
104 hospital discharge were 18.65%, 19.78%, and 20.59%, respectively. More details of the  
105 characteristics of the studied cohorts appeared in [Table 1](#).

## 106 **Subphenotypes derivation**

107 In the development cohort, the agglomerative hierarchical clustering model identified 4 distinct  
108 subphenotypes based on presenting clinical data of the patients (see [eResults](#), [eFigures 3 and](#)  
109 [4](#) in Supplement). Characteristics including demographics, clinical variables, comorbidities,

110 clinical outcomes, and medication treatments across the 4 subphenotypes were presented in  
111 [Table 2 and Figures 2 and 3](#).

112  
113 **Subphenotype I** consisted of 2707 (33.02%) patients. Compared to the others, it included more  
114 younger (median age 57.45 years, IQR [42.70, 70.02]) and female (N = 1601 [59.15%]) patients.  
115 Those patients had more normal values across all clinical variables and lower chronic  
116 comorbidity burden. The patients also had better clinical outcomes including a low 60-day  
117 mortality (N = 188 [6.94%]) and a low rates of mechanical ventilation (N = 190 [7.02%]) and ICU  
118 admission (N = 242 [8.94%]).

119  
120 **Subphenotype II** consisted of 3047 (37.16%) patients. Compared to other subphenotypes, it  
121 included more male patients (N = 1941 [63.70%]) and was likely to have more abnormal  
122 inflammatory markers (such as C-reactive protein, erythrocyte sedimentation rate, interleukin 6,  
123 lactate dehydrogenase, lymphocyte count, neutrophil count, white blood cell count, and ferritin)  
124 and markers of hepatic dysfunctions (such as ferritin, alanine aminotransferase, aspartate  
125 aminotransferase, and bilirubin). Overall comorbidity burden of Subphenotype II was low.  
126 Clinical outcomes including 60-day mortality (N = 528 [17.33%]), mechanical ventilation (N =  
127 527 [17.30%]), and ICU admission (N = 675 [22.15%]) of Subphenotype II were at a moderate  
128 level.

129  
130 **Subphenotype III** included 1486 (18.12%) patients, consisting of more older (median age 69.45  
131 years, IQR [57.05, 79.62]) and black (N = 503 [33.85%]) patients, compared to subphenotypes I  
132 and II. Those patients of Subphenotype III were likely to have more abnormal renal dysfunction  
133 markers (such as blood urea nitrogen, creatinine, chloride, and sodium) and hematologic  
134 dysfunction markers (such as d-dimer, hemoglobin, and red blood cell distribution width).  
135 Overall comorbidity burden of Subphenotype III was high. Clinical outcomes including 60-day

136 mortality (N = 337 [22.68%]), intubation (N = 195 [13.12%]), and ICU admission (N =242  
137 [16.29%]) of Subphenotype II were close to that of Subphenotype II and at a moderate level as  
138 well.

139  
140 **Subphenotype IV** included 959 (11.70%) patients. Compared to other subphenotypes, it  
141 included more older (median age 75.53 years, IQR [64.10, 84.83]) and male (N = 588 [61.31%])  
142 patients. Those patients of Subphenotype IV had more abnormal values across all clinical  
143 variables and higher chronic comorbidity burden than the others. Obesity burden is lower in  
144 Subphenotype IV than others. In line with its biological characteristics, Subphenotype IV had the  
145 worst clinical outcomes in 60-day mortality (N = 476 [49.64%]), intubation (N = 242 [25.23%]),  
146 and ICU admission (N =335 [34.93%]). In addition, the medications including antibiotics,  
147 corticosteroids, and vasopressor were more frequently used in Subphenotype IV.

148

## 149 **Subphenotype reproducibility and prediction**

150 In the development cohort, sensitivity analyses under two different settings (sensitivity  
151 to quality control and outliers and sensitivity to clustering methods) confirmed the  
152 underlying 4-cluster structure of the data (see [eResults](#), [eFigures 3 and 4](#), and [eTable 5](#) in  
153 Supplement). Patients' memberships of the 4 clusters re-derived by sensitivity analyses  
154 were highly consistent with those derived in the primary analysis (see [eFigure 6](#) in  
155 Supplement). Moreover, we did not find substantial changes in clinical characteristics of  
156 the subphenotypes in the sensitivity analyses (see [eTables 6 and 7](#) in Supplement).

157

158 Subphenotypes were also re-derived in the internal validation cohort, where the 4-  
159 cluster structure was found as the optimal fit as well (see [eResults](#) and [eFigure7](#) in

160 Supplement). Clinical characteristics of the re-derived subphenotypes in the internal  
161 validation cohort, including demographics, laboratory variables, comorbidities, and  
162 clinical outcomes, also showed very similar patterns with the subphenotypes derived in  
163 the primary analysis (see [Figure 3](#), and [eTable 8](#) and [eFigure 8](#) in Supplement).

164  
165 To further evaluate subphenotype robustness and usability, we trained a predictive  
166 model of subphenotypes in the development cohort and used it to predict subphenotype  
167 membership in the external validation cohort. Clinical variables of presenting laboratory  
168 tests for clustering analysis were used as candidate predictors. The trained predictive  
169 model (XGBoost classifier) achieved very high performance in predicting each  
170 subphenotype (see [eFigure 9](#) in Supplement). SHapley Additive exPlanation (SHAP) values  
171 illustrated contributions of the clinical variables in distinguishing each subphenotype from others  
172 (see [eFigure 10](#) in Supplement). Patterns of the produced SHAP values were highly in line with  
173 the subphenotype characteristics: 1) normal values of the clinical variables indicated  
174 Subphenotype I; 2) abnormal inflammatory and hepatic markers were predictive of  
175 Subphenotype II; 3) abnormal renal and hematologic markers indicated were likely to indicate  
176 Subphenotype III; 4) Subphenotype IV was associated with abnormal values of most variables.  
177 After that, the trained predictive model was used to predict subphenotype memberships of  
178 patients in the external validation cohort. The predicted subphenotypes in the external validation  
179 cohort were well separated in the UMAP space (see [eFigure 11](#) in Supplement) and showed  
180 clinical characteristics similar to findings in the primary analysis (see [Figure 3](#), and [eTable 9](#) and  
181 [eFigure 12](#) in Supplement).

182

183 Last, results from leave-one-center-out analysis also confirmed the four-cluster structure  
184 underlying our data (see [eFigure 13](#) in Supplement). Meanwhile, subphenotypes identified by  
185 the leave-one-center-out analysis among the whole population showed characteristics in line  
186 with those identified in the primary analysis (see [eTable 10](#) in Supplement). Those above  
187 demonstrated stability of the identified subphenotypes across the five centers.

## 188 **Temporal characteristics of subphenotypes**

189 Temporal patterns of the COVID-19 subphenotypes were illustrated by the bar charts, showing  
190 the composition of subphenotype memberships of patients confirmed per week, since the  
191 outbreak in NYC, i.e., March 1, 2020 (see [Figures 4a-c](#)). Except week 1 and week 14 that had  
192 few patients confirmed, the composition of the four subphenotypes per week evolved over time  
193 and showed similar patterns across the development, internal validation, and external validation  
194 cohorts. In general, patients with confirmed SARS-CoV-2 infection rapidly increased within the  
195 first month since the outbreak and reached the peak at week 5 (early April). Subphenotype I  
196 (mild symptom) and Subphenotype II (moderate symptom, low comorbidity burden) dominated  
197 the time period prior to the peak (first 4 weeks since outbreak). In contrast, Subphenotype IV  
198 (severe symptom, high comorbidity burden) had a low proportion within the first 4 weeks but  
199 showed a largely increased proportion from week 6 to week 9. Since week 10, the proportion of  
200 Subphenotype I gradually increased while others especially Subphenotype IV shrank.  
201 Subphenotype III (moderate symptom, high comorbidity burden) had a relatively stable  
202 proportion over time.

## 203 **Impact of SDoH on subphenotypes**

204 In general, worse SDoH in terms the socioeconomic variables were likely in Subphenotype IV  
205 (see [eTable 11](#) in Supplement). Moreover, logistic regression analysis identified similar patterns



206 of relationships between the SDoH variables with 60-day mortality risk across subphenotypes;  
207 however, absolute log odds and Hazard ratio of the SDoH variables varied across  
208 subphenotypes (see [Figure 4d](#) and [eTables 12-13](#) in Supplement). For example, low absolute  
209 log odds were observed in all six SDoH variables in Subphenotype I. In contrast, we did see  
210 increased absolute log odds of all six SDoH variables in Subphenotype IV. Hazard ratio showed  
211 similar pattern.

212  
213 Agglomerative hierarchical clustering based on the SDoH variables grouped the patients into a  
214 3-cluster model (see [eResults](#) and [eFigure 14](#) in Supplement), which can be interpreted as high  
215 (H), middle (M), and low (L) SDoH strata (see [eTable 14](#) in Supplement). Stratum L,  
216 representing disadvantaged SDoH status, accounted for a slightly higher mortality rate (H vs. M  
217 vs. L, 17.59% vs. 19.91% vs. 19.98%, P-value = 0.08). In addition, stratum L had a lower ICU  
218 admission rate (16.16%, P-value < 0.001). The relative high mortality but low ICU admission  
219 rate may be caused by critical care strain during periods of increased COVID-19 ICU demand,  
220 as suggested by a recent study<sup>31</sup>. Distributions of the SDoH strata by biological subphenotypes  
221 were shown in [eTable 15](#) in Supplement. In the analysis to further explore how SDoH strata  
222 affected the outcome of each biological subphenotype, we found varied patterns of correlations  
223 between SDoH strata and 60-day mortality (see [Figure 4e](#)) by subphenotypes. Notably, in line  
224 with the results of the univariate analysis above, SDoH strata were likely to have a strong  
225 impact on the 60-day mortality in Subphenotype IV. Particularly, in Subphenotype IV, SDoH  
226 stratum L was associated with a 55.19% 60-mortality rate, which was 5.55% higher than the  
227 subphenotype level (49.64%, see [Table 2](#)) and 8.52% higher than that of the SDoH stratum H.  
228 In subphenotypes I, II, and III, we didn't find mortality rate discrepancy higher than 3% between  
229 any pair of SDoH strata. Similarly, considering stratum H as reference, stratum L had largely

230 increased log odds of mortality in Subphenotype IV (log odds = 0.40, SD = 0.19, P-value = 0.04).

231 (see [eTable 16](#) in Supplement)

232

## 233 **Discussions**

234 We derived subphenotypes of COVID-19 patients treated at five major medical centers in NYC

235 across the whole course of the first wave of the pandemic, using the clinical data at the

236 presentation to the emergency department (ED) or hospital. Different from the previous

237 subphenotype studies of COVID-19<sup>22-24</sup>, we focused on a larger, more representative, and

238 diverse population presented at the ED and/or hospitalized without COVID-19 specific therapy.

239 We derived subphenotypes using clustering analysis in the development cohort and validated

240 them using a combination of multiple validation strategies, including the use of different data

241 processing, different data filtering, and different machine learning models (both unsupervised

242 clustering and supervised predictive models). All validation approaches confirmed the

243 reproducibility of the 4-cluster structure of the data and clinical characteristics of the identified

244 subphenotypes. We would also highlight that all machine learning models used for

245 subphenotype derivation and validation were performed only on the presenting clinical variables

246 that were routinely collected in daily patient care and are available to providers by ED or

247 hospital admission. This allows us to potentially capture the underlying variable mechanisms of

248 the complex disease, but also enhances the generalizability and feasibility of the identified

249 subphenotypes to be used in clinical practices and patient enrollment in clinical trials.

250

251 Importantly, the 4 subphenotypes identified were significantly separated in demographics,

252 clinical variables, and chronic comorbidities, and strongly predictive of the 60-day mortality

253 outcome. Subphenotype IV included more older, male patients, abnormal markers indicating

254 hyperinflammation, liver injury, cardiovascular problems, renal dysfunctions, and coagulation  
255 disorders, and a higher comorbidity burden (except for obesity) compared to the other  
256 subphenotypes. In contrast, Subphenotype I was composed of relatively healthy, younger  
257 females who had more normal values across all markers and comorbidity burdens compared to  
258 the other subphenotypes. There was a strong concordance between their clinical profiles and  
259 outcomes, such as Subphenotype IV showed the worst clinical outcome while Subphenotype I  
260 showed the best outcome among the 4 subphenotypes. These are in line with observations  
261 reported in a previous small cohort study<sup>23</sup>. Interestingly, Subphenotypes II and III showed  
262 similar, moderate-level 60-day mortality rates, but their clinical characteristic profiles suggested  
263 that they were likely to have distinct biological mechanisms. In particular, results from our  
264 primary analysis and validation approaches demonstrated that Subphenotype II was correlated  
265 with relative hyperinflammation, while Subphenotype III was associated with renal injury, lower  
266 platelet level and a high comorbidity burden (significantly higher than Subphenotypes I and II,  
267 and equivalent to Subphenotype IV). Moreover, in accordance with the clinical characteristics  
268 and outcomes, the worse subphenotypes (Subphenotypes III and IV) were more likely to receive  
269 medications in antibiotics, corticosteroids, and vasopressor than the others. These findings  
270 suggested that our identified subphenotypes offer insight into the varied mechanisms of COVID-  
271 19.

272  
273 Typically, data-driven approaches for the identification of subphenotypes of human disease are  
274 based on the unsupervised clustering methods<sup>12,14-16,22-24,32</sup>. The natural attributes of the  
275 unsupervised methodology in discovering underlying patterns from data make them the best fit  
276 for subphenotype identification. Once the subphenotypes were determined, there would be a  
277 need of subphenotype membership assignments for new patients. However, previous studies  
278 barely discussed such down-stream usability of the identified subphenotypes. In this analysis,  
279 we built a supervised predictive model of the identified subphenotypes. Our predictive model

280 achieved an ideal prediction performance in the development cohort and predicted  
281 subphenotypes in the external validation cohort that presented the same pattern of clinical  
282 characteristics with that of the originally derived subphenotypes. In this way, instead of  
283 validating the subphenotypes in a different route, the predictive model brought additional  
284 implications as: 1) it provides a feasible and accurate way to apply the identified subphenotypes  
285 to clinical practice; and 2) contributions of the clinical variables in subphenotype prediction  
286 calculated by the SHAP method showed concordant patterns with the subphenotypes' clinical  
287 characteristics and hence confirmed biological profiles of the subphenotypes in the multivariate  
288 prospective.

289  
290 Time is a crucial factor in the spread of COVID-19. Previous studies have examined the  
291 temporal trends of COVID-19 outcomes such as in-hospital mortality rate during the course of  
292 the pandemic<sup>26,27</sup>, but limited attention has been drawn on evolving patterns of COVID-19  
293 phenotypes. We filled this gap in the present study. Our observations suggested varied  
294 temporal trends of the identified subphenotypes during the first 14 weeks of the pandemic in  
295 NYC. Interestingly, since the COVID-19 outbreak in NYC on March 1, 2020, Subphenotypes I  
296 and II dominated the time period prior to the peak (first 4 weeks since outbreak), possibly as  
297 they contained more relatively younger patients who may have had more frequent social  
298 activities to be infected. Subphenotype IV, with older age, worse health conditions, and poorer  
299 outcomes, was boosted within the second month (April 2020) post spread peak, consistent with  
300 tremendous mortality rate of NYC in April<sup>33</sup>.

301 This would suggest that younger, biologically strong patients (Subphenotypes I and II) got  
302 infections early and boosted the spread, while older, biologically vulnerable patients  
303 (Subphenotype IV) accounted for the second infections within a population probably due to  
304 housing. After that, the proportion of Subphenotype I out of all patients confirmed per week  
305 gradually expanded while that of the others, especially Subphenotype IV shrank. The potential

306 reason would be that valuable experience (such as the improved use of masks and social  
307 distancing), reinforced health care systems, and announced health policies did protect the  
308 population who likely develop severe subphenotypes (Subphenotype IV). In general, such  
309 temporal trends of the biological subphenotypes would be a considerable, fine-grained  
310 explanation of the observed outcome (mortality rate) evolving trends in epidemiology<sup>26</sup>.

311  
312 SDoH such as vulnerable socioeconomic neighborhood status have been associated with poor  
313 outcomes of COVID-19<sup>26,30</sup>. In this work, we explored the impact of SDoH on different biological  
314 subphenotypes from both univariate and multivariate perspectives. We first examined the  
315 associations of individual socioeconomic characteristics with mortality risk in each  
316 subphenotype. We then derived comprehensive SDoH strata using the data-driven clustering  
317 method and evaluated their correlations with mortality risk in each subphenotype. The results  
318 confirmed our hypothesis that SDoH impacts biological subphenotypes differently. The highly  
319 expanded mortality risk log odds of individual SDoH variables and discrepancy of mortality rate  
320 among SDoH strata indicate that SDoH has a much stronger association with mortality  
321 outcomes in Subphenotype IV, compared to the others. In other words, once a sick, elderly  
322 patient shows up with COVID-19 (Subphenotype IV), the disadvantaged socioeconomic status  
323 significantly increased their mortality. In contrast, disadvantaged SDoH status was unlikely to  
324 lead to significantly increased mortality risk in Subphenotype I. This evidence further  
325 demonstrated that the COVID-19 pandemic has disproportionately affected patients with lower  
326 socioeconomic status. In general, our findings added new information on social disparities in the  
327 COVID-19 pandemic. Unlike previous studies<sup>29,30,34,35</sup> that focused on the entire population, we  
328 extended the study from a new angle by focusing on the biologically different populations (i.e.,  
329 subphenotypes). Our findings also showed evidence that the identified subphenotypes would  
330 provide considerable guidance in health policy to reduce social disparities in the pandemic.

331

332

## 333 **Limitations**

334 While this study presents a new contribution in the efforts to parse the biological heterogeneity  
335 of COVID-19, there remain several limitations. First of all, our data-driven approach relied on the  
336 availability of patient data. In this study, we identified subphenotypes using the routinely  
337 collected clinical variables that were correlated with COVID-19<sup>36</sup> and available in the INSIGHT  
338 database<sup>37</sup>. We were not able to extract presenting symptoms and vital data while the  
339 incorporation of such data would add in new insights.

340

341 Second, in our study, the analyzed data were collected at ED or hospital presentation, so the  
342 time between COVID-19 symptom onset to ED or hospital presentation could be a covariate of  
343 disease severity and clinical outcomes. However, such data was not available in the INSIGHT  
344 database.

345

346 Third, missing values may affect the robustness of the identified subphenotypes. In order to  
347 address this issue, we excluded variables with high missingness. For the remaining variables,  
348 we used the K-nearest neighbors imputation algorithm<sup>38</sup>. Even so, we still missed these real  
349 values hence may incorporate bias.

350

351 Fourth, our study was based on presenting clinical data, such that each patient was  
352 characterized in a snapshot. The full use of longitudinal data of patients may allow us to capture  
353 the complexity of the disease arc to identify interesting subphenotypes. Previous studies tried to  
354 derive COVID-19 subphenotypes based on longitudinal information<sup>22,24</sup>, yet they were based on  
355 univariate trajectory data in small cohorts. The collection of multivariate, longitudinal data in

356 large cohorts remains challenging and modeling such data to identify subphenotypes requires  
357 improved data-driven methods<sup>12,13,21</sup>.

358  
359 Fifth, this is a multiple institutional analysis in NYC. To evaluate the generalizability of the  
360 identified subphenotypes, further validation on data collected from other areas is needed in  
361 future work.

362  
363  
364  
365

## 366 **Methods**

### 367 **Study design and cohort description**

368 We used data of COVID-19 patients from INSIGHT Clinical Research Network (CRN)<sup>37</sup>.  
369 INSIGHT is funded by the Patient-Centered Outcomes Research Institute (PCORI) and  
370 aggregates clinical data of diverse patient populations across five academic medical centers in  
371 New York City (NYC), including New York University Langone Medical Center (NYU-LMC), New  
372 York Presbyterian - Weill Cornell Medical Center (NYP-WCMC), New York Presbyterian -  
373 Columbia University Medical Center (NYP-CUMC), Mount Sinai Health System (MSHS), and  
374 Montefiore Medical Center (MMC). COVID-19 diagnosis was defined as having at least one  
375 positive laboratory test result for SARS-CoV-2 infection or at least one ICD-10 diagnosis code  
376 for COVID-19 (see [eMethods](#) in Supplement). Study participants were adult patients who were  
377 diagnosed with COVID-19 and treated in ED or inpatient settings in these five health centers  
378 from March 1 to June 12, 2020. Criteria used to assess patient eligibility are illustrated in

379 [eFigure 1](#) in Supplement. Exclusion criteria include younger than 18 years old; duplicated  
380 patient IDs; having no emergency department (ED) or inpatient (IP) admission within 14 days  
381 after COVID-19 confirmation; or having missing values on all clinical variables. Considering the  
382 population diversity of the five medical centers (see [eTable 1](#) in Supplement), we combined  
383 patients of four centers and randomly divided them into the development cohort (70%) and  
384 internal validation cohort (30%). Patients of the last center were used as the external validation  
385 cohort.

386

### 387 **Candidate variables for subphenotype identification**

388 We considered 30 clinical variables associated with COVID-19 onset, symptoms, or outcomes<sup>36</sup>  
389 and available in the INSIGHT database as the candidate variables to derive subphenotypes.  
390 The variables included inflammatory markers (C-reactive protein, erythrocyte sedimentation rate  
391 [ESR], interleukin 6 [IL-6], procalcitonin, bands [i.e., premature neutrophil], lactate  
392 dehydrogenase [LDH], lymphocyte count, neutrophil count, and white blood cell count),  
393 inflammatory and hepatic markers (albumin and ferritin), hepatic markers (alanine  
394 aminotransferase [ALT], aspartate aminotransferase [AST], and bilirubin), markers of  
395 cardiovascular conditions (creatinine kinase [CK], lactate, troponin I, and troponin T), markers of  
396 renal dysfunctions (bicarbonate, blood urea nitrogen [BUN], creatinine, chloride, and sodium),  
397 markers of hematologic dysfunctions (d-dimer, hemoglobin, platelet count, prothrombin time  
398 [PT], red blood cell distribution width [RDW], and glucose), and oxygen saturation. For each  
399 patient, we extracted the first value of each clinical variable within the collection window, which  
400 was defined as: 1) time period from COVID-19 confirmation to the first inpatient encounter, if  
401 the patient has an inpatient admission within 14 days after confirmation; or 2) 14 days after  
402 COVID-19 confirmation if there was only ED encounters but no inpatient admissions following



403 the COVID-19 diagnosis. If there was no record in the collection window, we extracted the last  
404 value within 3 days before confirmation (see [eFigure 2](#) in Supplement).

## 405 **Other clinical characteristics, clinical outcomes, and medications**

406 We also examined other clinical characteristics of the patients, including demographics,  
407 comorbidities, and body mass index (BMI). Demographics included age, sex, and race. Baseline  
408 comorbidities included hypertension, diabetes, coronary artery disease (CAD), heart failure,  
409 chronic obstructive pulmonary disease (COPD), asthma, cancer, obesity, and hyperlipidemia.  
410 For each patient, the most recent BMI data was collected. We analyzed 60-day all-cause  
411 mortality as the primary outcome for the patients. Need for mechanical ventilation and  
412 admission to the intensive care unit (ICU) were the secondary outcomes. We also analyzed the  
413 treatments for COVID-19, including antibiotics (combining ceftriaxone, azithromycin, piperacillin  
414 tazobactam, meropenem, vancomycin, and doxycycline), corticosteroids (combining prednisone,  
415 methylprednisolone, dexamethasone, and hydrocortisone), hydroxychloroquine, enoxaparin,  
416 heparin, and vasopressor. These above data were collected from patient records available in  
417 the INSIGHT database as well.

## 418 **SDoH data**

419 To explore the impact of SDoH status on the subphenotypes, we extracted patients'  
420 neighborhood socioeconomic characteristics, including median household income, percentage  
421 of residents without a high school degree, percentage of residents who are essential workers,  
422 percentage of households with crowding housing conditions (i.e., households with >1 person  
423 per room), percentage of non-white residents, and unemployment rate. These characteristics  
424 were extracted from the 2018 American Community Survey<sup>39</sup>. Previous studies<sup>40-46</sup> have

425 indicated that these social conditions are associated with higher probability of infection,  
426 hospitalization, and other adverse outcomes, e.g., mortality, in COVID-19.

427

## 428 **Statistical methods**

### 429 Data preparation

430 We first assessed the value distributions and missingness of the 30 candidate clinical variables  
431 (see [eTables 2 and 3](#) in Supplement). For data quality control, 7 variables of high missingness  
432 (missing more than 70% values) were excluded and the remaining 23 variables were used for  
433 deriving subphenotypes. Logarithmic transformation was applied to the non-normal distributed  
434 variables (see [eTable 4](#) in Supplement). In order to eliminate the effects of value magnitude, all  
435 variables were scaled based on z-score. K-nearest neighbors (KNN) imputation<sup>38</sup> was used to  
436 address missing values (see [eMethods](#) in Supplement).

437

### 438 Subphenotype derivation, validation, and prediction

439 We originally derived subphenotypes using the development cohort. More specifically,  
440 agglomerative hierarchical clustering with Euclidean distance calculation and Ward linkage  
441 criterion<sup>47</sup> was applied to the 23 clinical variables after data preparation. We used agglomerative  
442 hierarchical clustering because it is robust to different types of data distributions and typically  
443 produces a dendrogram that visualizes data structure to help determine the optimal cluster  
444 number. Besides dendrogram, we calculated 21 measures of clustering models provided by  
445 'NbClust' software<sup>48</sup> to determine the optimal number of clusters, i.e., subphenotypes.

446

447 In order to evaluate the reproducibility, we validated our subphenotypes in four ways. First, we  
448 performed sensitivity analyses using the development cohort to evaluate 1) sensitivity to quality  
449 control and outliers and 2) sensitivity to clustering algorithms. To assess sensitivity to quality  
450 control and outliers, we incorporated all 30 candidate variables and excluded patients who have  
451 outlier values (see [eMethods](#) in Supplement). Then similar to the primary analysis, we  
452 performed agglomerative hierarchical clustering to re-derive subphenotypes and determined  
453 optimal cluster number using dendrogram and 'NbClust'. To assess sensitivity to clustering  
454 algorithms, we re-derived subphenotypes using the Gaussian mixture model (GMM)<sup>49</sup>, which is  
455 a probabilistic model for clustering analysis based on a mixture of Gaussian distributions. The  
456 optimal cluster number in GMM was determined by comprehensively considering Akaike  
457 information criterion (AIC), Bayesian information criterion (BIC), and median probability of group  
458 membership (see [eMethods](#) in Supplement).

459  
460 Second, we used the internal validation cohort and re-derived subphenotypes using the same  
461 agglomerative hierarchical clustering with the primary analysis for validation. The optimal cluster  
462 number was determined using dendrogram and 'NbClust' as well.

463  
464 Third, for the aims of confirming subphenotypes and their usability, we used the supervised  
465 predictive model. More specifically, considering subphenotype membership of each patient as  
466 the label to predict, we built a predictive model of subphenotypes based on the 23 clinical  
467 variables used for subphenotype derivation. The predictive model was based on the supervised  
468 XGBoost classifier<sup>50</sup>, a powerful tree-based machine learning model. The predictive model was  
469 trained in the development cohort using a 10-fold cross-validation strategy. To address the  
470 multi-label classification (since we identified more than 2 subphenotypes), a one-vs-the-rest  
471 strategy was used in model training. Prediction performance was measured by receiver  
472 operating characteristics curve (ROC) and area under ROC curve (AUC). We also engaged the

473 SHapley Additive exPlanation (SHAP) values to assess contributions of the clinical variables in  
474 distinguishing each subphenotype from the others. Once the predictive model was trained, it  
475 was performed on the external validation cohort to predict the patients' subphenotype  
476 memberships.

477  
478 Last, to assess stability of the subphenotypes across the five medical centers, we further  
479 performed leave-one-center-out analysis (see [eMethods](#) in Supplement).

480  
481 *Subphenotype interpretation*  
482 For the aim of subphenotype interpretation, we first visualized the subphenotypes in two ways: 1)  
483 2-D visualization calculated by Uniform Manifold Approximation and Projection (UMAP)  
484 algorithm<sup>51</sup> based on clinical variables for clustering (showing distributions of subphenotypes  
485 within low-dimensional space); 2) chord diagrams<sup>52</sup> showing differences of subphenotypes in  
486 terms of abnormal clinical variable groups and comorbidities (see [eMethods](#) in Supplement).

487  
488 We also characterized subphenotypes by evaluating their differences in demographics, all  
489 clinical variables, comorbidities, clinical outcomes, and medications prescribed after COVID-19  
490 confirmation. Data were presented as median (interquartile range [IQR]) for continuous  
491 variables and exact patient number (percentage) for categorical variables. To compare  
492 subphenotypes, we performed the Kruskal-Wallis test for continuous data and  $\chi^2$  test for  
493 categorical data. Analysis of covariance (ANCOVA) was also applied for between-  
494 subphenotypes comparisons, adjusting for age and gender. Two-tailed P-values smaller than  
495 0.05 were considered as the threshold for statistical significance. Survival analyses were  
496 performed to assess associations of subphenotypes to clinical outcomes, where Kaplan-Meier  
497 plots were created accordingly.

498

499 *Temporal pattern of subphenotypes*

500 To evaluate the temporal pattern of the subphenotypes during the course of the pandemic, we  
501 created bar charts to visualize the proportion of each subphenotype out of the total patients  
502 confirmed per week, since the COVID-19 outbreak in NYC (March 1, 2020).

503

504 *Impacts of SDoH on COVID-19 subphenotypes*

505 Multiple analyses were conducted to assess the impact of SDoH on COVID-19 subphenotypes.  
506 For each subphenotype, we first performed logistic regression analysis and Cox regression  
507 analysis to assess the association of each SDoH variable with 60-day mortality, adjusting for  
508 age, sex, and/or clinical variables. After that, we performed agglomerative hierarchical clustering  
509 on the 6 socioeconomic variables to derive comprehensive SDoH strata. Within each  
510 subphenotype, we compared 60-day mortality rates between the SDoH strata. We also used  
511 logistic regression analysis and Cox regression analysis to assess the association of SDoH  
512 strata with 60-day mortality, adjusting for age and sex, within each subphenotype.

513

514 **Ethical approval and patient consent.**

515 The Institutional Review Board of the Weill Cornell Medicine approved this study (Protocol  
516 number: 20-04021948).

517

518

519

## 520 **Data availability**

521 All data studied in this work can be downloaded from INSIGHT clinical research network  
522 at <https://insightcrn.org/our-data/>, via request.

523

## 524 **Code availability**

525 All computer codes in this study are available at [https://github.com/ChangSu10/COVID-](https://github.com/ChangSu10/COVID-Insight-subphenotyping)  
526 [Insight-subphenotyping](https://github.com/ChangSu10/COVID-Insight-subphenotyping). Implementation of our work is based on Python 3.7 and R 3.6.

527 More specifically, clustering models were implemented based on Python packages

528 ‘scikit-learn 0.23.2’ (<https://scikit-learn.org/stable/>) and ‘scipy 1.5.3’

529 (<https://www.scipy.org>). Supervised predictive modeling was based on ‘XGBoost 1.2.1’

530 (<https://xgboost.readthedocs.io/en/latest/>) and ‘SHAP 0.35.0’

531 (<https://shap.readthedocs.io/en/latest/>). Data dimension reduction and visualization

532 were performed based on Python package ‘UMAP-learn 0.3.9’ ([https://umap-](https://umap-learn.readthedocs.io/en/latest/)

533 [learn.readthedocs.io/en/latest/](https://umap-learn.readthedocs.io/en/latest/)). R package ‘NbClust’ ([534 \[project.org/web/packages/NbClust/NbClust.pdf\]\(https://cran.r-project.org/web/packages/NbClust/NbClust.pdf\)\) was used to calculate measures of](https://cran.r-</a></p></div><div data-bbox=)

535 clusters to determine the optimal cluster number in agglomerative hierarchical clustering.

536 Chord diagrams were created using R package ‘circlize’ ([537 \[project.org/web/packages/circlize/index.html\]\(https://cran.r-project.org/web/packages/circlize/index.html\)\). All statistical tests and survival analyses](https://cran.r-</a></p></div><div data-bbox=)

538 were performed based on R.

539

540

## 541 **Acknowledgements**

542 This study is funded by the COVID-19-Related Project Enhancement to the grant PCORI/HSD-  
543 1604-35187 (“Identifying and Predicting Patients with Preventable High Utilization”, PI: Kaushal)  
544 from the Patient-Centered Outcomes Research Institute.

545

546

## 547 **Competing interests**

548 The authors have declared that no conflict of interest exists.

549

550

## 551 **Author contributions**

552 E.S. and F.W. for conceptualization, investigation, writing, reviewing and editing of the  
553 manuscript. C.S. for investigation, data analysis, drafting, editing and reviewing manuscript.  
554 M.G.W. for providing data support, discussion and commenting the manuscript. Y.Z., J.H.F.,  
555 and R.K. for discussion, commenting and editing the manuscript.

556

## 557 References

- 558 1. Zhu, N., *et al.* A Novel Coronavirus from Patients with Pneumonia in China, 2019. *N.*  
559 *Engl. J. Med.* **382**, 727-733 (2020).
- 560 2. Richardson, S., *et al.* Presenting Characteristics, Comorbidities, and Outcomes Among  
561 5700 Patients Hospitalized With COVID-19 in the New York City Area. *JAMA* **323**, 2052-  
562 2059 (2020).
- 563 3. Tabata, S., *et al.* Clinical characteristics of COVID-19 in 104 people with SARS-CoV-2  
564 infection on the Diamond Princess cruise ship: a retrospective analysis. *The Lancet*  
565 *Infectious Diseases* **20**, 1043-1050 (2020).
- 566 4. Desai, N., *et al.* Temporal and spatial heterogeneity of host response to SARS-CoV-2  
567 pulmonary infection. *Nature Communications* **11**, 6319 (2020).
- 568 5. Wiersinga, W.J., Rhodes, A., Cheng, A.C., Peacock, S.J. & Prescott, H.C.  
569 Pathophysiology, Transmission, Diagnosis, and Treatment of Coronavirus Disease 2019  
570 (COVID-19): A Review. *JAMA* **324**, 782-793 (2020).
- 571 6. Takahashi, T., *et al.* Sex differences in immune responses that underlie COVID-19  
572 disease outcomes. *Nature* **588**, 315-320 (2020).
- 573 7. Gupta, S., *et al.* Association Between Early Treatment With Tocilizumab and Mortality  
574 Among Critically Ill Patients With COVID-19. *JAMA Internal Medicine* **181**, 41-51 (2021).
- 575 8. Domecq, J.P., *et al.* Outcomes of Patients With Coronavirus Disease 2019 Receiving  
576 Organ Support Therapies: The International Viral Infection and Respiratory Illness  
577 Universal Study Registry. *Crit. Care Med.* **Online First**(9000).
- 578 9. Dexamethasone in Hospitalized Patients with Covid-19 — Preliminary Report. *N. Engl. J.*  
579 *Med.* (2020).
- 580 10. Kreuziger, L.B., *et al.* COVID-19 and VTE/Anticoagulation: Frequently Asked Questions.  
581 Vol. 2021 (The American Society of Hematology, 2021).
- 582 11. Full-dose blood thinners decreased need for life support and improved outcome in  
583 hospitalized COVID-19 patients. Vol. 2021 (National Institutes of Health (NIH), 2021).
- 584 12. Weng, C., Shah, N.H. & Hripcsak, G. Deep phenotyping: Embracing complexity and  
585 temporality-Towards scalability, portability, and interoperability. *J. Biomed. Inform.* **105**,  
586 103433 (2020).
- 587 13. Reddy, K., *et al.* Subphenotypes in critical care: translation into clinical practice. *The*  
588 *Lancet Respiratory Medicine* **8**, 631-643 (2020).
- 589 14. Seymour, C.W., *et al.* Derivation, Validation, and Potential Treatment Implications of  
590 Novel Clinical Phenotypes for Sepsis. *JAMA* **321**, 2003-2017 (2019).
- 591 15. Bhavani, S.V., *et al.* Identifying Novel Sepsis Subphenotypes Using Temperature  
592 Trajectories. *Am. J. Respir. Crit. Care Med.* **200**, 327-335 (2019).
- 593 16. Calfee, C.S., *et al.* Subphenotypes in acute respiratory distress syndrome: latent class  
594 analysis of data from two randomised controlled trials. *The Lancet Respiratory Medicine*  
595 **2**, 611-620 (2014).
- 596 17. Ahmad, T., *et al.* Clinical Implications of Chronic Heart Failure Phenotypes  
597 Defined by Cluster Analysis. *J. Am. Coll. Cardiol.* **64**, 1765-1774 (2014).
- 598 18. Cikes, M., *et al.* Machine learning-based phenogrouping in heart failure to identify  
599 responders to cardiac resynchronization therapy. *Eur. J. Heart Fail.* **21**, 74-85 (2019).
- 600 19. Li, L., *et al.* Identification of type 2 diabetes subgroups through topological analysis of  
601 patient similarity. *Sci. Transl. Med.* **7**, 311ra174 (2015).
- 602 20. Neff, R.A., *et al.* Molecular subtyping of Alzheimer's disease using RNA sequencing data  
603 reveals novel mechanisms and targets. *Science Advances* **7**, eabb5398 (2021).



- 604 21. Bos, L.D.J., Sinha, P. & Dickson, R.P. The perils of premature phenotyping in COVID-19:  
605 a call for caution. *Eur. Respir. J.* **56**, 2001768 (2020).
- 606 22. Bhavani, S.V., Huang, E.S., Verhoef, P.A. & Churpek, M.M. Novel Temperature  
607 Trajectory Subphenotypes in COVID-19. *Chest* **158**, 2436-2439 (2020).
- 608 23. Legrand, M., *et al.* Differences in clinical deterioration among three sub-phenotypes of  
609 COVID-19 patients at the time of first positive test: results from a clustering analysis.  
610 *Intensive Care Med.* **47**, 113-115 (2021).
- 611 24. Su, C., *et al.* Identifying organ dysfunction trajectory-based subphenotypes in critically ill  
612 patients with COVID-19. *medRxiv*, 2020.2007.2016.20155382 (2020).
- 613 25. Rodríguez, A., *et al.* Deploying unsupervised clustering analysis to derive clinical  
614 phenotypes and risk factors associated with mortality risk in 2022 critically ill patients  
615 with COVID-19 in Spain. *Critical Care* **25**, 63 (2021).
- 616 26. Asch, D.A., *et al.* Variation in US Hospital Mortality Rates for Patients Admitted With  
617 COVID-19 During the First 6 Months of the Pandemic. *JAMA Internal Medicine* (2020).
- 618 27. Jorge, A., *et al.* Temporal trends in severe COVID-19 outcomes in patients with  
619 rheumatic disease: a cohort study. *The Lancet Rheumatology* **3**, e131-e137 (2021).
- 620 28. Gray, D.M., Anyane-Yeboah, A., Balzora, S., Issaka, R.B. & May, F.P. COVID-19 and the  
621 other pandemic: populations made vulnerable by systemic inequity. *Nature Reviews*  
622 *Gastroenterology & Hepatology* **17**, 520-522 (2020).
- 623 29. Wadhwa, R.K., *et al.* Variation in COVID-19 Hospitalizations and Deaths Across New  
624 York City Boroughs. *JAMA* **323**, 2192-2195 (2020).
- 625 30. Azar, K.M.J., *et al.* Disparities In Outcomes Among COVID-19 Patients In A Large  
626 Health Care System In California. *Health Aff. (Millwood)* **39**, 1253-1262 (2020).
- 627 31. Bravata, D.M., *et al.* Association of Intensive Care Unit Patient Load and Demand With  
628 Mortality Rates in US Department of Veterans Affairs Hospitals During the COVID-19  
629 Pandemic. *JAMA Netw. Open* **4**, e2034266-e2034266 (2021).
- 630 32. Knox, D.B., Lanspa, M.J., Kuttler, K.G., Brewer, S.C. & Brown, S.M. Phenotypic clusters  
631 within sepsis-associated multiple organ dysfunction syndrome. *Intensive Care Med.* **41**,  
632 814-822 (2015).
- 633 33. Thompson, C.N., *et al.* COVID-19 Outbreak - New York City, February 29-June 1, 2020.  
634 *Morbidity and mortality weekly report (MMWR)* **69**, 1725-1729 (2020).
- 635 34. Wang, Z., *et al.* Analysis of hospitalized COVID-19 patients in the Mount Sinai Health  
636 System using electronic medical records (EMR) reveals important prognostic factors for  
637 improved clinical outcomes. *medRxiv*, 2020.2004.2028.20075788 (2020).
- 638 35. Federgruen, A. & Naha, S. Variation in Covid-19 Cases Across New York City. *medRxiv*,  
639 2020.2005.2025.20112797 (2020).
- 640 36. Wynants, L., *et al.* Prediction models for diagnosis and prognosis of covid-19: systematic  
641 review and critical appraisal. *BMJ* **369**, m1328 (2020).
- 642 37. Kaushal, R., *et al.* Changing the research landscape: the New York City Clinical Data  
643 Research Network. *J. Am. Med. Inform. Assoc.* **21**, 587-590 (2014).
- 644 38. Troyanskaya, O., *et al.* Missing value estimation methods for DNA microarrays.  
645 *Bioinformatics* **17**, 520-525 (2001).
- 646 39. American Community Survey: 2018 Data Release New and Notable.
- 647 40. Kolak, M., Bhatt, J., Park, Y.H., Padrón, N.A. & Molefe, A. Quantification of  
648 Neighborhood-Level Social Determinants of Health in the Continental United States.  
649 *JAMA Netw. Open* **3**, e1919928-e1919928 (2020).
- 650 41. Whittle, R.S. & Diaz-Artiles, A. An ecological study of socioeconomic predictors in  
651 detection of COVID-19 cases across neighborhoods in New York City. *BMC Med.* **18**,  
652 271 (2020).

- 653 42. von Seidlein, L., Alabaster, G., Deen, J. & Knudsen, J. Crowding has consequences:  
654 Prevention and management of COVID-19 in informal urban settlements. *Build. Environ.*  
655 **188**, 107472 (2021).
- 656 43. Hawkins, R.B., Charles, E.J. & Mehaffey, J.H. Socio-economic status and COVID-19–  
657 related cases and fatalities. *Public Health* **189**, 129-134 (2020).
- 658 44. Lieberman-Cribbin, W., Tuminello, S., Flores, R.M. & Taioli, E. Disparities in COVID-19  
659 Testing and Positivity in New York City. *Am. J. Prev. Med.* **59**, 326-332 (2020).
- 660 45. Do, D.P. & Frank, R. Unequal burdens: assessing the determinants of elevated COVID-  
661 19 case and death rates in New York City’s racial/ethnic minority neighbourhoods. *J.*  
662 *Epidemiol. Community Health* **75**, 321 (2021).
- 663 46. Hong, B., Bonczak, B.J., Gupta, A., Thorpe, L.E. & Kontokosta, C.E. Exposure density  
664 and neighborhood disparities in COVID-19 infection risk. *Proc. Natl. Acad. Sci. U. S. A.*  
665 **118**, e2021258118 (2021).
- 666 47. Murtagh, F. & Legendre, P. Ward’s Hierarchical Agglomerative Clustering Method:  
667 Which Algorithms Implement Ward’s Criterion? *Journal of Classification* **31**, 274-295  
668 (2014).
- 669 48. Charrad, M., Ghazzali, N., Boiteau, V. & Niknafs, A. NbClust: An R Package for  
670 Determining the Relevant Number of Clusters in a Data Set. *Journal of Statistical*  
671 *Software* **61**(2014).
- 672 49. Reynolds, D.A. Gaussian Mixture Models. *Encyclopedia of biometrics* **741**, 659-663  
673 (2009).
- 674 50. Chen, T. & Guestrin, C. XGBoost: A Scalable Tree Boosting System. in *Proceedings of*  
675 *the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data*  
676 *Mining* 785–794 (Association for Computing Machinery, San Francisco, California, USA,  
677 2016).
- 678 51. McInnes, L., Healy, J. & Melville, J. Umap: Uniform manifold approximation and  
679 projection for dimension reduction. *arXiv preprint arXiv:1802.03426* (2018).
- 680 52. Gu, Z., Gu, L., Eils, R., Schlesner, M. & Brors, B. circlize implements and enhances  
681 circular visualization in R. *Bioinformatics* **30**, 2811-2812 (2014).
- 682

## Tables

**Table 1. Characteristics of the development, internal validation, and external validation cohorts**

Characteristics	Cohort		
	Development cohort	Internal validation cohort	External validation cohort
No. of patients	8,199	3,519	2,700
Construction method	70% patients (randomly selected) from 4 medical centers	Remaining 30% patients from 4 medical centers	Patients from the last center
Age, y, Median (IQR)	63.53 [50.57 - 75.15]	63.51 [50.95 - 75.17]	65.58 (51.08 - 77.39)
Sex female, N (%)	3,787 (46.2)	1,585 (45.0)	1,305 (48.3)
Race, N (%)			
White	2,036 (24.8)	838 (23.8)	675 (25.0)
Black	2,155 (26.3)	915 (26.0)	545 (20.2)
Asian	409 (5.0)	193 (5.5)	28 (1.0)
Other/unknown	3599 (43.9)	1573 (44.7)	1452 (53.8)
Outcomes (60 days), N (%)			
Mortality	1529 (18.65)	696 (19.78)	556 (20.59)
Mechanical ventilation (intubation)	1154 (14.07)	497 (14.12)	248 (9.19)
ICU admission	1494 (18.22)	661 (18.78)	-

Abbreviations: ICU = intensive care unit; IQR = Interquartile range; SDoH = social determinants of health.

**Table 2. Characteristics of the identified subphenotypes (development cohort)**

Variable	Total	Subphenotype I	Subphenotype II	Subphenotype III	Subphenotype IV	P-value <sup>1</sup>	P-value (age and sex adjusted) <sup>2</sup>
No. of patients (%)	8199 (100)	2707 (33.02)	3047 (37.16)	1486 (18.12)	959 (11.70)	-	-
<b>Demographics</b>							
Age, y, Median (IQR)	63.53 (50.57 - 75.15)	57.45 (42.70 - 70.02)	62.56 (51.63 - 72.77)	69.45 (57.05 - 79.62)	73.53 (64.10 - 82.83)	< 0.001	-
Sex female, N (%)	3787 (46.19)	1601 (59.14)	1106 (36.30)	709 (47.71)	371 (38.69)	< 0.001	-
<b>Race, N (%)</b>							
White	2036 (24.83)	695 (25.67)	777 (25.50)	367 (24.70)	197 (20.54)	< 0.001	-
Black	2155 (26.28)	697 (25.75)	611 (20.05)	503 (33.85)	344 (35.87)		
Asian	409 (4.99)	118 (4.36)	194 (6.37)	58 (3.90)	39 (4.07)		
Other/unknown	3599 (43.90)	1197 (44.22)	1465 (48.08)	558 (37.55)	379 (39.52)		
<b>Inflammatory markers</b>							
C-reactive protein, mg/L, Median (IQR)	9.40 (3.70 - 16.80)	4.32 (1.16 - 9.31)	12.74 (6.60 - 20.20)	8.20 (3.50 - 14.51)	14.90 (6.70 - 23.07)		
ESR, mm/hr, Median (IQR)	69.00 (42.00 - 97.00)	53.00 (34.00 - 81.00)	76.00 (50.00 - 100.00)	75.00 (45.25 - 102.75)	77.00 (41.75 - 106.25)	< 0.001	< 0.001
IL-6, pg/mL, Median (IQR)	19.00 (10.00 - 42.00)	13.00 (8.00 - 21.00)	21.00 (11.00 - 45.75)	17.00 (9.00 - 47.00)	27.00 (10.25 - 52.00)	< 0.001	< 0.001
Procalcitonin, ng/mL, Median (IQR)	0.20 (0.10 - 0.60)	0.10 (0.10 - 0.20)	0.20 (0.10 - 0.50)	0.30 (0.10 - 0.87)	0.60 (0.25 - 2.10)	< 0.001	0.26
Bands, %, Median (IQR)	2.00 (0.00 - 5.00)	3.00 (0.00 - 5.75)	2.00 (0.00 - 5.00)	2.00 (0.00 - 5.00)	2.00 (0.00 - 6.00)	< 0.001	0.04
LDH, U/L, Median (IQR)	377.00 (280.00 - 525.00)	292.00 (229.00 - 377.00)	437.00 (343.00 - 576.00)	349.00 (268.00 - 449.00)	565.50 (409.75 - 801.50)	0.37	0.14
Lymphocyte count, ×10 <sup>3</sup> /uL, Median (IQR)	1.00 (0.70 - 1.43)	1.20 (0.80 - 1.60)	1.00 (0.70 - 1.40)	0.80 (0.60 - 1.20)	0.90 (0.60 - 1.40)	< 0.001	< 0.001
Neutrophil count, ×10 <sup>3</sup> /uL, Median (IQR)	5.30 (3.70 - 7.90)	4.00 (2.90 - 5.40)	6.70 (4.80 - 9.50)	4.70 (3.40 - 6.60)	8.20 (5.90 - 11.00)	< 0.001	0.02
White blood cell count, ×10 <sup>3</sup> /uL, Median (IQR)	7.20 (5.30 - 9.90)	5.90 (4.60 - 7.60)	8.50 (6.50 - 11.50)	6.30 (4.70 - 8.30)	10.30 (7.60 - 13.57)	< 0.001	< 0.001
<b>Inflammation &amp; Hepatic markers</b>							
Albumin, g/dL, Median (IQR)	3.70 (3.30 - 4.10)	4.00 (3.60 - 4.30)	3.70 (3.20 - 4.00)	3.50 (3.10 - 3.90)	3.40 (2.90 - 3.80)		
Ferritin, ng/mL, Median (IQR)	645.00 (295.90 - 1347.00)	323.05 (157.75 - 594.33)	868.80 (454.00 - 1537.50)	599.00 (217.80 - 1380.50)	1174.00 (523.00 - 2284.00)	< 0.001	< 0.001
<b>Hepatic markers</b>							
Alanine aminotransferase, U/L, Median (IQR)	29.00 (19.00 - 48.00)	24.00 (17.00 - 36.00)	41.00 (26.00 - 68.00)	20.00 (13.00 - 29.00)	37.00 (22.00 - 65.00)	< 0.001	< 0.001
Aspartate aminotransferase, U/L, Median (IQR)	39.00 (26.00 - 63.00)	31.00 (23.00 - 42.00)	52.00 (35.00 - 80.00)	31.00 (22.00 - 46.00)	65.00 (36.00 - 118.00)	< 0.001	< 0.001
Bilirubin, mg/dL, Median (IQR)	0.30 (0.20 - 0.60)	0.20 (0.20 - 0.40)	0.40 (0.20 - 0.70)	0.30 (0.20 - 0.50)	0.40 (0.20 - 0.70)	< 0.001	< 0.001
<b>Cardiovascular markers</b>							
Creatine kinase, U/L, Median (IQR)	154.00 (78.00 - 359.00)	122.00 (72.00 - 227.00)	165.00 (83.00 - 387.50)	126.00 (63.00 - 288.00)	352.00 (137.00 - 1039.50)		
Lactate, mmol/L, Median (IQR)	1.90 (1.40 - 2.60)	1.50 (1.20 - 2.10)	2.00 (1.50 - 2.70)	1.60 (1.20 - 2.10)	3.10 (2.20 - 4.80)	< 0.001	< 0.001
Troponin I, ng/mL, Median (IQR)	0.10 (0.06 - 0.30)	0.10 (0.00 - 0.10)	0.10 (0.06 - 0.30)	0.10 (0.10 - 0.21)	0.20 (0.10 - 0.50)	< 0.001	< 0.001
Troponin T, ng/mL, Median (IQR)	0.01 (0.01 - 0.03)	0.01 (0.01 - 0.01)	0.01 (0.01 - 0.01)	0.03 (0.01 - 0.09)	0.05 (0.01 - 0.14)	< 0.001	0.16
<b>Renal markers</b>							
Bicarbonate, mmol/L, Median (IQR)	23.00 (21.00 - 26.00)	24.00 (22.00 - 27.00)	23.00 (21.00 - 25.00)	23.00 (20.00 - 25.00)	20.00 (17.00 - 23.00)		
BUN, mg/dL, Median (IQR)	17.00 (11.00 - 31.00)	12.00 (9.00 - 17.00)	16.00 (12.00 - 24.00)	31.00 (18.00 - 53.00)	52.00 (32.00 - 84.00)	< 0.001	< 0.001
Creatinine, mg/dL, Median (IQR)	1.00 (0.80 - 1.50)	0.86 (0.70 - 1.04)	1.00 (0.80 - 1.29)	1.70 (1.00 - 4.40)	2.10 (1.38 - 3.60)	< 0.001	< 0.001
Chloride, mmol/L, Median (IQR)	100.00 (97.00 - 104.00)	101.00 (98.00 - 104.00)	99.00 (95.00 - 102.00)	101.00 (97.00 - 105.00)	104.00 (98.00 - 113.00)	< 0.001	< 0.001
Sodium, mmol/L, Median (IQR)	137.00 (134.00 - 140.00)	138.00 (136.00 - 140.00)	136.00 (132.00 - 138.00)	138.00 (134.00 - 141.00)	141.00 (136.00 - 152.00)	< 0.001	< 0.001

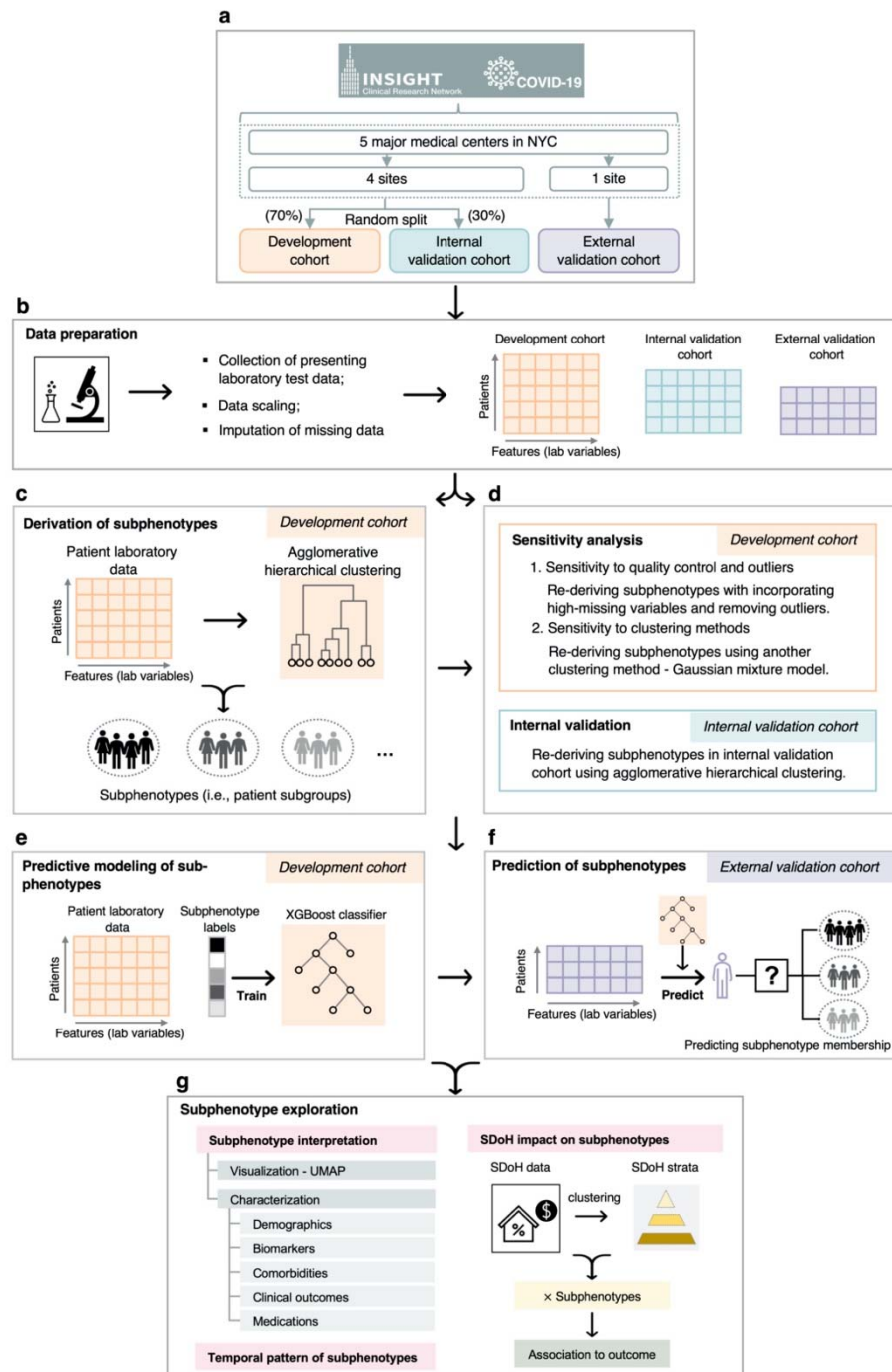
<b>Hematologic markers</b>						< 0.001	< 0.001
D-dimer, ng/mL, Median (IQR)	1360.00 (620.00 - 3370.00)	660.00 (370.00 - 1310.00)	1390.00 (690.00 - 3210.00)	1740.00 (836.50 - 3520.00)	4000.00 (2000.00 - 13582.50)		
Hemoglobin, g/dL, Median (IQR)	13.10 (11.50 - 14.60)	13.40 (12.30 - 14.60)	13.80 (12.50 - 15.10)	10.80 (9.00 - 12.30)	12.75 (10.70 - 15.00)	< 0.001	< 0.001
Platelet count, × 10 <sup>3</sup> /uL, Median (IQR)	211.00 (162.00 - 277.00)	204.00 (163.00 - 253.00)	225.00 (172.00 - 303.00)	194.00 (145.00 - 270.00)	217.00 (156.00 - 296.00)	< 0.001	< 0.001
Prothrombin time, s, Median (IQR)	13.30 (12.20 - 14.60)	12.70 (11.90 - 13.60)	13.50 (12.50 - 14.70)	13.20 (12.00 - 14.60)	14.80 (13.15 - 20.55)	< 0.001	< 0.001
Red blood cell distribution width, %, Median (IQR)	13.80 (12.90 - 15.00)	13.40 (12.80 - 14.40)	13.40 (12.70 - 14.20)	15.50 (14.00 - 17.50)	15.10 (13.80 - 16.70)	< 0.001	< 0.001
Glucose, mg/dL, Median (IQR)	121.00 (101.00 - 165.00)	108.00 (95.00 - 127.00)	133.00 (110.00 - 201.00)	117.00 (98.00 - 153.00)	164.00 (119.00 - 271.75)	< 0.001	< 0.001
<b>Other markers</b>						< 0.001	< 0.001
Oxygen saturation, %, Median (IQR)	69.00 (50.00 - 85.00)	65.00 (47.00 - 85.00)	69.00 (51.50 - 85.00)	69.00 (48.00 - 80.00)	76.50 (57.75 - 91.20)		
BMI, kg/m <sup>2</sup> , Median (IQR)	28.00 (25.00 - 33.00)	29.00 (25.00 - 34.00)	28.95 (25.00 - 33.00)	27.00 (23.00 - 32.00)	26.00 (23.00 - 31.00)	0.05	0.06
						< 0.001	0.73
<b>Comorbidity, (missing=590), N (%)</b>							
Hypertension	4744 (62.35)	1238 (49.68)	1696 (60.16)	1095 (78.44)	715 (79.27)	< 0.001	-
Diabetes	3104 (40.79)	666 (26.73)	1198 (42.50)	730 (52.29)	510 (56.54)	< 0.001	-
Coronary artery disease	1753 (23.04)	360 (14.45)	530 (18.80)	523 (37.46)	340 (37.69)	< 0.001	-
Heart failure	1132 (14.88)	176 (7.06)	286 (10.15)	430 (30.80)	240 (26.61)	< 0.001	-
COPD	972 (12.77)	264 (10.59)	259 (9.19)	290 (20.77)	159 (17.63)	< 0.001	-
Asthma	1091 (14.34)	392 (15.73)	372 (13.20)	232 (16.62)	95 (10.53)	< 0.001	-
Cancer	1438 (18.90)	363 (14.57)	444 (15.75)	423 (30.30)	208 (23.06)	< 0.001	-
Hyperlipidemia	3262 (42.87)	825 (33.11)	1169 (41.47)	779 (55.80)	489 (54.21)	< 0.001	-
Obesity	3039 (37.07)	1105 (40.82)	1179 (38.69)	495 (33.31)	260 (27.11)	< 0.001	-
<b>Outcomes (60 days), N (%)</b>							
Mortality	1529 (18.65)	188 (6.94)	528 (17.33)	337 (22.68)	476 (49.64)	< 0.001	-
Mechanical ventilation (intubation)	1154 (14.07)	190 (7.02)	527 (17.30)	195 (13.12)	242 (25.23)	< 0.001	-
ICU admission	1494 (18.22)	242 (8.94)	675 (22.15)	242 (16.29)	335 (34.93)	< 0.001	-
<b>Medications, N (%)</b>							
Antibiotics	2952 (36.00)	731 (27.00)	1219 (40.01)	559 (37.62)	443 (46.19)	< 0.001	-
Corticosteroids	1666 (20.32)	331 (12.23)	725 (23.79)	319 (21.47)	291 (30.34)	< 0.001	-
Enoxaparin	3312 (40.40)	1016 (37.53)	1582 (51.92)	418 (28.13)	296 (30.87)	< 0.001	-
Heparin	1310 (15.98)	255 (9.42)	585 (19.20)	304 (20.46)	166 (17.31)	< 0.001	-
Vasopressor	608 (7.42)	120 (4.43)	308 (10.11)	96 (6.46)	84 (8.76)	< 0.001	-

Abbreviations: BUN = blood urea nitrogen; COPD = chronic obstructive pulmonary disease; ESR = Erythrocyte sedimentation rate; ICU = intensive care unit; IL-6 = Interleukin 6; IQR = Interquartile range; LDH = Lactate dehydrogenase.

<sup>1</sup> Comparisons across all 4 subphenotypes were performed using the Kruskal-Wallis test (with Dunn's test for post-hoc pairwise comparisons) or  $\chi^2$  test.

<sup>2</sup> P-values, adjusting for age and sex, were calculated by analysis of covariance (ANCOVA) was performed based on General Linear Model.

## Figures

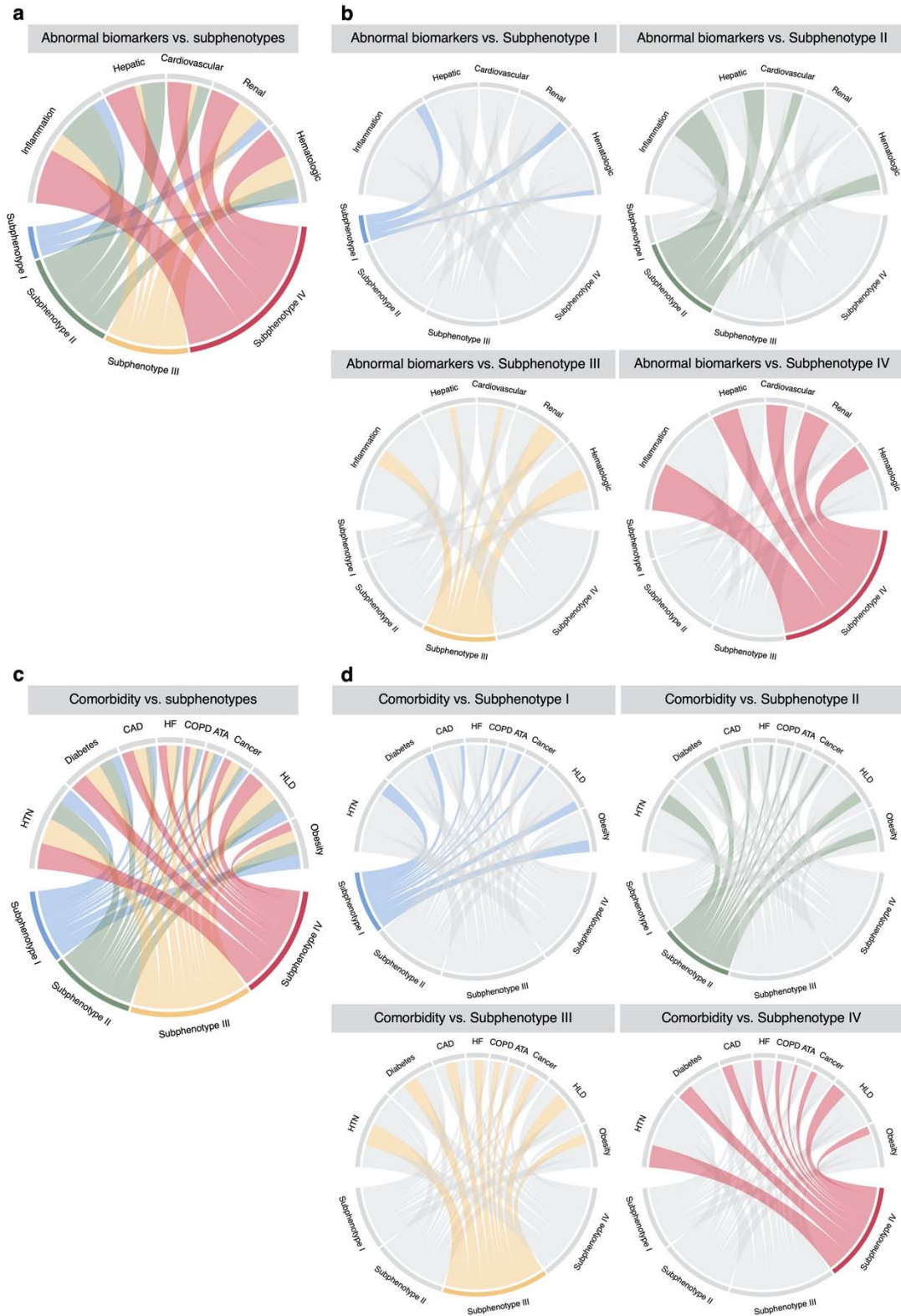


**Figure 1. A schematic of the analysis plan.**

**a.** Strategy for construction of development, internal validation, and external validation cohorts. Studied patients were treated in 5 major medical centers in New York City, including New York University Langone Medical Center, New York Presbyterian - Weill Cornell Medical Center, New

York Presbyterian - Columbia University Medical Center, Mount Sinai Health System, and Montefiore Medical Center. **b.** Data preparation for clustering analysis. **c.** Derivation of subphenotypes in the development cohort. Reproducibility of the identified subphenotypes were evaluated in multiple ways, including **(d)** sensitivity analyses in the development cohort and subphenotype re-derivation in the internal validation cohort; and **(e)** training subphenotype predictive model in the development cohort and **(f)** using it to predict subphenotype memberships of patients in the external validation cohort. Last, **(g)** further analyses were conducted to interpret subphenotypes, explore temporal patterns of subphenotypes during the pandemic, and evaluate impact of SDoH characteristics on subphenotypes.

Abbreviations: NYC = New York City; SDoH = social determinants of health; UMAP = Uniform Manifold Approximation and Projection

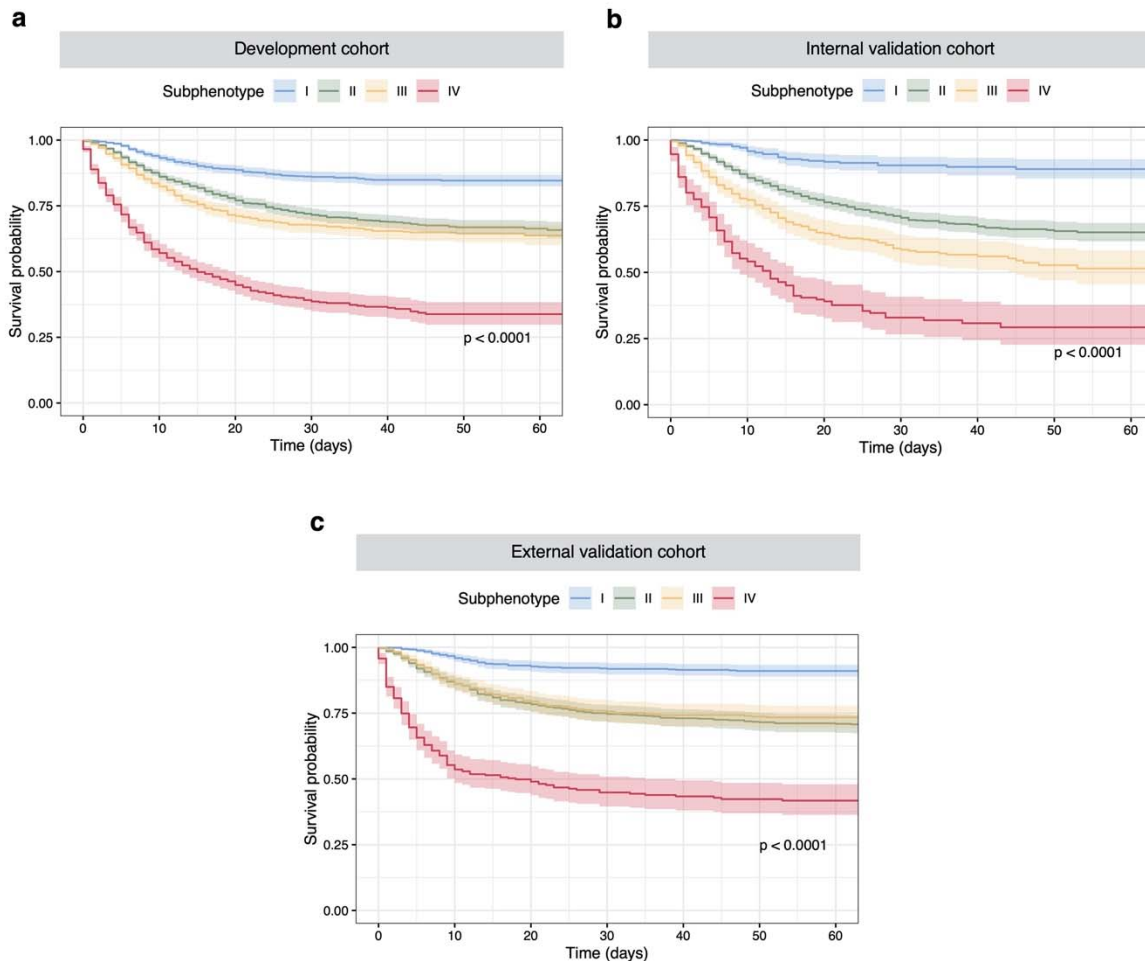


**Figure 2. Chord diagrams showing differences in abnormal clinical variables and comorbidity burden among subphenotypes.**

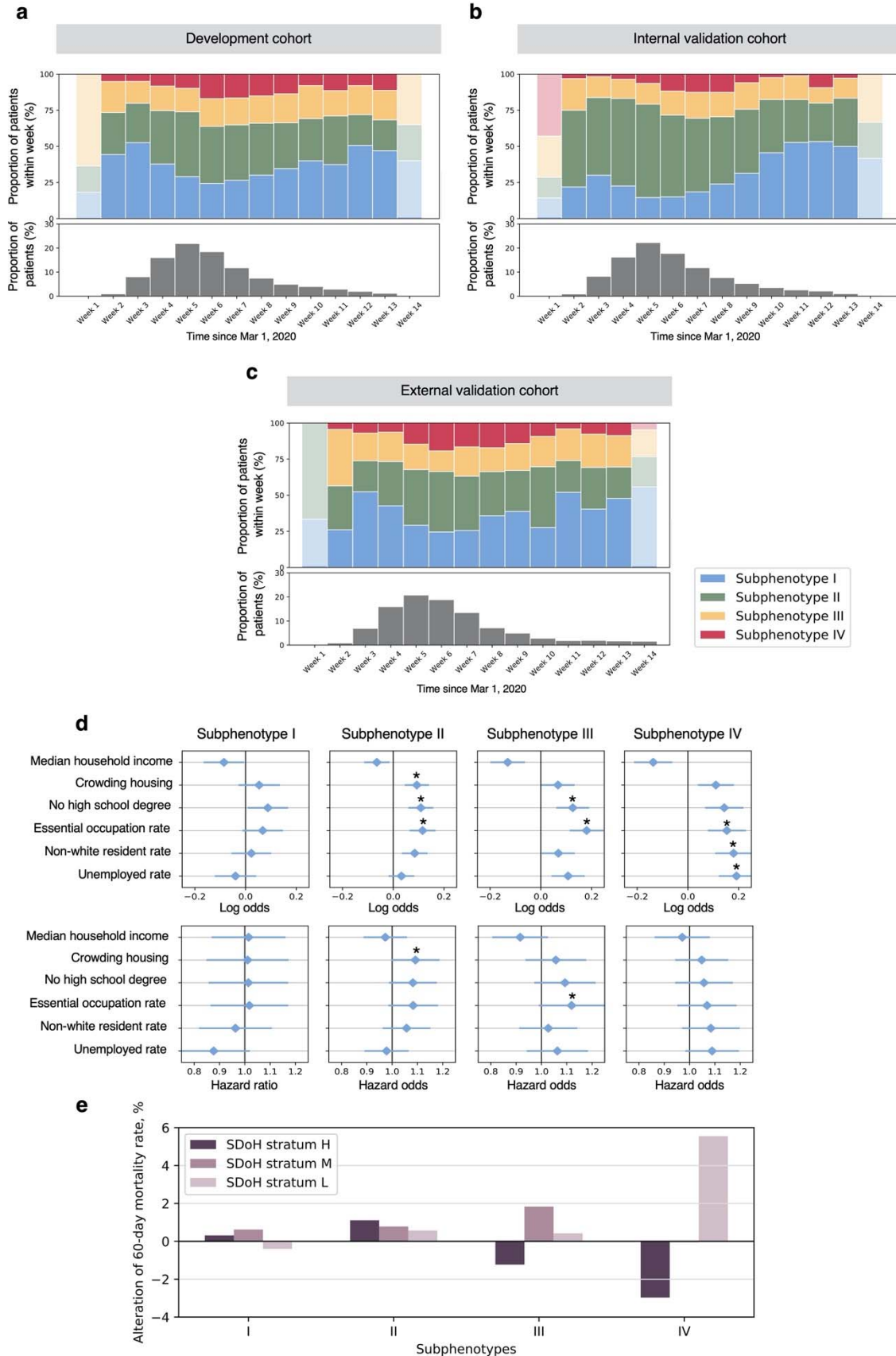


- a.** Abnormal biomarkers vs. all subphenotypes. **b.** Abnormal biomarkers vs. each subphenotype.  
**c.** Comorbidities vs. all subphenotypes. **d.** Comorbidities vs. each subphenotype

Abbreviations: ATA = asthma; CAD = coronary artery disease; COPD = chronic obstructive pulmonary disease; HF = heart failure; HLD = hyperlipidemia; HTN = hypertension.



**Figure 3. Kaplan-Meier (KM) plots for 60-day mortality by subphenotypes.** The survival probabilities were shown with 95% confidence interval. X-axis denotes time (days) after COVID-19 confirmation and Y-axis denotes the survival probability. **a-c.** KM plots by subphenotypes in the development, internal validation, and external validation cohorts, respectively.



**Figure 4. Plots showing temporal patterns and SDoH implications of subphenotypes.**

**a-c.** Proportions of subphenotype memberships of patients confirmed per week, since March 1, 2020. **d.** Log odds and Hazard ratio (mean values and standard deviation [error bar]) showing associations between individual SDoH characteristics and 60-day mortality risk, using logistic regression analysis and Cox regression analysis, adjusting for age and sex, respectively. **e.** Plot showing alteration of 60-day mortality rate (Y-axis) of each SDoH stratum to that of subphenotype level.

\* P-value < 0.05