



OPEN

Clustering of clinical and echocardiographic phenotypes of covid-19 patients

Eran Shpigelman^{1,5}, Aviram Hochstadt^{2,3,4,5}, Dan Coster¹, Ilan Merdler^{2,4}, Eihab Ghantous^{2,4}, Yishay Szekely^{2,4}, Yael Lichter^{2,4}, Philippe Taieb^{2,4}, Ariel Banai^{2,4}, Orly Sapir^{2,4}, Yoav Granot^{2,4}, Lior Lupu^{2,4}, Ariel Borohovitz^{2,4}, Sapir Sadon^{2,4}, Shmuel Banai^{2,4}, Ronen Rubinshtein^{3,4}, Yan Topilsky^{2,4} & Ron Shamir¹✉

We sought to divide COVID-19 patients into distinct phenotypical subgroups using echocardiography and clinical markers to elucidate the pathogenesis of the disease and its heterogeneous cardiac involvement. A total of 506 consecutive patients hospitalized with COVID-19 infection underwent complete evaluation, including echocardiography, at admission. A k-prototypes algorithm applied to patients' clinical and imaging data at admission partitioned the patients into four phenotypical clusters: Clusters 0 and 1 were younger and healthier, 2 and 3 were older with worse cardiac indexes, and clusters 1 and 3 had a stronger inflammatory response. The clusters manifested very distinct survival patterns (C-index for the Cox proportional hazard model 0.77), with survival best for cluster 0, intermediate for 1–2 and worst for 3. Interestingly, cluster 1 showed a harsher disease course than cluster 2 but with similar survival. Clusters obtained with echocardiography were more predictive of mortality than clusters obtained without echocardiography. Additionally, several echocardiography variables (E' lat, E' sept, E/e average) showed high discriminative power among the clusters. The results suggested that older infected males have a higher chance to deteriorate than older infected females. In conclusion, COVID-19 manifests differently for distinctive clusters of patients. These clusters reflect different disease manifestations and prognoses. Although including echocardiography improved the predictive power, its marginal contribution over clustering using clinical parameters only does not justify the burden of echocardiography data collection.

COVID-19 infection disease severity ranges widely, from asymptomatic or mild, self-limiting illness to severe progressive pneumonia, multiorgan failure, and death¹.

In addition to respiratory manifestations, several somewhat specific complications have been shown to be associated with COVID-19 illness, including cardiac and cardiovascular complications², thromboembolic complications³, neurologic complications⁴, and inflammatory manifestations⁵. As the clinical picture is quite variable, a question emerges regarding the factors that direct the disease in a specific course.

The case of cardiac involvement is considerably diverse. Although cardiac complications are common and are associated with increased mortality⁶, cardiac involvement is heterogeneous, including right ventricular (RV) dysfunction or dilatation, left ventricular (LV) diastolic dysfunction and systolic dysfunction (10%)².

As patients' baseline characteristics and disease manifestations vary, we hypothesized that different epitomes of disease manifestations may be identified. To identify those, we sought to use a strategy of machine learning-based clustering. Unsupervised discovery of subtypes of a single disease has been widely used in cardiology⁷, infectious disease⁸, and critical care medicine⁹ to find better pathologic explanations and improve current treatments for these conditions.

Although several trials have endeavored unsupervised clustering of COVID-19 illness^{10,11}, these included no comprehensive data regarding cardiac performance. As the echocardiographic data of patients with COVID-19 illness are essential for elucidation of both pathogenesis and prognosis¹², we find this addition imperative to an enhanced clustering of COVID-19 illness.

¹The Blavatnik School of Computer Science, Tel Aviv University, P.O. Box 39040, 6997801 Tel Aviv, Israel. ²Department of Cardiology, Tel Aviv Sourasky Medical Center, Dafna St 5, Tel Aviv-Yafo, Israel. ³Heart Institute, Edith Wolfson Medical Center, Ha-Lokhamim St 62, 5822012 Holon, Israel. ⁴The Sackler School of Medicine, The Tel-Aviv University, Tel Aviv, Israel. ⁵These authors contributed equally: E. Shpigelman and A. Hochstadt. ✉email: rshamir@tau.ac.il

Methods

Details about data acquisition and gathering have been specified previously². In brief, we prospectively studied consecutive adult patients (aged ≥ 18 years) admitted between March 21, 2020, and September 16, 2020, to the Tel Aviv Medical Center due to COVID-19 infection. All patients had a diagnosis of COVID-19 infection confirmed by a positive reverse-transcriptase polymerase chain reaction assay. Demographic data, comorbid conditions, medications, physical examination, laboratory, and ECG findings were systematically recorded. All patients underwent comprehensive transthoracic echocardiography within 48 h of admission as part of a predefined step-by-step protocol. Clinical and imaging data were collected prospectively. Mortality analysis started at the time of baseline echocardiographic examination and included in-hospital mortality. Mortality was ascertained until the end of follow-up, beyond hospitalization and irrespective of discharge date, for all patients by telephone calls and was complete for all the patients.

Ethics approval. Since data were evaluated retrospectively, pseudonymously and were solely obtained for treatment purposes, the ethics committee of the Tel Aviv Medical Center approved the study (institutional review board number 0196-20-TLV) and voided the requirement of informed consent. The research was performed in accordance with the Declaration of Helsinki.

Echocardiography. Echocardiography was performed in a standard manner with the same equipment (CX 50; Philips Medical Systems, Bothell, WA) by cardiologists with expertise in echocardiographic recording and interpretation. In accordance with current guidelines¹³, the following measures were undertaken to minimize the risk of infection: (1) All echocardiographic studies were bedside studies performed at the designated COVID-19 intensive care or internal ward units. (2) All echocardiographic examinations were performed with small, dedicated scanners because of their easier disinfection. (3) Echocardiographic scanners were set aside in each COVID-19-designated ward to minimize the risk of infection spread. (4) Personal protection at the time of echocardiographic recordings included N-95 respirator masks, fluid-resistant gowns, gloves, head covers, and eye shields. (5) Electrocardiographic monitoring during imaging was omitted, and all measurements were performed offline to reduce exposure and contamination. LV diameters, ejection fraction, and mass were measured as recommended¹⁴. Measurements of mitral inflow included the peak early filling (E wave) and late diastolic filling (A wave) velocities, E/A ratio, and deceleration time of early filling velocity. Early diastolic mitral septal and lateral annular velocities (e') were measured in the apical 4-chamber view¹⁵. Left atrial volume was calculated with the biplane area-length method at end systole. Forward stroke volume was calculated from the LV outflow tract with subsequent calculation of cardiac output.

From 4-chamber views encompassing the entire RV, end-systolic and end-diastolic RV areas and tricuspid annulus were measured. RV function was evaluated by tricuspid annular plane systolic excursion (TAPSE), systolic tricuspid lateral annular velocity measured in the apical 4-chamber view, and fractional area change^{14,16}. Hemodynamic right-sided assessment included the measurement of the pulmonic flow acceleration time to assess pulmonary vascular resistance¹⁷.

Patients underwent another comprehensive echocardiographic test whenever there was any clinical deterioration, i.e., the need for mechanical ventilation and hemodynamic support, according to the treating physician's judgment. This test was performed in the same manner as the first test performed upon arrival.

The cohort. Patients who signed the DNR/DNI ($n=24$) were removed from the cohort, resulting in 506 patients. For each patient, we used as input for clustering the measurements that were taken at admission and the first echocardiography result. Variables describing outcomes and treatments ($n=23$) and medications ($n=43$) were not part of the input data for clustering and were later used for evaluation of the clinical significance of the clusters. Variables that were missing in more than 2/3 of the cohort were also excluded ($n=23$). We allowed a relatively high missing rate, as the data are sparse, to keep a high number of variables. The missing rate of each variable is found in Supplementary 3. This process left 85 continuous variables and 56 categorical variables. Thirty-one of the continuous variables and two of the categorical variables were from echocardiography.

Computational methods. *Imputation and normalization.* Missing values in continuous variables were imputed using the Iterative Imputer algorithm based on MICE¹⁸. The continuous variables were normalized using the Yeo-Johnson power transform for nonnegative variables¹⁹ (see Supplementary 1). For the categorical variables, missing values were imputed with the most frequent value (another method for imputing the categorical variables achieved similar results. For details, see Supplementary 1).

Clustering. We used the k-prototypes²⁰ algorithm, which is a distance-based algorithm that allows the use of mixed data, i.e., both categorical and continuous variables. We chose k-prototypes because it was reported as one of the best performers in a recent benchmark study on mixed-data clustering algorithms²¹. The algorithm receives as input the number of clusters k and uses the distance function between variable vectors x, y :

$$d(x, y) = \sum_{j=1}^p (x_j - y_j)^2 + \gamma \sum_{j=p+1}^m \delta(x_j, y_j) \quad (1)$$

where x_1, \dots, x_p are numerical variables, x_{p+1}, \dots, x_m are categorical variables and δ is the Hamming distance function. γ defines the relative weight assigned to the categorical variables. We used $\gamma = 3$ and $k = 4$ (see Supplementary 2, 3).

To obtain robust clustering, we applied consensus clustering²² to multiple clustering results obtained on subsampled data. In each repetition, we randomly chose a fraction r of the patients and clustered them using k -prototypes. The process was repeated n times. For patients i and j , define $M(i, j)$ as the fraction of times in which they were in the same cluster, and the distance between them as $D(i, j) = 1 - M(i, j)$. Applying (regular) k -means on the distance matrix D gives the final clusters. We used $r = 0.85$ and $n = 50$ (see Supplementary 2).

Evaluation of clusters. To evaluate the quality of the clusters, we looked for clinically significant characteristics of the patients composing the different clusters. For each variable, we tested its significance under the null assumption that it does not vary between clusters. For the continuous variables, we performed ANOVA, and for the categorical variables, we performed the Chi2 test, both implemented in SciPy²³. P-values were corrected for multiple testing with FDR.

To find variables that are most discriminative between clusters, we computed the absolute standardized mean differences (ASMD) score (see Supplementary 1).

To visualize the different characteristics of the clusters, we created a radar plot of selected variables. Each variable was normalized to $[0.2-1]$, where 0.2 is the lowest cluster average or percentage and 1 is the highest.

To analyze the survival trends in different clusters, we plotted the Kaplan–Meier survival curves²⁴ for each cluster and computed the conditional multivariate log-rank test to compare the survival plots across the clusters. To measure the predictability of the clusters for in-hospital mortality, we estimated survival times from the Cox proportional hazard model²⁵ with binary variables that represent the cluster's membership as the covariates of the model and calculated the c-index²⁶.

We also wished to evaluate whether the clusters showed distinct survival patterns among the patients who received respiratory support. For that, we built another Cox proportional hazard model for in-hospital mortality using only the patients who received respiratory support, with cluster labels as the covariates. We excluded Cluster 0, where only one patient received such support, and computed hazard ratios per cluster (relative to cluster 1). Implementations of Kaplan–Meier, Cox proportional hazard, C-index and log-rank were taken from lifelines²⁷.

Evaluating the contribution of echocardiography. To evaluate the contribution of the 33 echocardiography variables to creating meaningful clusters, we tested the change in c-index obtained after randomly permuting across patients the values for each of the 33 variables, independently for each variable. In this way, we kept the distribution of each variable and the number of input variables unchanged. Next, we performed a similar process where we randomly chose 31 continuous variables and 2 categorical variables (the same numbers as for the echocardiography variables) and permuted their values to compare the contribution of the echocardiography data to randomly chosen variables. We also tested the change in echocardiography measurements over time using the results of a second echocardiography that some of the patients underwent. Full details of this analysis are found in Supplementary 6.

Results

The final cohort studied included 506 hospitalized patients with PCR-positive COVID-19 of average age 62.31 and 36.96% females ($n = 187$). Further demographic and clinical characteristics are shown in Table 1.

Identifying distinct patient subgroups. We clustered the patients into four clusters labeled 0 to 3, with sizes of 128, 195, 112, and 71, respectively. We analyzed several possible values for the number k of clusters and selected $k = 4$ as the appropriate number (see full analysis in Supplementary 2). As we demonstrate below, the echocardiography parameters had a major contribution to the quality of the results.

Feature	Average \pm STD/percentage
Number of individuals	506
Age	62.31 \pm 17.30
Sex (female)	36.96%
Total days in hospital	9.28 \pm 12.02
MEWS score at admission	4.66 \pm 3.15
CRP at admission	85.27 \pm 78.01
History of hypertension	45.45%
Diabetes	31.42%
Obesity	26.09%
Ischemic heart disease (IHD) or congestive heart failure (CHF)	19.96%
Chronic renal failure (CRF)	9.68%
Dementia/cognitive decline	6.92%
Liver disease	3.56%

Table 1. Demographics and selected clinical features of the cohort. MEWS COVID-19 Modified Early Warning Score for clinical deterioration²⁸, CRP C-reactive protein, a marker of inflammation.

The main parameters that significantly distinguish between the clusters are presented in Table 2, and the full list is shown in Supplementary 3. Based on these parameters, the most prominent characteristics of the clusters are as follows:

- **Cluster 0**—young patients (average age: 43.63) with less medical history of significant diseases, with mild signs of inflammation at admission (low CRP) and no deaths.
- **Cluster 1**—relatively young patients (average age: 61.27) with a more severe inflammatory process. Of these, 27% received respiratory support, and 10% died while hospitalized.
- **Cluster 2**—older patients (average age: 76.67) with fewer inflammatory markers upon admission, with a richer medical history (e.g., 75% with hypertension) and more females than any other cluster (69%).
- **Cluster 3**—older patients (average age: 76.22) with marked inflammatory processes and a significant medical history (e.g., 83% with hypertension). These patients had the highest in-hospital mortality rate (54%).

To further assess the discriminative value of the different variables between the clusters, we calculated the absolute standardized mean differences (ASMD) scores for all variables. Figure 1 shows the variables with the highest scores. Among the variables with high ASMD are some natural COVID-19 parameters, such as CRP for inflammatory state, pulmonary ultrasound findings and findings in chest X-ray for lung function, and MEWS and SOFA, established warning scores for clinical deterioration. Age had the highest score, and several echocardiography variables, in particular left ventricle variables (E' lat, E' sept, E/e average), were scored high. We also evaluated the clinical merit of the clusters by calculating their discriminative power with respect to outcomes, which were not part of the input for clustering. Among them, in-hospital mortality and respiratory support-related variables scored high.

	Cluster 0	Cluster 1	Cluster 2	Cluster 3	p-value
Number of individuals	128	128	112	71	–
Age*	43.63 ± 12.82	43.63 ± 12.82	76.67 ± 10.88	76.22 ± 10.41	6 × 10 ⁻⁸²
CRP*	33.71 ± 37.73	33.71 ± 37.73	36.31 ± 35.26	133.42 ± 81.62	6 × 10 ⁻⁴⁸
MEWS score at Admission*	1.91 ± 2.07	1.91 ± 2.07	5.07 ± 2.46	7.68 ± 2.99	2 × 10 ⁻³²
E/e average*	6.98 ± 1.59	6.98 ± 1.59	12.33 ± 5.56	14.35 ± 5.79	2 × 10 ⁻³⁵
At*	108.43 ± 25.59	108.43 ± 25.59	81.26 ± 27.93	68.80 ± 19.75	3 × 10 ⁻²⁴
O2 Saturation*	96.75 ± 3.13	96.75 ± 3.13	95.10 ± 4.77	86.45 ± 12.29	1 × 10 ⁻¹⁷
BNP	19.00 ± 17.85	19.00 ± 17.85	162.74 ± 206.94	609.94 ± 917.85	2 × 10 ⁻¹⁶
Troponin	8.49 ± 18.03	8.49 ± 18.03	20.84 ± 38.65	479.01 ± 1837.05	2 × 10 ⁻³
History of hypertension	9%	9%	75%	83%	7 × 10 ⁻³²
Sex (Female)	38%	38%	69%	32%	3 × 10 ⁻¹⁵
Respiratory support**	1%	1%	13%	62%	1 × 10 ⁻²²
In hospital mortality**	0%	0%	6%	54%	3 × 10 ⁻²⁷

Table 2. Statistics of selected variables that were significantly different among the clusters. P-values were computed using ANOVA for continuous variables and Chi2 for categorical variables, and FDR corrected for multiple testing. *Continuous variables, mean ± SD of each cluster. **This outcome variables were not used by the clustering algorithm. CRP C-reactive protein a marker of inflammation, MEWS COVID-19 Modified Early Warning Score for clinical deterioration²⁹. E/e average a left ventricle echocardiography parameter, At a right ventricle echocardiography parameter.

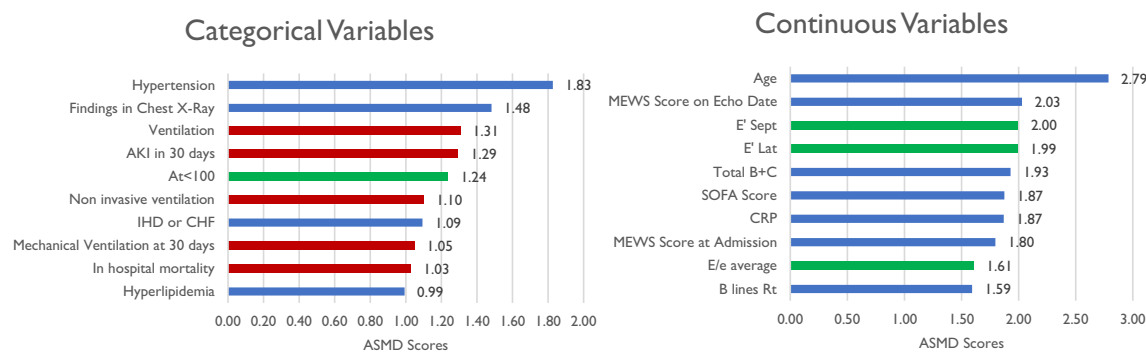


Figure 1. Variables with the highest ASMD scores. In red are outcomes that were not used for the clustering; in green are echocardiography variables. Total B + C a lung ultrasound parameter, combining consolidation and B line results.

Figure 2 shows a radar plot of significantly distinguishing variables (see Supplementary 4) that were selected to represent different aspects of the disease and the characteristics of the patients. Cluster 0 stands out as having the best results in most health-related parameters. In contrast, cluster 3 has the most severe results in most parameters. Cluster 2 patients have very similar age distributions to those in cluster 3, they have higher rates of dementia, and in both clusters we see high rates of hypertension and other comorbidities (likely due to high average age). However, remarkably, the COVID-19-related parameters of cluster 2 are much better: low CRP and high O₂ saturation (this variable was reversed, so a high value in the plot means low O₂ saturation). On the other hand, patients in cluster 1 were younger with fewer background diseases but worse COVID-19-related variables, such as high CRP, chest X-ray findings and O₂ saturation.

Interestingly, sex significantly differed among the clusters. Cluster 2 is significantly enriched with female patients and the only cluster with a majority of females, while clusters 1 and 3 have a percentage of females below the cohort average (37%). A comparison of outcomes between males and females aged 80 and above in the full cohort showed that males are significantly more likely to receive respiratory support (p-value = 0.02), which is a sign for deterioration. All other outcomes were worse in males but not significantly so, perhaps due to the small sample size. For full details see Supplementary 7.

We also wished to evaluate to what extent the clusters differ in terms of in-hospital mortality. Figure 3 shows the Kaplan-Meier survival curves of the four subgroups. We can see three distinct survival patterns, with no deaths in cluster 0, an intermediate survival pattern for clusters 1 and 2, and the worst survival of cluster 3. The log-rank test to assess differences in survival functions across clusters produced a highly significant p-value of 2.73×10^{-12} . The c-index for the Cox proportional hazard model (see “Methods”) was 0.77, which is relatively high.

Note that clusters 1 and 2 are very different in terms of age but have similar survival plots. For a detailed comparison of clusters 1 and 2, see Supplementary 5. While many parameters seem to be age related, Fig. 3 suggests that the course of the disease is not entirely dominated by age.

Figure 4 shows the fraction of patients in each cluster who received respiratory and hemodynamic support. Patients in Clusters 1 and 3 suffered from severe disease, and therefore, a higher percentage of them were treated with respiratory or hemodynamic support (drug or mechanical).

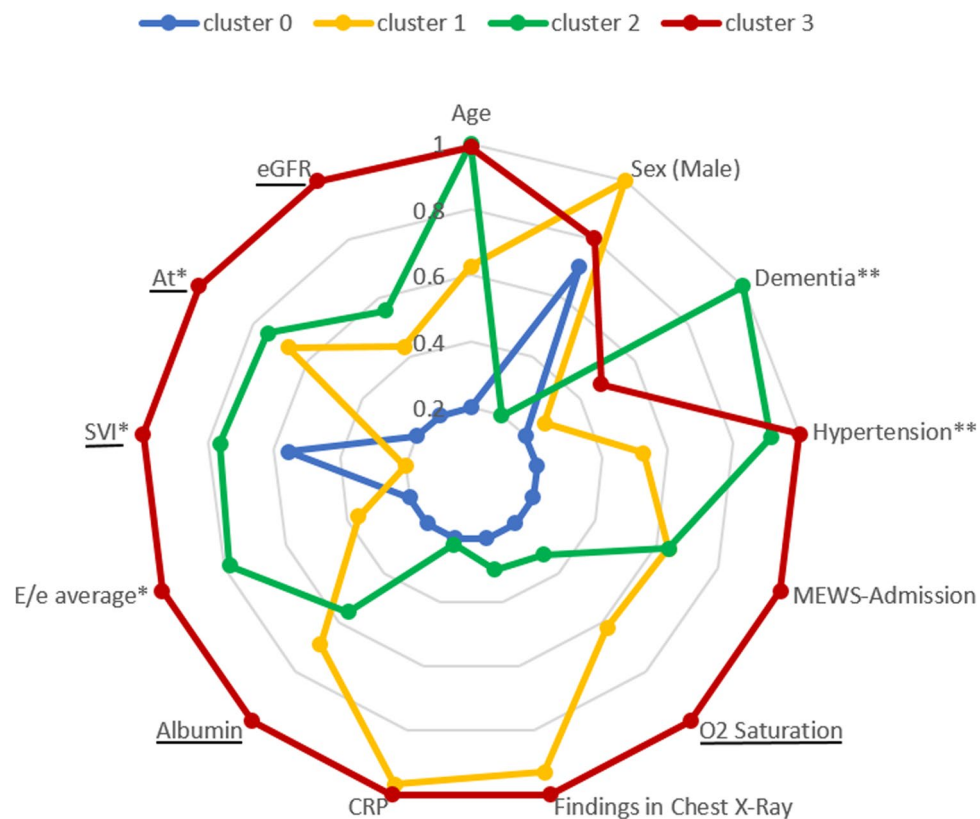


Figure 2. Radar plot for selected variables in clusters. Each variable is normalized to [0.2, 1], where 0.2 is the lowest cluster average/percentage and 1 is the highest. For ease of interpretability, the scale of some variables was reversed so that an increase from 0.2 to 1 always accounts for worse conditions. Underlined variables were reversed. *Echocardiography features: *E/e average* left ventricle feature, *At* right ventricle feature, *SVI* Stroke volume index. ***Hypertension* and *Dementia* are shown as examples for past diseases. Other past diseases showed similar trends (Supplementary 3). *CRP* C-reactive protein a marker of inflammation; *MEWS* COVID-19 Modified Early Warning Score for clinical deterioration²⁸, *eGFR* Estimated Glomerular Filtration Rate.

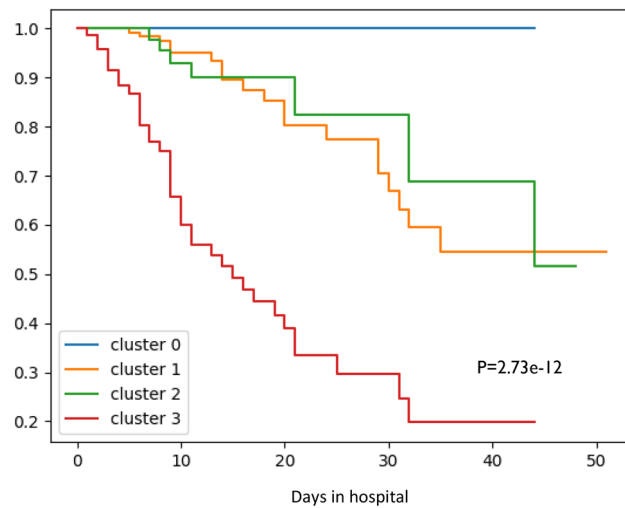


Figure 3. Kaplan–Meier survival curves for each cluster for the event “In-hospital mortality”. P is the log-rank p-value.

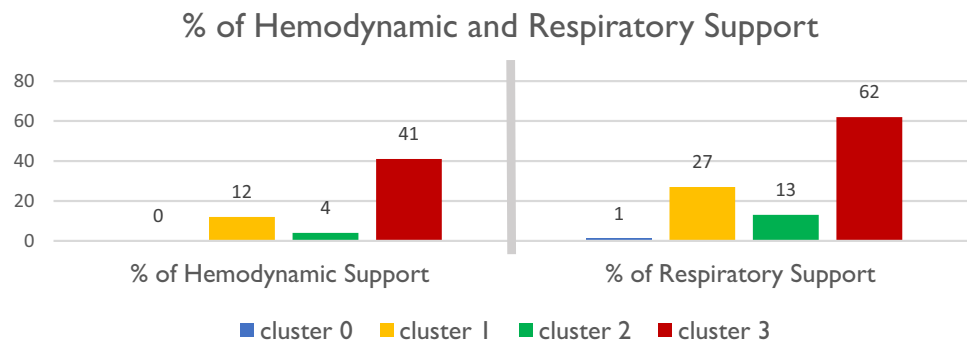


Figure 4. Percentage of patients who received respiratory support and hemodynamic support in each cluster.

Clusters 1–3 included a substantial number of patients who needed respiratory support. To test if their risk differed across clusters, we built a Cox proportional hazard ratio model for in-hospital mortality based only on those patients, using cluster labels as the covariates. Cluster 0 was excluded because only one patient in that group was ventilated. Cluster 1 was used as a reference. The hazard ratios were 1.25 (0.50, 3.14) for cluster 2 and 4.27 (2.42, 7.52) for cluster 3 (95% confidence interval in parentheses). Hence, ventilated patients in cluster 3 were at higher risk than those in cluster 1, although their inflammatory condition was similar at admission (CRP 129 vs. 133, see Table 1). Ventilated patients in cluster 2 were at a similar risk to those in cluster 1 despite the age difference between the groups. Although patients in cluster 2 were relatively older, they arrived in better inflammatory condition. The Cox model for the ventilated patients shows similar trends in survival to the full cohort, as was observed in the Kaplan–Meier curves in Fig. 3 (results not shown).

There were no significant differences between clusters in terms of treatment. For more details see Supplementary 8.

Echocardiography contribution. Full statistics of the echocardiography variables are presented in Table 3. To further assess the contribution of the echocardiography data to the formation of meaningful clusters, we shuffled the echocardiography values, reclustered the resulting data 50 times (see “Methods”) and recomputed the c-index and log-rank p-values in each case. The average log-rank p-value was $1.97e-06 \pm 6.54e-06$ with a median of $3.98e-07$, far less significant than on the original data ($2.73e-12$, Fig. 3). The average c-index was 0.74 ± 0.01 , a decline of 0.03 compared to the initial clustering with the echocardiography variables. A clustering solution obtained without the echocardiography variables obtained a c-index of 0.74 ± 0.01 and a log rank p-value of $4.66e-06 \pm 5.2e-06$. While the gaps are modest, they are statistically significant. Together, these tests suggest that clusters obtained with echocardiography data are more predictive of mortality.

As another test of the contribution, we repeated the process of choosing at random 31 of the 85 continuous variables and 2 of the 56 categorical and shuffling their values. The average of 50 random choices was 0.73 ± 0.03 , similar to just shuffling the echocardiography parameters. The average p-value of the log-rank test is $2.29e-04 \pm 8.88e-04$ with a median of $2.45e-08$. This means that the echocardiography data are roughly as meaningful for

Variable	Cluster 0		Cluster 1		Cluster 2		Cluster 3		p-value
	mean \pm std/ percentage	# of patients	mean \pm std/ percentage	# of patients	mean \pm std/ percentage	# of patients	mean \pm std/ percentage	# of patients	
<i>E' Lat</i>	11.32 \pm 2.76	126	8.82 \pm 2.61	172	6.50 \pm 2.01	110	6.55 \pm 1.73	57	1.16E-45
<i>E' Sept</i>	8.37 \pm 1.83	125	6.92 \pm 1.68	171	5.37 \pm 1.35	110	5.06 \pm 1.16	58	2.68E-45
<i>E/e average</i>	6.98 \pm 1.59	126	8.50 \pm 3.01	170	12.33 \pm 5.56	111	14.35 \pm 5.79	58	1.87E-35
<i>E/E' Sept</i>	8.02 \pm 1.89	115	9.34 \pm 3.09	165	13.15 \pm 6.39	104	16.21 \pm 6.88	50	3.54E-29
<i>E/E' Lat</i>	6.01 \pm 1.58	116	7.62 \pm 3.32	164	11.28 \pm 5.32	104	12.47 \pm 5.33	49	1.02E-28
A	49.99 \pm 10.23	125	62.76 \pm 17.22	173	76.31 \pm 20.72	98	69.13 \pm 20.06	45	1.35E-25
At	108.43 \pm 25.59	117	85.08 \pm 21.33	162	81.26 \pm 27.93	95	68.80 \pm 19.75	60	3.21E-24
IVSD	7.51 \pm 1.90	127	9.08 \pm 2.06	176	10.24 \pm 2.16	111	9.83 \pm 2.29	62	1.89E-21
<i>Diastolic Grade</i>	0.21 \pm 0.43	117	0.77 \pm 0.99	152	1.41 \pm 1.39	93	1.83 \pm 1.46	41	9.78E-20
<i>RA Pressure</i>	6.09 \pm 2.17	124	7.11 \pm 3.09	161	8.19 \pm 3.74	105	11.12 \pm 4.69	58	1.45E-18
<i>E/A</i>	1.35 \pm 0.40	125	1.03 \pm 0.32	173	0.87 \pm 0.26	98	1.20 \pm 0.67	45	1.95E-18
At < 100	33%(38)	115	77%(123)	160	77%(71)	92	93%(52)	56	2.45E-18
<i>LV mass</i>	111.06 \pm 36.50	125	151.03 \pm 51.31	169	140.04 \pm 50.25	108	171.47 \pm 67.81	61	3.18E-14
<i>TAPSE</i>	2.39 \pm 0.39	123	2.39 \pm 0.48	170	2.15 \pm 0.48	110	1.89 \pm 0.47	61	2.55E-13
<i>LA volume</i>	46.23 \pm 17.74	123	59.35 \pm 24.03	172	61.43 \pm 27.22	109	76.39 \pm 29.07	61	4.84E-13
<i>EF</i>	58.58 \pm 4.67	127	57.51 \pm 5.39	174	58.21 \pm 5.85	112	51.44 \pm 9.52	59	1.27E-12
<i>LAVI</i>	24.61 \pm 9.10	93	30.13 \pm 11.43	139	34.86 \pm 15.67	87	42.07 \pm 17.33	46	2.36E-12
<i>SV</i>	60.28 \pm 17.70	126	68.07 \pm 16.25	174	54.61 \pm 18.06	109	54.05 \pm 16.66	62	7.02E-11
<i>LVEDD</i>	43.74 \pm 5.72	126	45.90 \pm 5.53	177	40.83 \pm 6.87	110	46.75 \pm 8.23	62	7.62E-11
<i>RVED area</i>	19.89 \pm 4.32	99	21.99 \pm 4.58	143	18.82 \pm 4.70	90	23.65 \pm 5.39	56	1.04E-09
<i>LVESD</i>	27.89 \pm 4.98	124	29.92 \pm 5.81	175	26.77 \pm 5.84	106	32.33 \pm 9.63	62	9.95E-08
<i>RV S'</i>	11.08 \pm 1.94	122	12.03 \pm 2.59	173	10.80 \pm 3.12	104	9.73 \pm 3.02	60	1.06E-07
Bad heart condition (≥ 2)*	1%(1)	128	2%(4)	195	11%(12)	112	18%(13)	71	2.12E-07
<i>RVES area</i>	11.27 \pm 3.37	56	13.13 \pm 3.97	66	10.67 \pm 2.77	50	14.99 \pm 5.80	31	6.30E-06
<i>E</i>	64.89 \pm 14.21	127	62.52 \pm 17.36	178	67.68 \pm 25.86	112	77.82 \pm 22.56	61	7.66E-06
<i>CO</i>	4.38 \pm 1.22	126	5.69 \pm 4.45	170	4.03 \pm 1.28	107	4.23 \pm 1.45	59	7.58E-06
<i>Pericardial fluid</i>	0.07 \pm 0.26	127	0.08 \pm 0.28	179	0.22 \pm 0.42	112	0.30 \pm 0.49	61	7.60E-06
<i>E decel time</i>	166.00 \pm 35.04	115	180.34 \pm 50.32	167	198.18 \pm 56.85	102	165.64 \pm 55.32	53	1.19E-05
<i>SVI</i>	32.47 \pm 9.40	96	34.92 \pm 8.70	141	31.04 \pm 9.61	86	29.42 \pm 9.82	47	1.34E-03
<i>HR at Echo date</i>	73.75 \pm 12.35	126	78.92 \pm 14.33	174	74.51 \pm 15.09	108	80.00 \pm 17.51	61	2.98E-03
<i>CI</i>	2.35 \pm 0.66	96	3.00 \pm 2.74	138	2.31 \pm 0.76	84	2.26 \pm 0.77	45	8.29E-03
<i>EF Simp</i>	65.10 \pm 9.88	104	63.87 \pm 12.43	149	63.63 \pm 12.32	89	58.31 \pm 15.68	49	1.87E-02
<i>RVFAC CALC</i>	45.63 \pm 10.31	53	39.98 \pm 13.56	64	42.53 \pm 9.82	48	40.34 \pm 12.56	29	7.16E-02

Table 3. Statistics of all echocardiography variables. P-values were computed using ANOVA for continuous variables (in italics) and Chi2 for categorical variables (in bold), and FDR corrected for multiple testing. For the continuous variables mean \pm SD values in each cluster are presented. *Bad heart condition (≥ 2) value of 2 or above of at least one of the following: AS, AR, MS, MR, TS, TR, PS, PR.

forming the clustering as the rest of the parameters. It also shows redundancy in the contribution of different variables, as the changes are relatively small.

Echocardiography over time. Forty-eight patients who suffered clinical deterioration underwent multiple echocardiography measurements adjacent to the time of deterioration (only the first echocardiography measurement was used as input for the clustering). The differences in selected echocardiography variables between the first and second echocardiography results are shown in Supplementary 6. Due to the small sample size, the differences are not statistically significant, but some parameters show a trend.

Discussion

In the above trial, we have shown that hospitalized COVID-19 patients can be divided at admission into different clusters that have significant implications regarding disease course and final prognosis.

Although some aspects of the clusters seem self-evident (e.g., young patients have a better prognosis than old patients), others are less obvious, raising interest in the model's ability to influence our understanding of COVID-19's pathophysiology and progression and possibly help tailor treatments based on patient characteristics. The clusters significantly differed in natural COVID-19 parameters, such as CRP, chest X-ray findings and MEWS score, suggesting that the clusters are different in terms of the state of the disease. Moreover, the clusters

are also separated by mortality- and respiratory-related variables that were not part of the input data, showing good separation among the clusters in terms of these outcomes. A further analysis of the survival patterns shows that the clusters manifest very distinct survival patterns.

The main point of interest is the difference between clusters 1 and 2. Patients in cluster 1 were younger with fewer chronic diseases and less cardiac involvement on the one hand and higher levels of inflammatory markers and markers of pulmonary involvement on the other hand. Cluster 2 had the highest average patient age, but lower levels of inflammatory markers and notably was the only cluster enriched with females. This suggests that older infected males may be more prone to deteriorate, in comparison to females in the same age group. This is also supported by significantly higher rate of respiratory support provided to older males, in comparison to older females, across the whole cohort. An interesting finding is that pulmonary artery acceleration time was similar between clusters, possibly signifying that it is a marker of both pulmonary dysfunction due to increased afterload because of pulmonary vasoconstriction and cardiac intrinsic dysfunction. Although these clusters had very different starting points on admission, the patient's overall survival was similar in both. However, there were still significant differences in the need for hemodynamic support and respiratory support, with patients in cluster 1 necessitating more support. This difference might be explained by a different mechanism of deterioration between the two groups, possibly due to different pathological responses in the heart and lung systems between clusters. These differences show that the clusters are not only statistically and computationally justified but may also signify slightly different clinical entities.

Another noteworthy comparison is between cluster 1 and cluster 3. While the inflammatory state at admission was similar in the two clusters, the survival of cluster 3 was much worse, both when comparing the full clusters and when comparing only the ventilated patients. We can attribute the higher mortality to other factors, particularly age and comorbidity burden²⁸.

Variance in clinical deterioration of clusters. Another interesting finding is the difference between clusters of patients with clinical deterioration. Patients in clusters 1 and 3 had a decrease in stroke volume upon clinical deterioration together with an increase in right ventricular diastolic area (Table 3). This combination of an increase in RV preload (increase in right ventricular diastolic area) combined with a decrease in cardiac output (stroke volume) suggests that in patients in clusters 1 and 3, the right ventricle worked on the flat portion of the Frank–Starling curve³⁰, signifying significant right ventricular failure. This is in contrast with patients from cluster 2 who had an increase in stroke volume with unchanged RV area, suggesting normal RV filling pressure and high output, indicating that in patients from cluster 2, clinical deterioration was due to cytokine storm and vasodilatation and not due to RV failure.

Effect of echocardiography on clusters.. Although previous studies have tried to cluster COVID-19 patients, none have used echocardiography in the process. As previously demonstrated, echocardiography has added value over clinical characteristics alone in determining prognosis. In this study, we have shown that echocardiography data contributed to patient clustering in a modest but statistically significant manner and that clustering performed using echocardiographic data yields clusters that are better at predicting prognosis. However, conducting echocardiography is far more complex than obtaining the other clinical variables, the contribution of echocardiography, although significant, does not seem to justify the additional collection burden.

Limitations. The study was performed only on patients who were admitted to the hospital and thus cannot be generalized to the entire population of patients with COVID-19. However, hospitalized patients were needed to perform echocardiography on each patient, and admitted patients are usually higher risk patients and thus of higher interest to study.

Furthermore, the study was performed on patients admitted during the first wave of COVID-19, before the emergence of new treatments, the availability of vaccines and the appearance of new strains of COVID-19, such as the Delta and Omicron strains, and thus may not be relevant to the current manifestations of disease with newer treatments, vaccines and strains. Our group plans to perform a similar study on current patients to evaluate the clusters in a contemporary setting.

Additionally, this study was performed on patients admitted only to one center, and the results may not be generalizable to other cohorts of patients. Although we tried to find validation cohorts to compare the results and prevent overfitting, none were available that had comprehensive echocardiographic data. Furthermore, we used consensus clustering and carefully chose the clustering algorithm's hyperparameters to ensure robust clustering results.

Regarding the echocardiography contribution, although the current results do not support performing echocardiography at baseline in all patients, we did not include direct cardiac complications (such as myocardial infarction, myocarditis, and cardiac arrhythmias) as outcomes. Therefore, we cannot estimate the predictive ability of baseline echocardiography and cluster assignment of these conditions.

In conclusion, hospitalized COVID-19 patients were segregated into different clusters according to demographic, clinical and echocardiographic data at admission. These clusters of patients showed different disease courses and proved valuable in determining prognosis.

Data availability

The data described and analyzed in this study may be made available upon reasonable request to Y.T. (yant@tlvmc.gov.il).

Code availability

The code used in this study is available at https://github.com/Eranshp/clustering_echo_covid.git.

Received: 1 January 2023; Accepted: 18 May 2023

Published online: 31 May 2023

References

- Bhatraju, P. K. *et al.* Covid-19 in critically ill patients in the Seattle region: Case series. *N. Engl. J. Med.* **382**(21), 20132–20232. <https://doi.org/10.1056/nejmoa2004500> (2020).
- Szekely, Y. *et al.* Spectrum of cardiac manifestations in COVID-19: A systematic echocardiographic study. *Circulation* **142**(4), 342–353. <https://doi.org/10.1161/CIRCULATIONAHA.120.047971> (2020).
- Helms, J. *et al.* High risk of thrombosis in patients with severe SARS-CoV-2 infection: A multicenter prospective cohort study. *Intensive Care Med.* **46**(6), 1089–1098. <https://doi.org/10.1007/s00134-020-06062-x> (2020).
- Liotta, E. M. *et al.* Frequent neurologic manifestations and encephalopathy-associated morbidity in Covid-19 patients. *Ann. Clin. Transl. Neurol.* **7**(11), 2221–2230. <https://doi.org/10.1002/acn3.51210> (2020).
- Huang, C. *et al.* Clinical features of patients infected with 2019 novel coronavirus in Wuhan, China. *Lancet* **395**(10223), 497–506. [https://doi.org/10.1016/S0140-6736\(20\)30183-5](https://doi.org/10.1016/S0140-6736(20)30183-5) (2020).
- Clarkin, K. J. *et al.* COVID-19 and cardiovascular disease. *Circulation* **2020**, 1648–1655. <https://doi.org/10.1161/CIRCULATIONAHA.120.046941> (2020).
- Zweck, E. *et al.* Phenotyping cardiogenic shock. *J. Am. Heart Assoc.* **10**, 14. <https://doi.org/10.1161/JAHA.120.020085> (2021).
- Lalani, K., Yildirim, I., Phadke, V. K., Bednarczyk, R. A. & Omer, S. B. Assessment and validation of syndromic case definitions for respiratory syncytial virus infections in young infants: A latent class analysis. *Pediatr. Infect. Dis. J.* **38**, 1177–1182. <https://doi.org/10.1097/INF.0000000000002468> (2020).
- Vranas, K. C. *et al.* Identifying distinct subgroups of ICU patients: A machine learning approach. *Crit. Care Med.* **45**(10), 1607–1615. <https://doi.org/10.1097/CCM.0000000000002548> (2017).
- Sinha, P. *et al.* Prevalence of phenotypes of acute respiratory distress syndrome in critically ill patients with COVID-19: A prospective observational study. *Lancet Respir. Med.* **8**(12), 1209–1218. [https://doi.org/10.1016/S2213-2600\(20\)30366-0](https://doi.org/10.1016/S2213-2600(20)30366-0) (2020).
- Essay, P., Mosier, J. & Subbian, V. Phenotyping COVID-19 patients by ventilation therapy: Data quality challenges and cohort characterization. in *Public Health and Informatics: Proceedings of MIE 2021*, 198–202 (IOS Press, 2021). <https://doi.org/10.3233/SHTI210148>.
- Fauvel, C. *et al.* Cardiovascular manifestations secondary to COVID-19: A narrative review. *Respir. Med. Res.* **81**, 100904. <https://doi.org/10.1016/j.resmer.2022.100904> (2022).
- Kirkpatrick, J. N. *et al.* ASE statement on protection of patients and echocardiography service providers during the 2019 novel coronavirus outbreak: Endorsed by the American College of Cardiology. *J. Am. Coll. Cardiol.* **75**(24), 3078–3084. <https://doi.org/10.1016/j.jacc.2020.04.002> (2020).
- Lang, R. M. *et al.* Recommendations for cardiac chamber quantification by echocardiography in adults: An update from the American Society of Echocardiography and the European Association of Cardiovascular Imaging. *J. Am. Soc. Echocardiogr.* **28**(1), 1–39.e14. <https://doi.org/10.1016/j.echo.2014.10.003> (2015).
- Nagueh, S. F. *et al.* Recommendations for the evaluation of left ventricular diastolic function by echocardiography: An update from the American Society of Echocardiography and the European Association of Cardiovascular Imaging. *J. Am. Soc. Echocardiogr.* **29**(4), 277–314. <https://doi.org/10.1016/j.echo.2016.01.011> (2016).
- Topilsky, Y. *et al.* Preoperative factors associated with adverse outcome after tricuspid valve replacement. *Circulation* **123**(18), 1929–1939. <https://doi.org/10.1161/CIRCULATIONAHA.110.991018> (2011).
- Kitabatake, A. *et al.* Noninvasive evaluation of pulmonary hypertension by a pulsed Doppler technique. *Circulation* **68**(2), 302–309. <https://doi.org/10.1161/01.CIR.68.2.302> (1983).
- White, I. R., Royston, P. & Wood, A. M. Multiple imputation using chained equations: Issues and guidance for practice. *Stat. Med.* **30**(4), 377–399. <https://doi.org/10.1002/sim.4067> (2011).
- Yeo, I. & Johnson, R. A. *A New Family of Power Transformations to Improve Normality or Symmetry*, vol. 87. <https://academic.oup.com/biomet/article/87/4/954/232908>. (2000).
- Huang, Z. Extensions to the k-means algorithm for clustering large data sets with categorical values. *Data Min. Knowl. Disc.* **12**, 283–304. <https://doi.org/10.1023/A:1009769707641> (1998).
- Preud'homme, G. *et al.* Head-to-head comparison of clustering methods for heterogeneous data: A simulation-driven benchmark. *Sci. Rep.* **11**(1), 83340. <https://doi.org/10.1038/s41598-021-83340-8> (2021).
- Monti, S. *et al.* Consensus clustering: A resampling-based method for class discovery and visualization of gene expression microarray data. *Mach. Learn.* **52**, 91–118. <https://doi.org/10.1023/A:1023949509487> (2003).
- Virtanen, P. *et al.* SciPy 1.0: Fundamental algorithms for scientific computing in Python. *Nat. Methods.* **17**(3), 261–272. <https://doi.org/10.1038/s41592-019-0686-2> (2020).
- Kaplan, E. L. & Meier, P. Nonparametric estimation from incomplete observations. *J. Am. Stat. Assoc.* **53**(282), 457–481. <https://doi.org/10.1080/01621459.1958.10501452> (1958).
- Cox, D. R. Regression models and life-tables. *J. R. Stat. Soc. B* **34**(2), 187–220 (1972).
- Harrell, F. E., Califf, R. M., Pryor, D. B., Lee, K. L. & Rosati, R. A. Evaluating the yield of medical tests. *JAMA* **247**(18), 2543–2546. <https://doi.org/10.1001/jama.1982.03320430047030> (1982).
- Davidson-Pilon, C. Lifelines: Survival analysis in Python. *J. Open Source Softw.* **4**(40), 1317. <https://doi.org/10.21105/joss.01317> (2019).
- Sonaglioni, A. *et al.* Charlson comorbidity index, neutrophil-to-lymphocyte ratio and undertreatment with renin-angiotensin-aldosterone system inhibitors predict in-hospital mortality of hospitalized COVID-19 patients during the omicron dominant period. *Front. Immunol.* **13**, 958418. <https://doi.org/10.3389/fimmu.2022.958418> (2022).
- Barnett, W. R. *et al.* Initial MEWS score to predict ICU admission or transfer of hospitalized patients with COVID-19: A retrospective study. *J. Infect.* **82**(2), 282–327 (2021).
- Taieb, P. *et al.* Risk prediction in patients with COVID-19 based on haemodynamic assessment of left and right ventricular function. *Eur. Heart J. Cardiovasc. Imaging* **22**(11), 1241–1254. <https://doi.org/10.1093/EHJCI/JEAB169> (2021).

Acknowledgements

This study was supported in part by the Israel Science Foundation (Grant No. 3165/19), within the Israel Precision Medicine Partnership program, and by a grant from the Tel Aviv University Center for AI and Data Science (TAD). ES was supported in part by a fellowship from the Edmond J. Safra Center for Bioinformatics at Tel-Aviv University.

Author contributions

Design, analysis and writing: E.S., A.H.; assistance in design: D.C.; data collection: I.M., E.G., Y.S., Y.L., P.T., A.B., O.S., Y.G., L.L., A.B., S.S., S.B.; assistance in interpretation of results: R.R.; design, writing and supervision: Y.T., R.S.

Competing interests

The authors declare no competing interests.

Additional information

Supplementary Information The online version contains supplementary material available at <https://doi.org/10.1038/s41598-023-35449-1>.

Correspondence and requests for materials should be addressed to R.S.

Reprints and permissions information is available at www.nature.com/reprints.

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2023