

1 **Regional performance variation in external validation of four prediction models**
2 **for severity of COVID-19 at hospital admission: An observational multi-centre**
3 **cohort study**

4
5
6 **Authors**

7 Kristin E. Wickstrøm^{1,2}, Valeria Vitelli³, Ewan Carr⁴, Aleksander R. Holten^{2,5}, Rebecca Bendayan^{4,6}
8 Andrew H. Reiner⁷, Daniel Bean^{4,8}, Tom Searle^{4,6}, Anthony Shek⁹, Zeljko Kraljevic⁴, James Teo^{9,10},
9 Richard Dobson^{4,6,8, 11,12}, Kristian Tonby^{2,13}, Alvaro Köhn- Luque³, Erik K. Amundsen^{1,14}.

10

11 **Affiliations**

12 1. Department of Medical Biochemistry, Blood Cell Research Group, Oslo University Hospital,
13 Oslo, Norway.

14 2. Institute of Clinical Medicine, University of Oslo, Oslo, Norway.

15 3. Oslo Centre for Biostatistics and Epidemiology, Faculty of Medicine, University of Oslo, Oslo,
16 Norway.

17 4. Department of Biostatistics and Health Informatics, Institute of Psychiatry, Psychology and
18 Neuroscience, King's College London, London, U.K.

19 5. Department of Acute Medicine, Oslo University Hospital and Institute of Clinical Medicine,
20 University of Oslo, Oslo, Norway.

21 6. NIHR Biomedical Research Centre at South London and Maudsley NHS Foundation Trust and
22 King's College London, London, U.K.

23 7. Oslo Centre for Biostatistics and Epidemiology, Oslo University Hospital, Oslo, Norway.

24 8. Health Data Research UK London, University College London, London, U.K.

25 9. Department of Clinical Neuroscience, Institute of Psychiatry, Psychology and Neuroscience,
26 King's College London, London, U.K.

27 10. King's College Hospital NHS Foundation Trust, London, U.K.

28 11. Institute of Health Informatics, University College London, London, U.K.

NOTE: This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.

29 12. NIHR Biomedical Research Centre at University College London Hospitals NHS Foundation
30 Trust, London, U.K.

31 13. Dept. of Infectious diseases, Oslo University Hospital, Oslo, Norway.

32 14. Oslo Metropolitan University, Department of Life Sciences and Health

33

34

35 **Correspondence**

36 Erik K. Amundsen, Department of Medical Biochemistry, Oslo University Hospital, Oslo, Norway.

37 Mail: uxamue@ous-hf.no

38

39

40

41

42

43

44

45

46

47

48

49

50

51

52

53

54 **Abstract**

55 **Background:**

56 Several prediction models for coronavirus disease-19 (COVID-19) have been published. Prediction
57 models should be externally validated to assess their performance before implementation. This
58 observational cohort study aimed to validate published models of severity for hospitalized patients
59 with COVID-19 using clinical and laboratory predictors.

60

61 **Methods:**

62 Prediction models fitting relevant inclusion criteria were chosen for validation. The outcome was
63 either mortality or a composite outcome of mortality and ICU admission (severe disease). 1295
64 patients admitted with symptoms of COVID-19 at Kings Cross Hospital (KCH) in London, United
65 Kingdom, and 307 patients at Oslo University Hospital (OUH) in Oslo, Norway were included. The
66 performance of the models was assessed in terms of discrimination and calibration.

67 **Results:**

68 We identified two models for prediction of mortality (referred to as Xie and Zhang¹) and two models
69 for prediction of severe disease (Allenbach and Zhang²).

70 The performance of the models was variable. For prediction of mortality Xie had good
71 discrimination at OUH with an area under the receiver-operating characteristic (AUROC) 0.87 [95 %
72 confidence interval (CI) 0.79-0.95] and acceptable discrimination at KCH, AUROC 0.79 [0.76-0.82].

73 In prediction of severe disease, Allenbach had acceptable discrimination (OUH AUROC 0.81 [0.74-
74 0.88] and KCH AUROC 0.72 [0.68-0.75]). The Zhang models had moderate to poor discrimination.

75 Initial calibration was poor for all models but improved with recalibration.

76 **Conclusions:** The performance of the four prediction models was variable. The Xie model had the
77 best discrimination for mortality, while the Allenbach model had acceptable results for prediction of
78 severe disease.

79 **Introduction**

80 Severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) was discovered in Wuhan, China in
81 December 2019. The virus was shown to cause viral pneumonia, later designated as coronavirus
82 disease 2019 (COVID-19) [1]. The disease has evolved as a pandemic with an extensive amount of
83 severe cases with high mortality [2]. Several biomarkers, clinical and epidemiological parameters
84 have been associated with disease severity [3, 4]. Practical tools for prediction of prognosis in
85 COVID-19 patients are still lacking in clinical practice [5, 6]. Prediction models can be crucial to
86 prioritize patients needing hospitalization, intensive care treatment, or future individualized therapy.

87 Since the onset of the pandemic, the number of prediction models for COVID-19 patients has
88 been continuously growing [7]. Prediction models should be validated in different populations with a
89 sufficient number of patients reaching the outcome before implementation [8-10]. A validation study
90 of 22 prediction models at one site was recently published [6]. Interestingly, this study found that
91 none of the models performed better than oxygen saturation alone, even though the performance at
92 the original study sites in most cases was much better.

93 This study aimed to validate published prediction models of severity and mortality for hospitalized
94 patients based on laboratory and clinical values in COVID-19 cohorts from London (United
95 Kingdom) and Oslo (Norway).

96 The study is reported according to the guidelines in “Transparent reporting of a multivariable
97 prediction model for individual prognosis or diagnosis” (TRIPOD) [11] and has also followed
98 recommendations from “Prediction Model Risk of Bias Assessment Tool” (PROBAST) [12].

99

100

101 **Methods**

102

103 Selection of prediction models

104 A literature search was performed to select prediction models for validation. Published articles or
105 preprint manuscripts were included until 29.05.2020. A structured search was performed in PubMed
106 with the words “COVID-19” and “prediction model” or “machine learning” or “prognosis model”.
107 Prediction models included in the review by Wynants et. al. [7] in May 2020 were also investigated,
108 as well as search for articles/preprints citing Wynants et. al. using Google Scholar 18.05.2020.
109 The inclusion criteria for selection of multivariable prediction models were: (1) Symptomatic
110 hospitalized patients over 18 years with PCR confirmed COVID-19; (2) outcomes including
111 respiratory failure or intensive care unit (ICU) admission or death or composite outcomes of these.
112 (3) The predictive models had to include at least one laboratory test as we wanted to explore models
113 that combined clinical and laboratory variables (4). All variables had to be available in the OUH
114 dataset and the model had to be described in adequate detail.

115

116 Study design and participants.

117

118 The study was performed as a retrospective validation study with adult patients hospitalized with
119 COVID-19. Two cohorts were included: (1) Oslo University Hospital (OUH) in Norway, (2) Kings
120 Cross Hospital (KCH) in London, United Kingdom. The patients included were all adult inpatients
121 testing positive for SARS-CoV-2 by real-time polymerase chain reaction (RT-PCR) with symptoms
122 consistent with COVID-19 at admission. SARS-CoV-2 -positive patients admitted for conditions not
123 related to COVID-19 were excluded, e.g. pregnancy-related conditions or trauma. Patients referred
124 from other hospitals were also excluded, as we did not have access to measurements from the first
125 hospital admission.

126

127 OUH cohort

128 OUH is a large urban university hospital. Patients admitted between 6th March and 31th December
129 2020 were included. The OUH project protocol was approved by the Regional Ethical Committee of
130 South East Norway (Reference 137045). All patients with confirmed COVID-19 were included in
131 the quality registry “COVID19 OUS”, approved by the data protection officer (Reference 20/08822).
132 Informed consent was waived because of the strictly observational nature of the project.
133 Demographics, clinical variables and hospital stay information were manually recorded in the
134 registry and merged with laboratory results exported from the laboratory information system in
135 Microsoft Excel.

136

137 KCH cohort

138 In the KCH cohort patients were admitted between 23rd February to 1st May 2020 at two hospitals
139 (King’s College Hospital and Princess Royal University Hospital) in South East London (UK) of
140 Kings College Hospital NHS Foundation Trust.

141 Data (demographics, emergency department letters, discharge summaries, lab results, vital signs)
142 were retrieved from components of the electronic health record (EHR) using a variety of natural
143 language processing (NLP) informatics tools belonging to the CogStack ecosystem [13]. The project
144 operated under London South East Research Ethics Committee (reference 18/LO/2048) approval
145 granted to the King’s Electronic Records Research Interface (KERRI); specific work on COVID-19
146 research was reviewed with expert patient input on a virtual committee with Caldicott Guardian
147 oversight. Data from this cohort has been published in prior studies [14, 15].

148

149 Missing values

150 Predictive variables were collected from the admission to the emergency department (ED). If not
151 available in the ED, the first available values within 24 hours from hospital admission were used.
152 Missing values (i.e. no recorded values within 24 hours) were generally imputed using k-nearest

153 neighbors (KNN) although we tested more advanced techniques based on Python's scikit-learn
154 IterativeImputer, including random forest-based imputation, and multiple imputation using Bayesian
155 ridge and Gaussian process methods [16, 17].

156

157 Statistical analyses and performance measurements for the prediction models

158 Univariate comparisons between patients with 'mild' versus 'severe' disease were carried out for
159 continuous (Wilcoxon rank-sum test) and binary (X^2 test) measures. Severe disease was defined as
160 transfer to ICU or in-hospital mortality.

161 Validation of the selected prediction models was assessed with discrimination and calibration
162 as recommended in TRIPOD [11]. Discrimination is the ability of the model to differentiate between
163 those who do or do not experience the outcome. It is commonly estimated by concordance index (c-
164 index) which is identical to the area under the receiver-operating characteristic curve (AUROC) for
165 models with binary endpoints. The discrimination for the models at OUH and KCH was also
166 compared to the discrimination in the original development cohort and to the external validation by
167 Gupta et. al [6]. Calibration is the agreement between the observed outcomes and the outcome
168 predictions from the model. It is preferably reported by a calibration plot, intercept and slope.
169 Models were recalibrated by adjusting the intercept of the logistic regression models according to the
170 frequency of outcomes at each study site [18]. All statistical analyses were conducted in Python 3.7
171 and R 3.4 [19].

172

173 **Results**

174

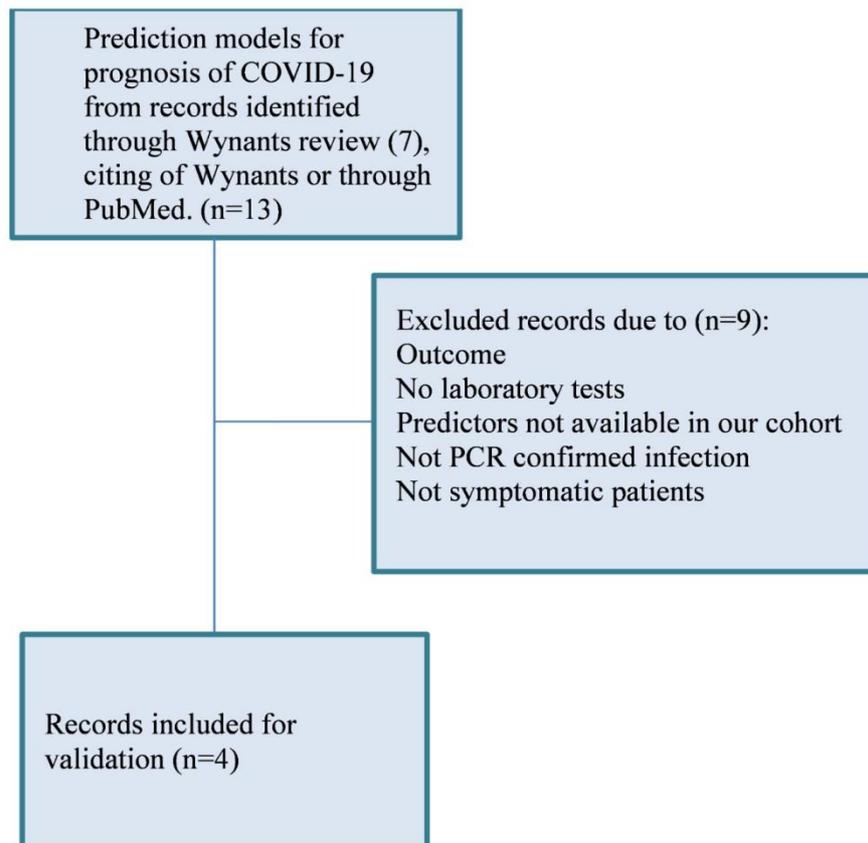
175 Selection of prediction models

176 Four publications comprising five prediction models fit our inclusion criteria [14, 20-22]. The
177 inclusion process is illustrated in Figure 1. However, since one of the models was developed at KCH

178 and validated at OUH in a previous publication [14], only four models are presented here. The four
179 models are referred to as ‘Xie’[20], ‘Zhang1’, ‘Zhang2’[21] and ‘Allenbach’[22].

180

181 **Figure 1: Selection of prediction models for validation.**



182

183

184 Information on the predictor variables and outcomes of the four models are summarized in
185 Table 1.

186 All predictors were measured at hospital admission. Missing values and imputation methods used in
187 the development cohorts were not well described. The Xie model had hospital mortality as the only
188 outcome. Zhang presented two models with different outcomes: (1) Mortality and (2) Composite
189 outcome of mortality or ‘poor outcome’. Poor outcome was defined as acute respiratory distress
190 syndrome (ARDS), intubation or extracorporeal membrane oxygenation (ECMO) treatment, ICU

191 admission or death. The Allenbach model used a composite outcome of transfer to ICU or mortality
 192 within 14 days of hospital admission. There were no details of the censoring date in the original
 193 studies. Mortality during the hospital stay was used for the OUH cohort and for the KCH cohort
 194 hospital mortality at data collection time.

195 All prediction models were based on multiple logistic regression and presented coefficients
 196 and intercepts for the different variables that enabled the calculation of risk prediction for our
 197 cohorts. Allenbach additionally provided an 8-point scoring system derived from the logistic
 198 regression model. However, we chose to use the regression model for calculation as this retains as
 199 much information as possible.

200 **Table 1: Predictors and outcomes in the three prediction models.**

	Zhang models	Xie model	Allenbach model
Country of development cohort	China	China	France
Predictors	Age, sex, neutrophil count, lymphocyte count, platelets count, CRP and creatinine at admission.	Age, LDH, SpO ₂ , lymphocyte count (log, due to extreme value).	CRP (per 100mg/L), age, lymphocyte count, WHO scale (22) by admission.
Outcome	1. Mortality 2. Poor outcome, defined as developing ARDS, receiving intubation or ECMO treatment, ICU admission and death.	1. Hospital mortality	1. ICU transfer or death by 14 days after admission.

201 CRP; C-reactive protein, LDH; lactate dehydrogenase, SpO₂; Peripheral oxygen saturation, WHO;
 202 World Health Organization, ICU; Intensive care unit, NEWS2; National Early Warning score 2,
 203 eGFR; estimated glomerular filtration rate, ECMO; extracorporeal membrane oxygenation, ARDS;
 204 acute respiratory distress syndrome.

205

206 Description of the cohorts

207 Patient characteristics for the three development cohorts and the KCH and OUH cohorts are shown
208 in Table S1 (supplementary material).

209 Since the three models use different outcomes and timeframes, the number of patients included in
210 each validation is not the same. An overview of missing values is presented in Table 2. Missing
211 values were imputed via simple imputation and multiple imputations [17]. Preliminary analyses
212 showed no differences between AUROCs calculated with different imputation methods (see Table
213 S2). Thus, the simple imputation method k-nearest neighbor was used for the rest of this paper. At
214 KCH the number of missing values was very high for LDH (87.8 %) and relatively high for SpO₂
215 (33.3 %) and WHO scale (33.8 %).

216 The OUH cohort consisted of 307 patients while the KCH cohort consisted of 1295 patients
217 (Figure S1). For the OUH cohort median age was 60 years with 57 % males, while in the KCH
218 cohort the median age was 69 with 59 % males. In the OUH cohort, 32 patients died in the hospital
219 (10.4 %), while 333 (26.8 %) had died at the hospital by data collection time in the KCH cohort. For
220 the composite outcome death or ICU transfer, the number of patients with the outcome was 66 (21.5
221 %) at OUH and 419 (33.7 %) at KCH.

222 The percentage of patients with hypertension and diabetes was higher in the KCH cohort (54
223 % and 35 %, respectively) than in the OUH cohort (34 % and 21 %, respectively). The patients at
224 KCH also had higher levels of CRP, creatinine, LDH, and possibly a lower number of lymphocytes
225 than the OUH patients; all of which are known predictors for severe COVID-19.

226 In Table 3, univariate associations are presented for mild/moderate and severe groups for the
227 KCH and OUH cohorts. In general, the same variables were predictive for severe disease at KCH
228 and OUH; except for ischemic heart disease, temperature and platelets which were associated with
229 severe disease at OUH, but not KCH.

230

231 **Table 2: Missing values, and results for discrimination and calibration.**

	Zhang1			Zhang2			Xie			Allenbach		
	Validation		Dev.	Validation		Dev.	Validation		Dev.	Validation		Dev.
	Oslo	London	Wuhan	Oslo	London	Wuhan	Oslo	London	Wuhan	Oslo	London	Paris
Participants, n	307	1244	775	307	1244	775	307	1286	299	307	1248	152
Outcome, n (%)	32 (10)	333 (27)	33 (4.3)	66 (22)	419 (34)	75 (9.7)	32 (10)	333 (26)	155 (52)	62 (20)	389 (31)	47 (32)
Missing values Predictors (%)												
-ALC	3.9	4.6	*	3.9	4.6	*	3.9	7.7	*	3.9	4.9	*
-ANC	3.9	4.7	NA	3.9	4.7	NA	NA	NA	NA	NA	NA	NA
-Platelets	NA	4.5	NA	NA	4.5	NA	NA	NA	NA	NA	NA	NA
-CRP	NA	3.3	NA	NA	3.3	NA	NA	NA	NA	NA	3.6	NA
-LDH	NA	NA	NA	NA	NA	NA	12.4	87.8	NA	NA	NA	NA
-Crea.	NA	3.5	NA	NA	3.5	NA	NA	NA	NA	NA	NA	NA
-SaO2	NA	NA	NA	NA	NA	NA	NA	33.3	NA	NA	NA	NA
-WHO	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	33.8	NA
Outcome (%)	None	0.07	NA	None	4.1	0.065	None	0.07	NA	None	3.8	0.03
AUROC (C-index)	0.72 [0.62-0.82]	0.64 [0.60-0.68]	0.91	0.77 [0.70-0.84]	0.67 [0.64-0.70]	0.88	0.87 [0.79-0.95]	0.79 [0.76-0.82]	0.89 [0.86-0.93]	0.81 [0.74-0.88]	0.72 [0.68-0.75]	0.79
Calibration Slope	0.57 [0.27-0.87]	0.37 [0.18-0.56]	0.98 [0.24-1.72]	1.18 [0.59-1.77]	0.75 [0.58-0.92]	1.04 [0.79-1.29]	0.86 [0.71-1.00]	1.03 [0.89-1.17]	1.00 [0.77-1.26]	1.03 [0.79-1.28]	0.76 [0.62-0.89]	0.89
Calibration intercept	0.04 [-0.01-0.09]	0.17 [0.10-0.23]	0.19 [-0.10-0.48]	-0.02 [-0.17-0.12]	0.10 [0.04-0.16]	0.01 [-0.13-0.15]	-0.02 [-0.05-0.01]	-0.04 [-0.09-0.01]	0.00 [-0.33-0.33]	0.00 [-0.06-0.07]	0.09 [0.04-0.14]	-0.06
Calibration before recal.; slope	0.47 [0.24-0.71]	0.38 [0.18-0.58]	-	1.56 [0.68-2.43]	0.92 [0.66-1.17]	-	0.53 [0.29-0.77]	0.87 [0.73-1.00]	-	1.19 [0.94-1.45]	0.87 [0.70-1.00]	-
Calibration before recal.; intercept	0.03 [-0.02-0.08]	0.17 [0.11-0.24]	-	0.01 [-0.13-0.15]	0.16 [0.10-0.22]	-	-0.06 [-0.16-0.03]	-0.12 [-0.19-0.06]	-	0.02 [-0.04-0.07]	0.12 [0.07-0.17]	-

232 * Information missing. Dev.=Development, ALC= Absolute lymphocyte count, ANC=Absolute

233 neutrophil count, Crea.=Creatinine, recal.=recalibration

234

235

236

237

238 **Table 3: Univariate analysis of predictors at KCH and OUH**

	OUH cohort				KCH cohort			
	N	Mild/ Moderate disease	Severe disease	P-value	N	Mild/ moderate disease	Severe disease	P- value
Age	307	55 [46-70]	68 [58-78]	<0.01	1295	67 [53-82]	75 [62- 86]	<0.01
Male sex (%)	307	129 (54)	46 (70)	0.02	1295	463 (56)	271(65)	0.01
Hypertension (%)	307	75 (31)	29 (44)	0.05	1295	428 (52)	244 (58)	0.04
Diabetes (%)	307	46 (19)	18 (27)	0.15	1295	282 (34)	154 (37)	0.40
Ischemic heart disease (%)	307	20 (8)	13 (20)	0.01	1295	105 (13)	66 (16)	0.17
Chronic lung disease (%)	307	61 (25)	22 (33)	0.19	1295	82 (10)	52 (12)	0.22
Days at hospital	307	5.0 [2.0-9.0]	16.5 [8.0-24.0]	<0.01	854	7.0 [3.0- 12.0]	16.0 [10.5- 31.1]	<0.01
Temperature (celcius)	306	37.1 [36.5-37.8]	37.8 [36.8-38.8]	<0.01	864	36.9 [36.6- 37.4]	37.0 [36.6- 37.5]	0.22
Resp/min (highest)	303	22 [18-28]	28 [22-32]	<0.01	860	19 [18- 20]	20 [19- 24]	<0.01
NEWS2 score	299	4 [2-6]	7 [5-10]	<0.01	815	2 [1- 4]	4 [2- 6]	<0.01
SpO2 ¹	307	96.0 [93.0-98.0]	92.0 [87.3-95.0]	<0.01	858	96.0 [95.0, 98.0]	96.0 [94.0, 97.0]	<0.01
CRP (mg/L)	307	34 [10-74]	93 [45-154]	<0.01	1203	73 [33- 128]	118 [59- 196]	<0.01
Creatinine (µmol/L)	307	77 [64-94]	97 [71-128]	<0.01	1200	87 [69- 118]	108 [83- 166]	<0.01
LDH (U/L)	269	237 [188-305]	329 [242-499]	<0.01	157	349 [277- 431]	532 [393- 706]	<0.01
Neutrophils (10 ⁹ /L)	295	4.0 [2.0-5.9]	5.7 [2.8-7.9]	<0.01	1186	5.1 [3.6- 7.2]	6.4 [4.5- 8.8]	<0.01
Lymphocytes (10 ⁹ /L)	295	1.1 [0.8-1.6]	0.8 [0.6-1.1]	<0.01	1187	1.0 [0.7-1.3]	0.9 [0.6-1.3]	<0.01
Platelets (10 ⁹ /L)	307	215 [173-279]	179 [135-244]	<0.01	1188	214 [165- 269]	205 [153- 272]	0.15

239 ICU; intensive care unit, ACE; angiotensin converting enzyme, NEWS; National early warning
240 score, CRP; C-reactive protein. LDH; Lactate dehydrogenase, eGFR; estimated glomerular filtration
241 rate. Continuous variables in median [IQR] and categorical variables in number (percent). P-values
242 are calculated with the Pearson X^2 test for categorical variables, and with the Wilcoxon rank-sum test
243 for continuous variables. SaO₂ value under oxygen treatment was registered if oxygen was applied,
244 there are also values for patients without oxygen in this registration.

245

246 Performance of the prediction models

247 The validation of the four prediction models with both the OUH and KCH cohorts is presented in
248 terms of discrimination (AUROC) and calibration (slope and intercept) in Table 2 and Figures 2 and
249 3, respectively. For the models predicting mortality, the Xie model had the highest AUROC both in
250 the KCH cohort (0.79; 95 % CI 0.76-0.82) and the OUH cohort (0.87; 95 % CI 0.79-0.95). The
251 Zhang1 model had a lower AUROC at both KHC (0.64; 95 % CI 0.60-0.68) and OUH (0.72; 95 %
252 CI 0.62-0.82).

253 For ‘severe disease’, discrimination was highest in the Allenbach model with AUROCs 0.72 (95 %
254 CI 0.68-0.75) for KCH and 0.81 (95 % CI 0.74-0.88) for OUH. For the Zhang2 model, the AUROC
255 was 0.67 (95 % CI 0.64-0.70) for KCH and 0.77 (95 % CI 0.70-0.84) for OUH. For the Xie and
256 Allenbach models, discrimination at OUH was similar to the development cohorts (Figure 2). And,
257 although the difference was not statistically significant at the 0.05 confidence level, we found better
258 discrimination for both of these models at OUH compared to KCH.

259 The calibration plots are shown in Figure 3 (after recalibration). Figure S3 in supplementary
260 shows the calibration results before and after recalibration for the Xie and Allenbach models.
261 Recalibration will not render models with poor discrimination more useful. Thus, we focused on the
262 recalibration of the Xie and Allenbach models as these had the best discrimination. Recalibration

263 improved the predictions for both the Xie and Allenbach models at OUH and the Xie model at KCH,
264 and the slope and intercept were acceptable for both models at both hospitals after recalibration.

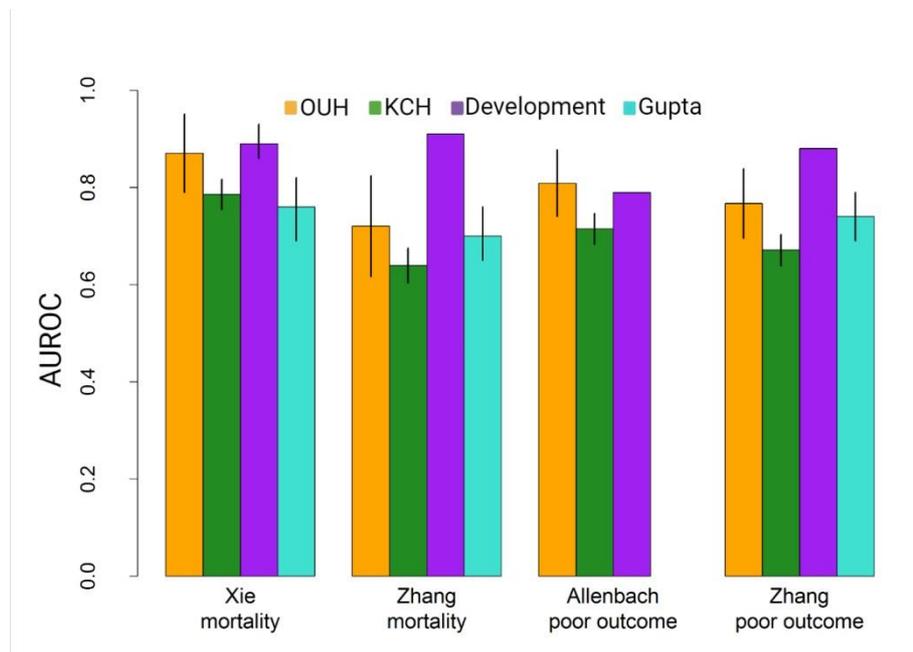
265

266 **Figure 2: AUROCs from validation of the four models at the KCH and OUH cohorts, and the**

267 **original AUROC from development cohorts [20-22]. Also shown are the results from the**

268 **external validation of the Xie and Zhang models by Gupta et al [6]. Lines represent the 95%**

269 **CI of the AUROCs. For the development cohorts only Xie reported confidence intervals.**



270

271

272

273

274

275

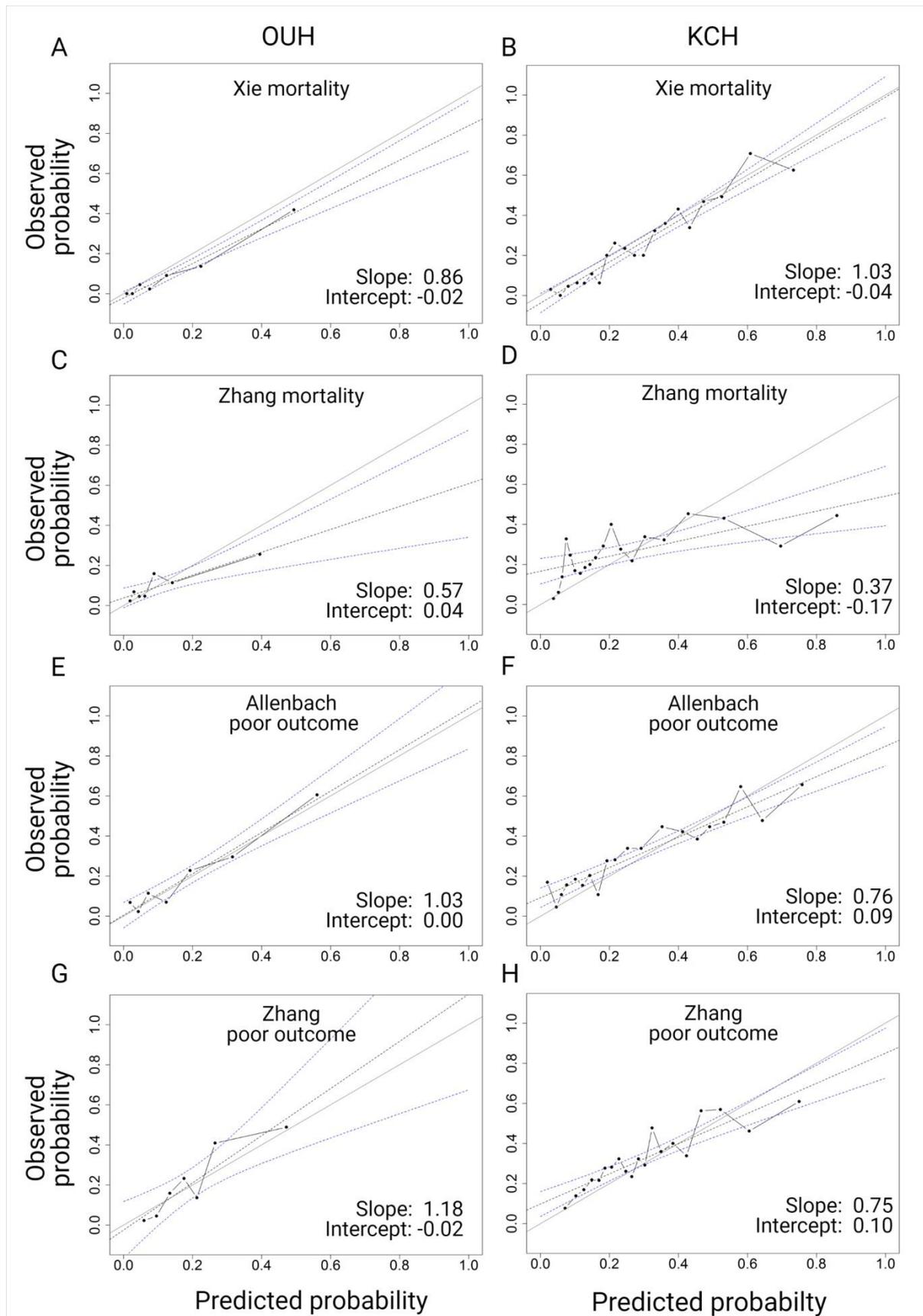
276

277

278

279

280 **Figure 3: Calibration plots for OUH and KCH after recalibration**



281

282 **Discussion**

283 In this study, we validated four prediction models for prognosis in hospitalized COVID-19 patients
284 from London, UK and Oslo, Norway. We found varying performance of the models in the two
285 cohorts. The models performed better in the OUH cohort with similar discrimination to the original
286 studies. The Xie and Allenbach models had the best performance for prediction of death and severe
287 disease, respectively.

288 Initial calibration was poor for all models, but improved after recalibration of the intercept
289 according to the frequency of the outcome in our cohorts. This improves the accuracy of the
290 prediction for each patient without affecting the discrimination and is recommended in several
291 publications [5, 11, 18]. Local or possibly regional/national recalibration is likely to be important for
292 COVID-19 prediction models since there is a large variation in the frequency of severe disease and
293 death in different studies.

294 In some cases, we found poorer discrimination in the validation cohorts compared to the
295 development cohorts. This is consistent with past evidence showing discrimination in development
296 cohorts to be better than at external validation due to overfitting and differences in characteristics of
297 the cohorts [23, 24]. The cohorts in the original studies and at KCH and OUH had many differences
298 such as mortality, age and frequencies of severe disease and comorbidities. UK and Norway differ in
299 the structures of their healthcare systems, and the incidence of COVID-19 has been far higher in the
300 UK. These factors may have affected the selection of patients for hospital and ICU admission, which
301 might have resulted in a more homogenous patient population in regards to severity at KCH. It is to
302 be expected that discrimination will be less good when the population is more homogenous.

303 The findings underline the importance of validation at several external sites. This is
304 particularly true for a new disease like COVID-19, with rapidly developing treatment guidelines, and
305 with an overwhelming effect on healthcare resources in some locations, but not at others.

306 The Xie model had the best results compared to the other models. The differences in the
307 performance of the prediction models might have several reasons. Firstly, the predictors used in one
308 model might have better predictive value than predictors used in others. SaO₂, which is included in
309 the Xie model, is a strong clinical indicator of the severity of disease, and often indicates a need for
310 ICU transfer. Secondly, there might be weaknesses in the models, as bias is common in prediction
311 models [12]. To date, only the Allenbach study is published in a peer-reviewed journal, while Xie
312 and Zhang are preprints. Thirdly, criteria for ICU admittance might vary across sites. The fact that
313 we and other studies generally find better discrimination for mortality than for severe disease (often
314 defined by ICU admittance) supports this hypothesis. For instance, patients with short life
315 expectancy will often not be admitted to the ICU, but given oxygen therapy in a hospital ward and
316 transferred to nursing homes for palliative care. These patients, not fulfilling the criteria for severe
317 disease, often have predictors that indicate severe disease at admission.

318 Many prediction models have been published, but few have been systematically validated
319 [24]. To our knowledge, only one study to date has validated COVID-19 prediction models; Gupta
320 et. al recently validated 22 prognostic models [6], including the Xie and Xhang models. For the OUH
321 cohort, we found substantially better discrimination for the Xie and Allenbach models for the
322 prediction of mortality and severe disease, respectively. The performance of the models at KCH was
323 more similar to the results in the Gupta study, also performed at a London hospital. The rate of
324 severe disease, mortality and the characteristics of the London cohorts are quite similar which might
325 explain the similar performance at these two sites.

326 Several other prediction models have been recently published, such as models based on
327 NEWS2 or the ISARIC model [14, 25]. The AUROCs of the models are in the range of 0.75 to 0.80,
328 which is not a substantial improvement over single univariate predictors of severity. Thus, the
329 finding that the Xie and Allenbach models perform well at both the original study site and at our

330 validation cohort at OUH might indicate that it is possible to achieve higher AUROCs with relatively
331 simple prediction models.

332 Our study has several strengths. Validation was performed at two sites in different countries
333 with consistent inclusion and exclusion criteria. We included all eligible patients admitted to the
334 hospital during the study period therefore the cohorts should be representative of the study sites.
335 Moreover, the study was conducted and reported according to the TRIPOD guidelines. However,
336 there are also some weaknesses. Firstly, the OUH cohort is not very large with relatively few patients
337 meeting the outcomes. Some publications recommend including at least 100 patients with the
338 relevant outcome [10]. However, studies with lower numbers are frequently published, and may still
339 contain useful information. Furthermore, the KCH cohort is probably one of the largest cohorts
340 analyzed in prediction models for severe COVID-19. Secondly, Gupta et. al. included 22 models in
341 their validation study, while we ended our inclusion of models in May, and included only four
342 models in this study. Whereas it could be interesting to include more models we think that the results
343 for the Xie and Allenbach models at OUH indicate that further studies of these models could be
344 interesting. Thirdly, there was a relatively high number of missing values for LDH and SpO₂ at
345 KCH. It is uncertain how much this affected the results. Both are included in the Xie model and
346 SpO₂ is a strong predictor for mortality, while LDH is probably a weaker predictor (6). The number
347 of missing values at OUH was low and probably did not affect the validation.

348 In conclusion, following the TRIPOD guidelines, our study validated developed models for
349 prediction of prognosis in COVID-19, and showed that these models have a variable performance in
350 different cohorts. The Xie model and Allenbach model clearly had the best performance, and we
351 suggest that these models should be included in future studies of COVID-19 prediction models.
352 However, the performance of these models at our two validation sites was not similar, which

353 underlines the importance of external validation of prediction models at several study sites before
354 their implementation in the clinical practice.

355 **Supporting information**

356 **S1 File. Supplementary tables**

357 **S2 File. Supplementary figures**

358 **S2 File TRIPOD checklist**

359 **Acknowledgments**

360 We would like to thank Prof. Anne Ma Dyrhol Riise and Dr. Ane M. Andersson at the Department of
361 Infectious Diseases, OUH, for their support with the quality registry “COVID19 OUS”.

362

363 **Author contributions**

364 **Conceptualization:** Kristin Wickstrøm, Erik K. Amundsen, Kristian Tonby, Aleksander R. Holten,
365 Valeria Vitelli, Alvaro Köhn-Luque

366 **Data Curation:** Kristin Wickstrøm, Erik K. Amundsen, Kristian Tonby, Valeria Vitelli, Alvaro
367 Köhn-Luque, Andrew H. Reiner, Ewan Carr, Rebecca Bendayan, Daniel Bean, Anthony Shek,
368 Zeljko Kraljevic

369 **Formal Analysis:** Valeria Vitelli, Alvaro Köhn-Luque, Andrew H. Reiner, Ewan Carr

370 **Methodology:** Kristin Wickstrøm, Erik K. Amundsen, Kristian Tonby, Valeria Vitelli, Alvaro
371 Köhn-Luque,

372 **Project Administration:** Erik K. Amundsen, Kristian Tonby

373 **Resources:** Erik K. Amundsen, James Teo, Richard Dobson

374 **Software:** Valeria Vitelli, Alvaro Köhn-Luque, Andrew H. Reiner, Ewan Carr, Daniel Bean,
375 Anthony Shek, Zeljko Kraljevic

376 **Supervision:** Erik K. Amundsen, James Teo, Richard Dobson

377 **Validation:** Kristin Wickstrøm, Valeria Vitelli, Alvaro Köhn-Luque, Andrew H. Reiner, Ewan Carr

378 **Visualization:** Kristin Wickstrøm, Valeria Vitelli, Alvaro Köhn-Luque, Andrew H. Reiner, Erik K.

379 Amundsen

380 **Writing (original draft):** Kristin Wickstrøm, Erik K. Amundsen

381 **Writing (Review):** Kristian Tonby, Aleksander R. Holten, Valeria Vitelli, Alvaro Köhn-Luque,

382 Ewan Carr, Rebecca Bendayan, Daniel Bean, Anthony Shek, Zeljko Kraljevic, Tom Searle, James

383 Teo, Richard Dobson

384

385 **References**

386 1. Zhou F, Yu T, Du R, Fan G, Liu Y, Liu Z, et al. Clinical course and risk factors for mortality of adult
387 inpatients with COVID-19 in Wuhan, China: a retrospective cohort study. *Lancet*.

388 2020;395(10229):1054-62. doi: 10.1016/S0140-6736(20)30566-3. PubMed PMID: 32171076.

389 2. Weiss P, Murdoch DR. Clinical course and mortality risk of severe COVID-19. *Lancet*.

390 2020;395(10229):1014-5. doi: 10.1016/S0140-6736(20)30633-4. PubMed PMID: 32197108;

391 PubMed Central PMCID: PMC7138151.

392 3. Henry BM, de Oliveira MHS, Benoit S, Plebani M, Lippi G. Hematologic, biochemical and immune
393 biomarker abnormalities associated with severe illness and mortality in coronavirus disease 2019
394 (COVID-19): a meta-analysis. *Clinical chemistry and laboratory medicine*. 2020. Epub 2020/04/15.

395 doi: 10.1515/cclm-2020-0369. PubMed PMID: 32286245.

396 4. Zeng F, Li L, Zeng J, Deng Y, Huang H, Chen B, et al. Can we predict the severity of COVID-19 with
397 a routine blood test? *Polish archives of internal medicine*. 2020. Epub 2020/05/02. doi:

398 10.20452/pamw.15331. PubMed PMID: 32356642.

399 5. Martin GP, Sperrin M, Sotgiu G. Performance of prediction models for COVID-19: the Caudine
400 Forks of the external validation. *Eur Respir J*. 2020;56(6). doi: 10.1183/13993003.03728-2020.

401 PubMed PMID: 33060155.

402 6. Gupta RK, Marks M, Samuels THA, Luintel A, Rampling T, Chowdhury H, et al. Systematic
403 evaluation and external validation of 22 prognostic models among hospitalised adults with COVID-
404 19: An observational cohort study. *Eur Respir J*. 2020. doi: 10.1183/13993003.03498-2020. PubMed
405 PMID: 32978307.

406 7. Wynants L, Van Calster B, Bonten MMJ, Collins GS, Debray TPA, De Vos M, et al. Prediction
407 models for diagnosis and prognosis of covid-19 infection: systematic review and critical appraisal.

408 2020;369:m1328. doi: 10.1136/bmj.m1328 %J BMJ.

409 8. Steyerberg EW, Bleeker SE, Moll HA, Grobbee DE, Moons KG. Internal and external validation of
410 predictive models: a simulation study of bias and precision in small samples. *Journal of clinical
411 epidemiology*. 2003;56(5):441-7. doi: 10.1016/s0895-4356(03)00047-7. PubMed PMID: 12812818.

412 9. Altman DG, Vergouwe Y, Royston P, Moons KG. Prognosis and prognostic research: validating a
413 prognostic model. *BMJ (Clinical research ed)*. 2009;338:b605. Epub 2009/05/30. doi:

414 10.1136/bmj.b605. PubMed PMID: 19477892.

- 415 10. Collins GS, Ogundimu EO, Altman DG. Sample size considerations for the external validation of a
416 multivariable prognostic model: a resampling study. *Statistics in medicine*. 2016;35(2):214-26. Epub
417 2015/11/11. doi: 10.1002/sim.6787. PubMed PMID: 26553135; PubMed Central PMCID:
418 PMCPMC4738418.
- 419 11. Moons KGM, Altman DG, Reitsma JB, Ioannidis JPA, Macaskill P, Steyerberg EW, et al.
420 Transparent Reporting of a multivariable prediction model for Individual Prognosis Or Diagnosis
421 (TRIPOD): Explanation and Elaboration. *Annals of Internal Medicine*. 2015;162(1):W1-W73. doi:
422 10.7326/M14-0698 %J *Annals of Internal Medicine*.
- 423 12. Wolff RF, Moons KGM, Riley RD, Whiting PF, Westwood M, Collins GS, et al. PROBAST: A Tool to
424 Assess the Risk of Bias and Applicability of Prediction Model Studies. *Annals of Internal Medicine*.
425 2019;170(1):51-8. doi: 10.7326/M18-1376 %J *Annals of Internal Medicine*.
- 426 13. Jackson R, Kartoglu I, Stringer C, Gorrell G, Roberts A, Song X, et al. CogStack - Experiences of
427 deploying integrated information retrieval and extraction services in a large National Health Service
428 Foundation Trust hospital. *BMC Medical Informatics and Decision Making*. 2018;18. doi:
429 10.1186/s12911-018-0623-9.
- 430 14. Carr E, Bendayan R, Bean D, Stammers M, Wang W, Zhang H, et al. Evaluation and improvement
431 of the National Early Warning Score (NEWS2) for COVID-19: a multi-hospital study. *BMC Medicine*.
432 2021;19(1):23. doi: 10.1186/s12916-020-01893-3.
- 433 15. Zakeri R, Pickles A, Carr E, Bean DM, O'Gallagher K, Kraljewic Z, et al. Biological responses to
434 COVID-19: Insights from physiological and blood biomarker profiles. *Curr Res Transl Med*.
435 2021;69(2):103276. doi: 10.1016/j.retram.2021.103276. PubMed PMID: 33588321; PubMed
436 Central PMCID: PMCPMC7857048.
- 437 16. Donders AR, van der Heijden GJ, Stijnen T, Moons KG. Review: a gentle introduction to
438 imputation of missing values. *Journal of clinical epidemiology*. 2006;59(10):1087-91. Epub
439 2006/09/19. doi: 10.1016/j.jclinepi.2006.01.014. PubMed PMID: 16980149.
- 440 17. White IR, Royston P, Wood AM. Multiple imputation using chained equations: Issues and
441 guidance for practice. 2011;30(4):377-99. doi: 10.1002/sim.4067.
- 442 18. Janssen KJ, Moons KG, Kalkman CJ, Grobbee DE, Vergouwe Y. Updating methods improved the
443 performance of a clinical prediction model in new patients. *Journal of clinical epidemiology*.
444 2008;61(1):76-86. doi: 10.1016/j.jclinepi.2007.04.018. PubMed PMID: 18083464.
- 445 19. team Rc. A language and environment for statistical computing. R Foundation for Statistical
446 Computing, Vienna, Austria. 2020. Available from : <https://www.R-project.org/>.
- 447 20. Xie J, Hungerford D, Chen H, Abrams ST, Li S, Wang G, et al. Development and external
448 validation of a prognostic multivariable model on admission for hospitalized patients with COVID-
449 19. 2020:2020.03.28.20045997. doi: 10.1101/2020.03.28.20045997 %J *medRxiv*.
- 450 21. Zhang H, Shi T, Wu X, Zhang X, Wang K, Bean D, et al. Risk prediction for poor outcome and
451 death in hospital in-patients with COVID-19: derivation in Wuhan, China and external validation in
452 London, UK. 2020:2020.04.28.20082222. doi: 10.1101/2020.04.28.20082222 %J *medRxiv*.
- 453 22. Allenbach Y, Saadoun D, Maalouf G, Vieira M, Hellio A, Boddaert J, et al. Multivariable
454 prediction model of intensive care unit transfer and death: a French prospective cohort study of
455 COVID-19 patients. 2020:2020.05.04.20090118. doi: 10.1101/2020.05.04.20090118 %J *medRxiv*.
- 456 23. Siontis GC, Tzoulaki I, Castaldi PJ, Ioannidis JP. External validation of new risk prediction models
457 is infrequent and reveals worse prognostic discrimination. *Journal of clinical epidemiology*.
458 2015;68(1):25-34. doi: 10.1016/j.jclinepi.2014.09.007. PubMed PMID: 25441703.
- 459 24. Ramspek CL, Jager KJ, Dekker FW, Zoccali C, van Diepen M. External validation of prognostic
460 models: what, why, how, when and where? *Clin Kidney J*. 2021;14(1):49-58. doi:
461 10.1093/ckj/sfaa188. PubMed PMID: 33564405; PubMed Central PMCID: PMCPMC7857818.

462 25. Knight SR, Ho A, Pius R, Buchan I, Carson G, Drake TM, et al. Risk stratification of patients
463 admitted to hospital with covid-19 using the ISARIC WHO Clinical Characterisation Protocol:
464 development and validation of the 4C Mortality Score. *BMJ (Clinical research ed)*. 2020;370:m3339.
465 doi: 10.1136/bmj.m3339. PubMed PMID: 32907855; PubMed Central PMCID: PMC7116472.
466