



## Supplementary Materials for

### **Symptom clusters in Covid19: A potential clinical prediction tool from the COVID Symptom study app**

**Authors:** Carole H. Sudre<sup>1\*†</sup> & Karla A. Lee<sup>2†</sup> & Mary Ni Lochlainn<sup>2†</sup>, Thomas Varsavsky<sup>1</sup>, Benjamin Murray<sup>1</sup>, Mark S. Graham<sup>1</sup>, Cristina Menni<sup>2</sup>, Marc Modat<sup>1</sup>, Ruth C E Bowyer<sup>2</sup>, Long Nguyen<sup>3</sup>, David Drew<sup>3</sup>, Amit D. Joshi<sup>3</sup>, Wenjie Ma<sup>3</sup>, Chuan-Guo Guo<sup>3</sup>, Chun-Han Lo<sup>3</sup>, Sajaysurya Ganesh<sup>4</sup>, Abubakar Buwe<sup>4</sup>, Joan Capdevila Pujol<sup>4</sup>, Julien Lavigne du Cadet<sup>4</sup>, Alessia Visconti<sup>2</sup>, Maxim B Freydin<sup>2</sup>, Julia S. El-Sayed Moustafa<sup>2</sup>, Mario Falchi<sup>2</sup>, Richard Davies<sup>4</sup>, Maria F. Gomez<sup>5</sup>, Tove Fall<sup>5</sup>, M.Jorge Cardoso<sup>1</sup>, Jonathan Wolf<sup>4</sup>, Paul W. Franks<sup>2,5</sup>, Andrew T. Chan<sup>3</sup>, Tim D. Spector<sup>2</sup>, Claire J Steves<sup>2†</sup>, Sébastien Ourselin<sup>1\*†</sup>

Correspondence to: [carole.sudre@kcl.ac.uk](mailto:carole.sudre@kcl.ac.uk) / [sebastien.ourselin@kcl.ac.uk](mailto:sebastien.ourselin@kcl.ac.uk)

#### **This PDF file includes:**

Materials and Methods  
Figs. S1 to S4

#### **Material and Methods**

##### **Assessment of exposure, Ascertainment of outcomes, Ascertainment of covariates**

Exposure, outcome and covariates were all ascertained via the app as previously described (3,4). A subset of individuals reported being tested for COVID-19. BMI was calculated as kg/m<sup>2</sup>. Visit to hospital was recorded if the location was ever recorded as hospital or “back from hospital”.

The analysis of the disease course was separated into the unsupervised clustering and a predictive analysis for need of respiratory support using the projected clusters.

### **Subject Selection**

All participants included in the analysis were required to:

- 1) Report a hospital visit or show sign of recovery (recovery was defined as a significant drop (at least 2) in number of reported symptoms at the day of last report compared to the day where sum of symptom was maximal)
- 2) Record their symptoms on the app at least three times over four days or more between the time their symptoms start and either a hospital visit or the start of symptom decline.
- 3) For the recovery group: be tested positive
- 4) For the group visiting hospital: either a) be tested positive, b) be imputed positive from the day where sum of symptoms was maximal following the imputation method described in Menni et al (4) c) Reported classic symptoms of COVID19 (2 days or more of fever and cough)

All participants fulfilling the criteria of selection were included with no selection done on the country (UK, US or Sweden) in which they reported.

### **Statistical analysis**

Data from the app were downloaded into a server and only records where the self-reported characteristics fell within the following ranges were utilized for further analyses: age between 16 (18 in the US) and 100 years; BMI between 16 and 55 kg/m<sup>2</sup>.

14 symptoms were recorded: abdominal pain, chest pain, sore throat, shortness of breath, fatigue, hoarse voice, headache, loss of smell or taste, confusion, diarrhea, fever, persistent cough, unusual muscle pains or skipped meals.

For the assessment of risk factors, reported respiratory support was used as dependent variable in a logistic regression, using age, BMI, as continuous variables and sex, frailty score, diabetes, lung, heart or kidney disease as binary variables. The frailty score was binarized at a threshold of 3 points.

### Clustering analysis

Unsupervised time series clustering was performed using Mc2PCA (5) with 6 dimensions of projection from the 14 recorded binarized symptom course. This method allows for the clustering of time series with non-equal duration, using the covariance matrix of the time series. Optimization of the clustering is performed using a K-means iterative process as follows: for each cluster a singular value decomposition is performed over the average of covariance matrices and the first  $n$  (here 6) dimensions are used to calculate the projection. Attribution to each cluster is then chosen to minimize the residual error after projection. The process was iterated until the change in error ratio was below  $10^{-4}$  (convergence criteria)

Linear interpolation between time points was used in the case of missing data and a limit of 5 days of interrupted record was filled in. Where more than one record was present for a single day, the latest record was considered.

To determine the optimal number of clusters to consider over the disease continuum, for each number of clusters, the K-Means clustering Mc2PCA was run with 20 random initialization and the attempt with final minimal average distortion was selected. The Bayesian information criteria (BIC) was applied to balance model fit and model complexity leading to a choice of 6 clusters. Supplementary figure 2 presents the output of the BIC selection:

After separation of the different clusters, an average symptom course was calculated as the interpolated frequency of reported symptom for each given day over the mean duration for the given cluster. Ability to predict final cluster classification with a reduced number of days of reporting was assessed via weighted precision and recall score in the 6-class problem and reported for the training and test set.

### Predictive analysis

To demonstrate the relevance of the clustering to predict need for respiratory support, the performance of two random forest models on the testing set were compared. The first one used only the demographic characteristics while the second used the projected cluster at 5 days and the associated aggregated sum of symptoms over the five days. Youden Index obtained from the training set (0.066 and 0.059) for the model including symptomatology or only personal characteristics was applied on the test set for binary classification.

The random forests were trained to optimize the receiver operating curve area in a five-fold cross-validation setting with randomized grid search over hyperparameters.

A bootstrap analysis with 1000 samples was used to provide confidence intervals in the reported precision and false positive rates.

**Figure S1.** Flowchart showing entry of participants into analysis

**Abbreviations:** T0 no test; T+/-: Tested positive/negative; H+/H-: Attended hospital or did not attend hospital; I+: imputed COVID-19; SR+: self-reported classical symptoms of COVID-19 (2 days or more of fever and cough)

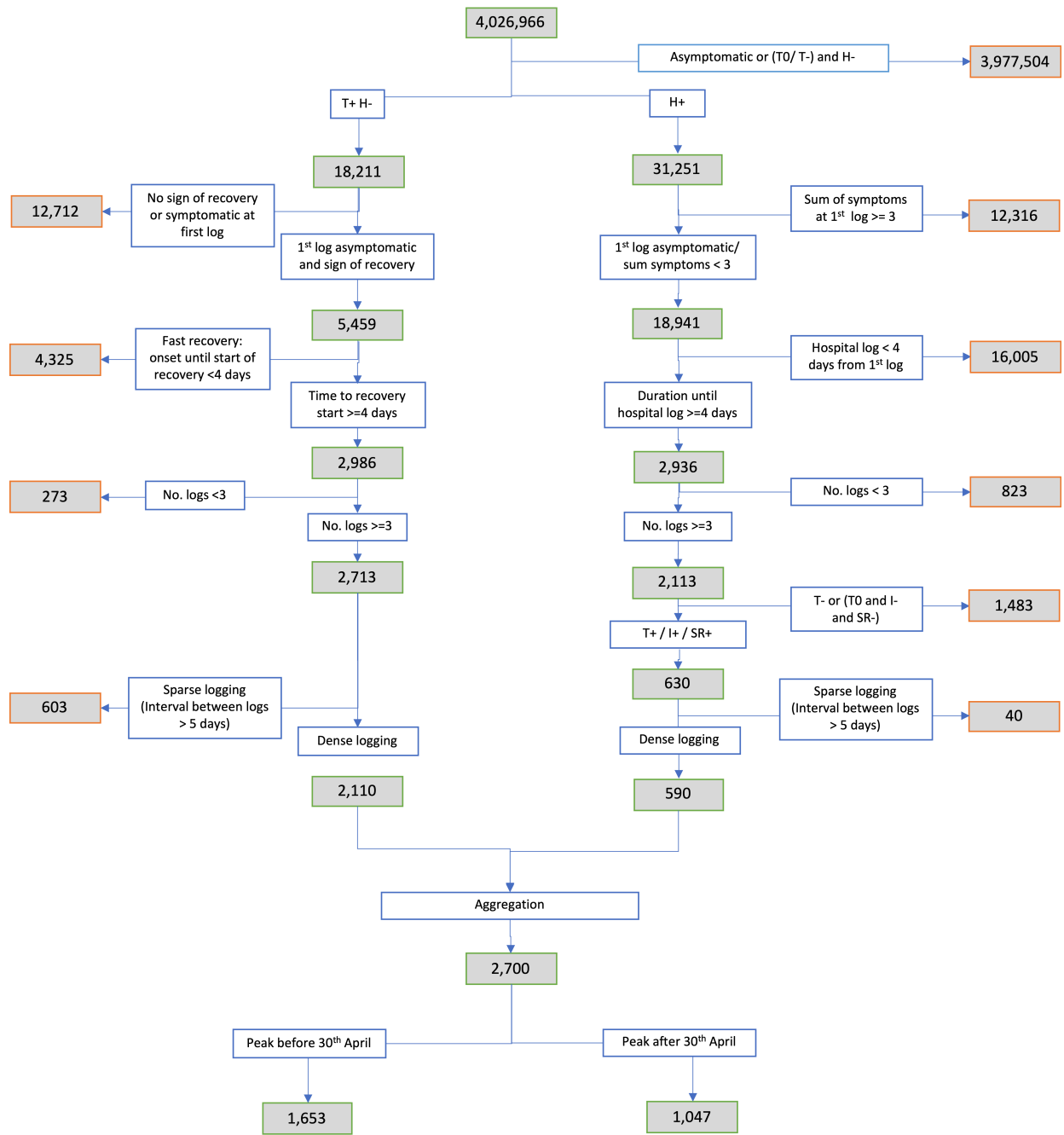
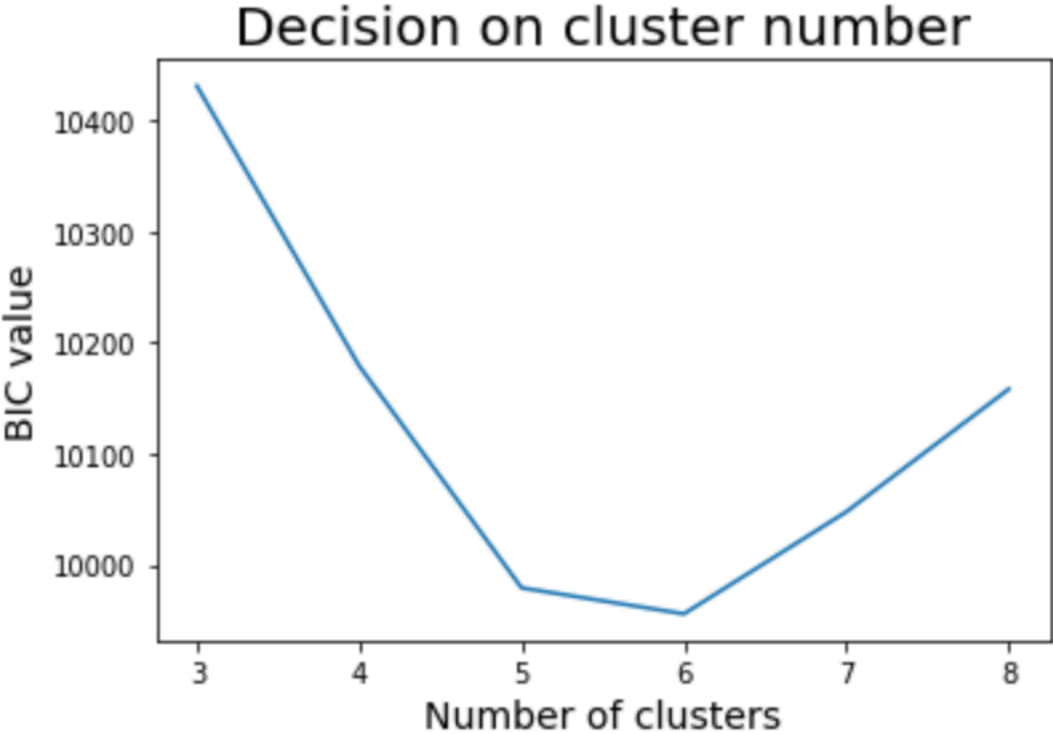


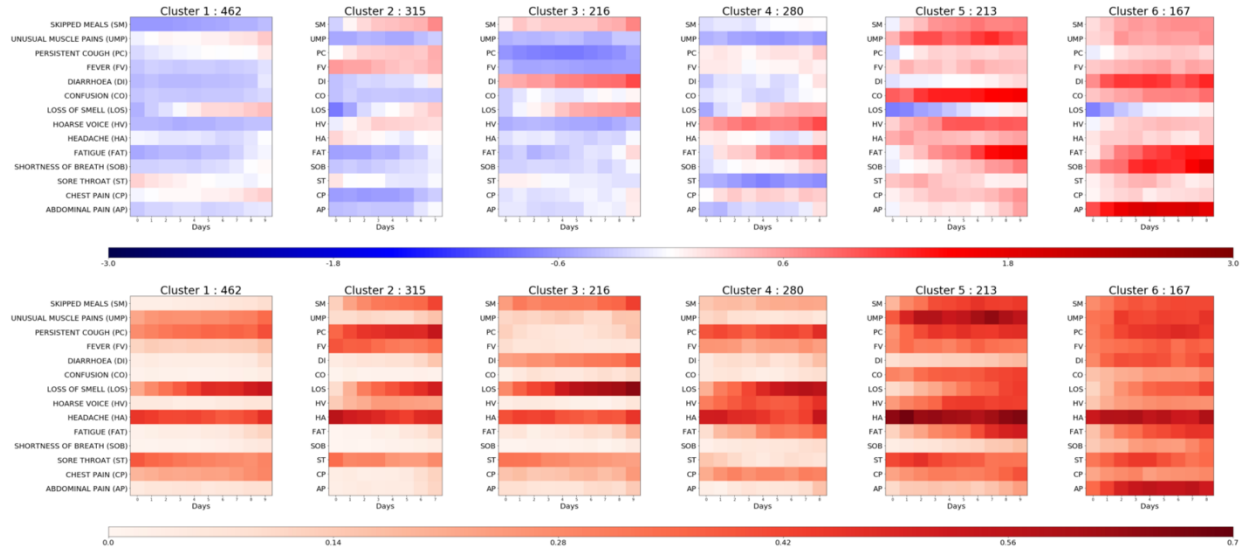
Figure S2 Bayesian information criterion for the choice of number of clusters



**Figure S3**

**Top: associated Z-Score of presentation of symptoms over overall symptom distribution (red = reported more than average; blue = reported less than average)**

**Bottom: Frequency of positive answers per symptoms across days for each cluster (darker = reported more frequently)**





**Fig S4 Confusion matrix for cluster assignment on the independent replication set using the projections learnt with limited number of days from the training set**

