

Artificial Intelligence-assisted quantification of COVID-19 pneumonia burden from computed tomography improves prediction of adverse outcomes over visual scoring systems

Short title: AI quantification of COVID-19 pneumonia

Kajetan Grodecki MD,^{1,2} Aditya Killekar BS,³ Judit Simon MD,^{4,5} Andrew Lin MD,¹ Sebastien Cadet MS,³ Priscilla McElhinney BS,¹ Cato Chan MD,⁶ Michelle C. Williams, MBChB, PhD,⁷ Barry D Pressman MD,⁶ Peter Julien MD,⁶ Debiao Li MD,¹ Peter Chen MD,⁸ Nicola Gaibazzi MD,⁹ Udit Thakur MD,¹⁰ Elisabetta Mancini MD,¹¹ Cecilia Agalbato MD,¹¹ Jiro Munechika MD,¹² Hidenari Matsumoto MD,¹³ Roberto Menè MD,^{14,15} Gianfranco Parati MD,^{14,15} Franco Cernigliaro MD,^{14,15} Nitesh Nerlekar MD,¹⁰ Camilla Torlasco MD,^{14,15} Gianluca Pontone MD,¹¹ Pal Maurovich-Horvat MD,^{4,5} Piotr J Slomka PhD,³ Damini Dey PhD¹

¹ Biomedical Imaging Research Institute, Cedars-Sinai Medical Center, Los Angeles, CA, USA

² 1st Department of Cardiology, Medical University of Warsaw, Warsaw, Poland

³ Division of Artificial Intelligence in Medicine, Department of Medicine, Cedars-Sinai Medical Center, Los Angeles, CA, USA

⁴ Department of Radiology, Medical Imaging Centre, Semmelweis University, Budapest, Hungary

⁵ MTA-SE Cardiovascular Imaging Research Group, Semmelweis University, Budapest, Hungary

⁶ Department of Imaging, Cedars-Sinai Medical Center, USA

⁷ BHF Centre for Cardiovascular Science, University of Edinburgh, Edinburgh, United Kingdom

⁸ Department of Medicine, Women's Guild Lung Institute, Cedars-Sinai Medical Center, Los Angeles, CA, USA

⁹ Cardiology, Azienda Ospedaliero-Universitaria di Parma, Parma, Italy

¹⁰ Monash Health, Melbourne, Australia

¹¹ Centro Cardiologico Monzino IRCCS, University of Milan, Italy

¹² Division of Radiology, Showa University School of Medicine, Tokyo, Japan

¹³ Division of Cardiology, Showa University School of Medicine, Tokyo, Japan

¹⁴ Department of Cardiovascular, Neural and Metabolic Sciences, IRCCS Istituto Auxologico Italiano, Milan, Italy

¹⁵ Department of Medicine and Surgery, University of Milano-Bicocca, Italy

Address for correspondence:

Damini Dey, PhD

Professor and Research Scientist, Director of Quantitative Image Analysis,

Biomedical Imaging Research Institute, Cedars-Sinai Medical Center
116 N Robertson Boulevard, Los Angeles, CA 90048
Email: damini.dey@cshs.org

Type of Manuscript: Full Paper

Authors report no conflict of interest.

Funding

This research was supported by Cedars-Sinai COVID-19 funding. This research was also supported by the National Heart, Lung, and Blood Institute of the National Institutes of Health (R01HL133616). The content is solely the responsibility of the authors and does not necessarily represent the official views of the National Institutes of Health. Dr. Kajetan Grodecki was supported by the Foundation for Polish Science (FNP).

Acknowledgements: We would like to thank all the individuals involved in the collection, processing, and analysis of data in this international multicenter study.

Abstract (250/250 words)

Objectives:

We sought to assess the performance of artificial intelligence (AI)-assisted quantification of pneumonia burden from chest computed tomography (CT) for predicting clinical deterioration or death in patients hospitalized with COVID-19 in comparison to semi-quantitative visual scoring systems.

Methods:

Total pneumonia burden was quantified using a deep-learning algorithm and semi-quantitative pneumonia severity scores were visually estimated. The primary outcome was

clinical deterioration (intensive care unit admission, invasive mechanical ventilation, or vasopressor therapy) or in-hospital death.

Results:

The final population comprised 743 patients (mean age 65 ± 17 years, 55% men), of whom 175 (23.5%) experienced clinical deterioration or death. The area under the receiver operating characteristic curve (AUC) for predicting the primary outcome was significantly higher for AI-assisted quantitative pneumonia burden (0.739, $p=0.021$) compared with the visual lobar severity score (0.711, $p<0.001$) and visual segmental severity score (0.722, $p=0.042$). AI-assisted pneumonia assessment exhibited lower performance when applied for calculation of the lobar severity score (AUC of 0.723, $p=0.021$). Time taken for AI-assisted quantification of pneumonia burden was lower (38 ± 10 seconds) compared to that of visual lobar (328 ± 54 seconds, $p<0.001$) and segmental (698 ± 147 seconds, $p<0.001$) severity scores.

Conclusions:

AI-assisted quantification of pneumonia burden from chest CT improves prediction of clinical deterioration in COVID-19 patients over semi-quantitative severity scores, at a fraction of the analysis time.

Advances in knowledge:

Quantitative pneumonia burden assessed using AI demonstrated higher performance for predicting clinical deterioration compared to current semi-quantitative scoring systems. Such an AI system has the potential to be applied for image-based triage of COVID-19 patients in clinical practice.

**Artificial Intelligence-assisted quantification of COVID-19 pneumonia burden from
computed tomography improves prediction of adverse outcomes over visual scoring
systems**

Short title: AI quantification of COVID-19 pneumonia

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

BJR UNCORRECTED PROOFS

Abstract (250/250 words)

Objectives:

We aimed to evaluate the effectiveness of utilizing artificial intelligence (AI) to quantify the extent of pneumonia from chest computed tomography (CT) scans, and to determine its ability to predict clinical deterioration or mortality in patients admitted to the hospital with COVID-19 in comparison to semi-quantitative visual scoring systems.

Methods:

A deep-learning algorithm was utilized to quantify the pneumonia burden, while semi-quantitative pneumonia severity scores were estimated through visual means. The primary outcome was clinical deterioration, the composite endpoint including admission to the intensive care unit, need for invasive mechanical ventilation, or vasopressor therapy, as well as in-hospital death.

Results:

The final population comprised 743 patients (mean age 65 ± 17 years, 55% men), of whom 175 (23.5%) experienced clinical deterioration or death. The area under the receiver operating characteristic curve (AUC) for predicting the primary outcome was significantly higher for AI-assisted quantitative pneumonia burden (0.739, $p=0.021$) compared with the visual lobar severity score (0.711, $p<0.001$) and visual segmental severity score (0.722, $p=0.042$). AI-assisted pneumonia assessment exhibited lower performance when applied for calculation of the lobar severity score (AUC of 0.723, $p=0.021$). Time taken for AI-assisted quantification of pneumonia burden was lower (38 ± 10 seconds) compared to that of visual lobar (328 ± 54 seconds, $p<0.001$) and segmental (698 ± 147 seconds, $p<0.001$) severity scores.

Conclusions:

1 Utilizing AI-assisted quantification of pneumonia burden from chest CT scans offers a more
2 accurate prediction of clinical deterioration in patients with COVID-19 compared to semi-
3 quantitative severity scores, while requiring only a fraction of the analysis time.
4
5

6
7 **Advances in knowledge:**
8

9 Quantitative pneumonia burden assessed using AI demonstrated higher performance for
10 predicting clinical deterioration compared to current semi-quantitative scoring systems. Such
11 an AI system has the potential to be applied for image-based triage of COVID-19 patients in
12 clinical practice.
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1. Introduction

Coronavirus disease 2019 (COVID-19) is a global pandemic and has caused public health crises of catastrophic proportions, with over 395 million confirmed cases worldwide as of February 7, 2022. Although a reverse transcription-polymerase chain reaction (RT-PCR) test in respiratory tract specimens is necessary for diagnosing COVID-19, computed tomography (CT) remains the primary modality used to assess the extent of the disease and predict its progression.^{1,2} Several semi-quantitative scoring systems have been proposed to visually evaluate parenchymal opacifications associated with COVID-19, in which the final metric is obtained by summing the scores describing the extent of abnormalities in respective pulmonary lobes or segments.^{3,4} While conventional visual scoring of COVID-19 pneumonia extent correlates with clinical disease severity, its routine application is time-consuming and requires proficiency in cardiothoracic imaging.⁵ To aid radiologists and enhance their diagnostic accuracy, various artificial intelligence (AI) solutions have been developed.⁶ Deep learning algorithms, a form of AI, have demonstrated high performance in fully automated segmentation of lung lesions associated with COVID-19 and identifying patients at risk of experiencing adverse clinical outcomes.^{7,8} Whereas several studies have evaluated the diagnostic accuracy of AI-assisted and conventional pneumonia scoring systems, there is a paucity of data comparing the prognostic value of these approaches. In this retrospective analysis of a large, international, multicenter registry, our aim was to evaluate the effectiveness of using AI-assisted quantitative pneumonia burden from chest CT scans to predict clinical deterioration or mortality in patients hospitalized with COVID-19. We compared this approach to semi-quantitative visual scoring systems.

2. Materials and methods

2.1 Study design

This prospective, international, multicenter registry included patients enrolled consecutively from: North America (Cedars Sinai Medical Center, Los Angeles, USA [n = 41]), Europe (Semmelweis University, Budapest Hungary [n = 579]; Centro Cardiologico Monzino [n = 75], and Istituto Auxologico Italiano [n = 17 both Milan, Italy), Asia (Showa University Hospital, Tokyo, Japan [n = 25]), and Australia (Monash Medical Centre, Victoria, Australia [n = 6]). All patients underwent baseline chest CT and had a positive RT-PCR test result for SARS-CoV-2 during their index admission between January 10 and November 15, 2020 (Figure 1). For patients with serial chest CT imaging, we included only the results of their initial scan. The CT images from each patient and the clinical database were fully anonymized and transferred to Cedars-Sinai Medical Center for core lab analysis. The study was conducted with the approval of local institutional review boards (Cedars-Sinai Medical Center IRB# study 617), and written informed consent was waived for fully anonymized data analysis.

2.2 Scan Protocol and Image Reconstruction

Chest CT scans were conducted using various multi-slice CT systems, including the Aquilion ONE (Toshiba Medical Systems, Otawara, Japan), GE Revolution, GE Discovery CT750 HD, or LightSpeed VCT (GE Healthcare, Milwaukee, WI, USA), and Brilliance iCT and Incisive CT (Philips Healthcare, Cleveland, OH, USA). Scans without intravenous contrast utilized a peak x-ray tube voltage of 120 kV, automatic tube current modulation (300-500 mAs), and a slice thickness of 0.625 to 1.25 mm. The contrast-enhanced protocol consisted of a peak x-ray tube voltage of 120 kV, automatic tube current modulation (500-650 mAs), and a slice thickness of 0.625 to 1.0 mm. Iodinated contrast material (Iomeron 400 and 350, Bracco Imaging SpA, Milan, Italy; or Ominpaque 350, GE Healthcare, United States) totaling 80-100 ml was injected

1 intravenously at a rate of 5 ml/s and followed by 20-30ml of saline chaser at a flow rate of 4-5
2 ml/s. Standard lung filters specific to each CT vendor were used to reconstruct the images. All
3 scans were obtained while patients were in the supine position during an inspiratory breath-
4 hold.
5
6
7
8

9 **2.3 CT Image Analysis**

10 Images were analyzed at the Cedars-Sinai Medical Center core laboratory by two physicians
11 (K.G. and A.L.) with 3 and 8 years of experience in chest CT, respectively, and who were
12 blinded to clinical data. A standard lung window (width of 1500 Hounsfield units [HU] and
13 level of -400 HU) was used (Figure 2A-B).
14
15
16
17
18
19

20 At the core laboratory of Cedars-Sinai Medical Center, two physicians (K.G. and A.L.) with 3
21 and 8 years of experience in chest CT, respectively, analyzed the images and were not provided
22 with clinical data. The standard lung window with a width of 1500 Hounsfield units [HU] and
23 a level of -400 HU was utilized (as shown in Figure 2A-B).
24
25
26
27
28
29

30 For AI-assisted pneumonia burden quantification, deep-learning research software
31 (LungQuant v.1.0, Cedars-Sinai Medical Center, Los Angeles, CA, USA) was used. First
32 ground glass opacities (GGO) and high-opacities (comprising consolidation and pleural
33 effusion) were segmented using convolutional Long Short-Term Memory (ConvLSTM)
34 network (Figure 2C).⁸ ConvLSTM was utilized in our study as it operates directly on images,
35 allowing for quick segmentation and precise 3D quantification of lung lesions involved in
36 COVID-19 pneumonia from a stack of both contrast and non-contrast CT images. The
37 ConvLSTM networks can preserve relevant features while dismissing irrelevant ones through
38 the feedback loop, resulting in a memory-efficient approach for the comprehensive analysis of
39 the images.
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54

55 Following the acquisition of lesion masks, they were modified using a semi-automated
56 brush-like tool to distinguish consolidation from pleural effusion, with the boundaries delimited
57
58
59
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

by a region-growing algorithm. Adaptive thresholds were implemented, determined by a fixed window surrounding the attenuation of the clicked pixel by the operator. Lobe segmentation was performed using a deep convolutional neural network, an alternative type of neural network that models spatial and temporal correlation while reducing translational variance in signals. The size of input images is used to construct deep convolutional neural networks. The networks were trained with the Lung Tissue Research Consortium dataset.⁹ The right lung was divided into upper, middle and lower lobes by the horizontal and oblique fissures, and the left lung was divided into upper and lower lobes by the oblique fissure (Figure 2D). GGO was defined as hazy opacities that do not obscure the underlying bronchial or vascular structures, consolidation as opacification obscuring the underlying bronchial and vascular structures, and pleural effusion as a fluid collection in the pleural cavity.¹⁰ Chronic lung abnormalities such as emphysema or fibrosis were excluded from segmentation. Volumes of lesion components and total lesion volumes were automatically calculated by the software. Total pneumonia burden was calculated as total lesion volume / total lung volume x 100% (Figure 2E-F). AI calculations were performed on Nvidia Titan RTX 24GB graphics processing unit.

For lobar severity score, the extent of the parenchymal opacities involving GGO, and consolidations were visually assessed for each of the 5 pulmonary lobes and scores ranging from 0 to 5 were attributed accordingly: 0 for no involvement; 1 for involvement <0%;5%); 2 for involvement <5%;25%); 3 for involvement <25%;50%); 4 for involvement <50%;75%); and 5 for involvement $\geq 75\%$. The total lobar severity score ranged between 0 and 25 points.³

To compare the performance of the lobar severity score between expert reader and AI, lobar involvement of the opacifications calculated using a deep-learning algorithm was translated into semi-quantitative scores as described.

For segmental severity score, the extent of the parenchymal opacities involving GGO and consolidations were visually assessed for each of the 20 pulmonary segments and scores

1 ranging from 0 to 2 were attributed accordingly: 0 for no involvement; 1 for involvement $\leq 50\%$;
2 and 2 for involvement $> 50\%$. The total segmental severity score ranged between 0 and 40
3 points.⁴
4
5
6

7 All the cases were evaluated using each of the approaches. To limit the bias, images
8 were first scored using a semi-quantitative approach. The minimal interval between the repeated
9 evaluation of the case was 4 weeks. The time necessary to score the case using each of the three
10 approaches was noted for all the patients.
11
12
13
14
15
16
17
18

19 **2.4 Statistical analysis**

20 Normal distribution of the data was assessed using the Shapiro-Wilk test. Continuous variables
21 were reported as mean \pm standard deviation or median (interquartile range [IQR]), while
22 categorical variables were expressed as absolute numbers (percentage). Student's t-test or
23 nonparametric Mann-Whitney U-test was used to compare continuous variables as appropriate,
24 while categorical variables were compared using the Chi-square test. Discriminatory
25 performance of the scores was determined by the C-statistic, and compared using the method
26 of DeLong et al.¹¹ Optimal sensitivity and specificity were determined by the Youden index to
27 facilitate selection of the cut-off values. Additionally, the scoring systems were divided into
28 quartiles and the frequency of clinical deterioration was compared using the odds ratio. A
29 reclassification table was constructed only for the comparison of pulmonary lobes evaluation
30 since AI software does not provide per-segment quantification. A reclassification table was
31 constructed only for the comparison of pulmonary lobes evaluation since AI software does not
32 provide per-segment quantification. Reclassification table was constructed to visualize the
33 directions (described as up- and down-reclassifications) and frequencies of reclassifications as
34 performed by the AI algorithm in relation to the clinical standard being visual assessment. Net
35 reclassification improvement (NRI) was calculated using the method described by Pencina et
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

al. to objectivize its frequency.¹² Agreement between AI-assisted and visual analysis of lobar involvement was evaluated with a weighted kappa statistic.¹³ The degree of agreement was considered excellent for kappa >0.80; substantial for kappa 0.61-0.80; moderate for kappa 0.41-0.60; fair for kappa 0.21-0.40; and poor for kappa ≤0.20. Correlations between continuous variables were assessed using Spearman's rank correlation coefficient. All probability values were 2-tailed, and a p-value of <0.05 was considered statistically significant. Data were processed using the SPSS software, version 23 (IBM SPSS Statistics, Newer York, USA) and SAS 9.4 (SAS Institute, Cary, USA).

3. Results

3.1 Patient characteristics

A total of 743 patients (age 65 ± 17 years; 55% male) with laboratory-confirmed COVID-19 who underwent chest CT during their admission were included. The primary outcome occurred in 175 (23.5%) patients: 93 (53.2%) were admitted to ICU, 70 (40.0%) required mechanical ventilation, 64 (36.5%) required vasopressors, and 121 (69.1%) experienced in-hospital death. The chest CT was performed at a median time of 6 days (IQR 4-8 days) from self-reported onset of symptoms, and the median time from chest CT to occurrence of primary outcome was 3 days (IQR 1-13 days). Of the total patients, 79 (10.6%) patients experienced clinical deterioration or died, while the remaining patients (n=568; 76.5%) did not require critical care or had been discharged alive at the time of data collection. Patients who experienced deterioration or died were older and had a higher number of comorbidities, as shown in Table 1.

3.2 Chest CT measurements

Lung measurements from chest CT are summarized in Table 2 and Figure 3. Patients with deterioration or death had a higher total burden of COVID-19 pneumonia compared to patients that did not experience deterioration or death (16.0% [IQR, 4.5-39.3%] vs 3.7% [IQR, 0.3-

10.3%], $p < 0.001$). Similarly, patients who deteriorated or died were characterized with a higher visual lobar (10 [IQR, 6-15] vs 6 [IQR, 3-9], $p < 0.001$) as well as segmental (18 [IQR, 10-27] vs 9 [IQR, 4-15], $p < 0.001$) severity scores compared to patients who did not require critical care or were discharged alive. Data regarding the distribution of individual score components for lobar and segmental severity scores are presented in Supplementary Tables 1 and 2. Time required for calculation of AI-assisted pneumonia burden (38 ± 10 sec) was significantly lower compared to both visual lobar (328 ± 54 sec, $p < 0.001$) and segmental (698 ± 147 sec, $p < 0.001$) severity scores.

3.3 Predictive accuracy of pneumonia scoring systems

For the prediction of the primary outcome, the area under the receiver operating characteristic curve (AUC) for AI-assisted pneumonia burden (0.739, $p = 0.021$) was significantly higher than that of the visual lobar severity score (0.711, $p < 0.001$) and visual segmental severity score (0.722, $p = 0.042$; Figure 4). The sensitivities and specificities were: 68% and 70% for AI-assisted pneumonia burden; 48% and 84% for lobar severity score; was 51% and 83% for segmental severity score.

The frequency of clinical deterioration in each of the quartiles for respective scoring systems are presented in Table 3. For AI-assisted pneumonia burden, the odds of clinical deterioration were: 8.46 (95% CI: 4.85-14.73) for quartile 4 versus quartile 1, 5.70 (95% CI: 3.46-9.40) for quartile 4 versus quartile 2, and 3.80 (95% CI: 2.40-6.00) for quartile 4 versus 3. For visual lobar severity score, the odds of clinical deterioration were: 6.50 (95% CI: 3.86-10.93) for quartile 4 versus quartile 1, 5.90 (95% CI: 3.52-9.88) for quartile 4 versus quartile 2, and 3.40 (95% CI: 2.16-5.34) for quartile 4 versus 3. For visual segmental severity score, the odds of clinical deterioration were: 7.27 (95% CI: 4.31-12.26) for quartile 4 versus quartile 1, 5.74 (95% CI: 3.54-9.31) for quartile 4 versus quartile 2, and 3.56 (95% CI: 2.27-5.58) for quartile 4 versus 3.

3.4 Agreement between visual and AI-assisted analysis

To compare the performance of the lobar severity score between expert reader and AI, lobar involvement of the opacifications calculated using a deep-learning algorithm was translated into semi-quantitative scores as described. The agreement between visual and AI-assisted lobar severity scores was substantial (weighted kappa = 0.609). AI-derived lobar severity score achieved higher predictive accuracy for clinical deterioration than visual expert reading (AUC of 0.723 vs 0.711, $p = 0.043$; Figure 5A), but underperformed compared to the input pneumonia burden (AUC of 0.723 vs 0.739, $p = 0.021$). The estimation of pneumonia involvement for individual lobes showed an excellent level of agreement (weighted kappa = 0.862; Figure 5B). Out of 3715 lobes, discordant classification was noted in 387 (10.4%) of them. The NRI values were 0.3% and 0.4% in patients without and with clinical deterioration, respectively. Thus, the total NRI was -0.1% ($p = 0.994$, Table 4).

3.5 Correlation of scoring systems with serum biomarkers

Bivariate correlations between pneumonia scoring systems and serum biomarkers are presented in Table 5. Serum biomarkers were more strongly correlated with AI-assisted pneumonia burden than the semi-quantitative severity scores. The pneumonia scoring systems had a moderate correlation with lactate dehydrogenase and C-reactive protein levels, and weak correlations with lymphocytes, ferritin, D-dimer, and creatine kinase-MB.

4. Discussion

In this international multicenter study of patients with COVID-19, we compared the accuracy of an AI-assisted pneumonia burden with conventional pneumonia severity scores derived from chest CT for predicting clinical deterioration. We demonstrate that quantitative pneumonia burden determined with an AI achieved higher predictability of clinical deterioration to the semi-quantitative visual scoring systems and significantly reduced the time required for pneumonia evaluation from chest CT.

Chest CT is presently recommended for COVID-19 patients who exhibit moderate or severe respiratory symptoms, have a high pretest probability of infection, or require urgent triage for other clinical scenarios.¹⁴⁻¹⁶ To facilitate the standardized evaluation of pulmonary involvement from CT, several different severity scores have been proposed.^{3,4,9} The semi-quantitative scoring systems – developed originally to describe the idiopathic pulmonary fibrosis and adapted later for CT examination of patients recovering from severe acute respiratory syndrome – have been recently shown to associate with the clinical disease severity and adverse outcomes in COVID-19 patients.^{3,17,18} Although visual analysis of lung involvement is the only available approach to many of the institutions, its application remains limited to the staff proficient at cardiothoracic imaging. Moreover, the reproducibility of the measurements may be dependent on the experience of the individual reader, and the scoring was showed to differ significantly between radiologists and clinicians.^{7,19}

Alternatively, the extent of pneumonia can be characterized using quantitative measurements, which require segmentation of both lungs and parenchymal lesions.²⁰ While the manual approach is prohibitively time-consuming and could not be employed in a routine clinical setting, the application of deep learning – a class of AI – has been demonstrated as a robust tool generating results with an accuracy similar to the experts.²¹ Previous studies have

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

proved AI to increase the performance of junior radiologists to the senior level as well as improve diagnostic accuracy in comparison to visual severity scores.⁷

In spite of the central role of CT in prognostication, a paucity of data remains on the incremental value of AI-assisted pneumonia burden quantification in comparison to semi-quantitative scores for prediction of clinical deterioration in COVID-19 patients admitted to the hospital. The superior performance of the AI in comparison to the expert reader was previously found by the Gieraerts et al. in a single-center study involving 250 COVID-19 patients, although they did not report a significant difference in prognostic accuracy between quantitative pneumonia burden (AUC of 0.878) and semi-quantitative (AUC of 0.888) measured with an AI.²² This could be associated with a relatively small number of events (n = 39) limiting their statistical power since such a high accuracy could not be achieved even by the model combining multiple quantitative lesion features in a landmark study by Zhang et al (AUC of 0.848).⁷ In a three-fold larger study involving real-life data from five continents, we report improved prognostication with AI-assisted pneumonia burden as compared to visual estimation, but also AI-derived semi-quantitative severity scores. Although the agreement between AI and visual estimation of pulmonary involvement in COVID-19 was excellent, the reclassification table for the comparison of pulmonary lobes evaluation showed the tendency of visual scoring to overestimate the disease burden in severe cases. While the NRI remained unaffected, the results confirm previously observed bias in the visual estimation of lung abnormalities, which may negatively affect the overall performance of the scoring.²³ Further, observed decrease in the prognostic value of the AI-assisted measurements following the translation of pneumonia burden into the lobar severity score, suggests semi-quantitative scales being naturally limited by the categorization of the continuous data.²⁴

Our results also showed the correlation of pneumonia severity scores with blood biomarkers related to systemic inflammation, thus underscoring the importance of lung

1 involvement as the key parameter in the overall prognostic implications. The strongest
2 correlations were found for C-reactive protein and lactate dehydrogenase. The first indicates
3 the association of lung injury with acute inflammation, and the latter – being the marker of liver
4 function – may suggest its role in the pathogenesis of multi-organ failure.²⁵⁻²⁷ Although the
5 relation between pulmonary inflammation in COVID-19 and the pathogenetic sequelae
6 resulting in clinical deterioration is not fully understood, this may mechanistically explain the
7 prognostic value of chest CT imaging.
8
9

10
11
12
13
14
15
16
17 There are several limitations to our study. Firstly, there may have been heterogeneity in
18 COVID-19 pneumonia severity or in-hospital outcomes due to different patient profiles and
19 treatment protocols across countries. Secondly, information on patients' respiratory status upon
20 admission to the intensive care unit was not consistently available, but we included intubation
21 and invasive ventilation as a hard endpoint. Finally, we did not investigate the impact of
22 treatment on outcomes; however, supportive care remains the cornerstone of COVID-19
23 therapy, and only a small number of patients received targeted interventions in our study.
24
25
26
27
28
29
30
31
32

33 34 **5. Conclusion**

35
36 We show that the AI-assisted quantitative pneumonia burden outperforms semi-quantitative
37 severity scores for prediction of clinical deterioration in COVID-19 patients, which also
38 validates the application of AI for lessening the workload in the radiology departments.²⁸ The
39 presented deep-learning algorithm requires little to no interaction, facilitating, therefore, the
40 rapid risk assessment by clinicians with limited experience in cardiothoracic imaging.
41 Quantification of the parenchymal opacification on chest CT might be applied for image-based
42 triage to optimize the distribution of resources during the pandemic.
43
44
45
46
47
48
49
50

51 In conclusion, AI-assisted pneumonia burden improves the prediction of clinical
52 deterioration in COVID-19 patients as compared to semi-quantitative severity scores and may
53 significantly expedite CT-based triage in the emergency environment.
54
55
56
57
58
59
60
61
62
63
64
65

References

1. Khatami F, Saatchi M, Zadeh SST, et al. A meta-analysis of accuracy and sensitivity of chest CT and RT-PCR in COVID-19 diagnosis. *Sci Rep* 2020; **10**(1): 22402.
2. Pontone G, Scafuri S, Mancini ME, et al. Role of computed tomography in COVID-19. *J Cardiovasc Comput Tomogr* 2021; **15**(1): 27-36.
3. Li K, Fang Y, Li W, et al. CT image visual quantitative evaluation and clinical classification of coronavirus disease (COVID-19). *Eur Radiol* 2020; **30**(8): 4407-16.
4. Yang R, Li X, Liu H, et al. Chest CT Severity Score: An Imaging Tool for Assessing Severe COVID-19. *Radiology: Cardiothoracic Imaging* 2020; **2**(2): e200047.
5. Guan X, Yao L, Tan Y, et al. Quantitative and semi-quantitative CT assessments of lung lesion burden in COVID-19 pneumonia. *Scientific Reports* 2021; **11**(1): 5148.
6. Chartrand G, Cheng PM, Vorontsov E, et al. Deep Learning: A Primer for Radiologists. *RadioGraphics* 2017; **37**(7): 2113-31.
7. Zhang K, Liu X, Shen J, et al. Clinically Applicable AI System for Accurate Diagnosis, Quantitative Measurements, and Prognosis of COVID-19 Pneumonia Using Computed Tomography. *Cell* 2020; **181**(6): 1423-33.e11.
8. Grodecki K, Killekar A, Lin A, et al. Rapid quantification of COVID-19 pneumonia burden from computed tomography with convolutional LSTM networks. *ArXiv* 2021: arXiv:2104.00138v1.
9. Grodecki K, Lin A, Cadet S, et al. Quantitative Burden of COVID-19 Pneumonia on Chest CT Predicts Adverse Outcomes: A Post-Hoc Analysis of a Prospective International Registry. *Radiology: Cardiothoracic Imaging* 2020; **2**(5): e200389.
10. Hansell DM, Bankier AA, MacMahon H, McLoud TC, Müller NL, Remy J. Fleischner Society: glossary of terms for thoracic imaging. *Radiology* 2008; **246**(3): 697-722.
11. DeLong ER, DeLong DM, Clarke-Pearson DL. Comparing the areas under two or more correlated receiver operating characteristic curves: a nonparametric approach. *Biometrics* 1988; **44**(3): 837-45.
12. Pencina MJ, D'Agostino RB, Sr., Steyerberg EW. Extensions of net reclassification improvement calculations to measure usefulness of new biomarkers. *Stat Med* 2011; **30**(1): 11-21.
13. Landis JR, Koch GG. The Measurement of Observer Agreement for Categorical Data. *Biometrics* 1977; **33**(1): 159-74.
14. Rubin GD, Ryerson CJ, Haramati LB, et al. The Role of Chest Imaging in Patient Management during the COVID-19 Pandemic: A Multinational Consensus Statement from the Fleischner Society. *Radiology* 2020; **296**(1): 172-80.
15. Choi AD, Abbara S, Branch KR, et al. Society of Cardiovascular Computed Tomography guidance for use of cardiac computed tomography amidst the COVID-19 pandemic Endorsed by the American College of Cardiology. *Journal of Cardiovascular Computed Tomography* 2020; **14**(2): 101-4.
16. Pontone G, Baggiano A, Conte E, et al. "Quadruple Rule-Out" With Computed Tomography in a COVID-19 Patient With Equivocal Acute Coronary Syndrome Presentation. *JACC Cardiovasc Imaging* 2020; **13**(8): 1854-6.
17. Kazerooni EA, Martinez FJ, Flint A, et al. Thin-section CT obtained at 10-mm increments versus limited three-level thin-section CT for idiopathic pulmonary fibrosis: correlation with pathologic scoring. *AJR Am J Roentgenol* 1997; **169**(4): 977-83.
18. Chang Y-C, Yu C-J, Chang S-C, et al. Pulmonary Sequelae in Convalescent Patients after Severe Acute Respiratory Syndrome: Evaluation with Thin-Section CT. *Radiology* 2005; **236**(3): 1067-75.

19. Li S, Liu S, Wang B, et al. Predictive value of chest CT scoring in COVID-19 patients in Wuhan, China: A retrospective cohort study. *Respiratory Medicine* 2021; **176**: 106271.
20. Colombi D, Bodini FC, Petrini M, et al. Well-aerated Lung on Admitting Chest CT to Predict Adverse Outcome in COVID-19 Pneumonia. *Radiology* 2020; **296**(2): E86-e96.
21. Commandeur F, Goeller M, Betancur J, et al. Deep Learning for Quantification of Epicardial and Thoracic Adipose Tissue From Non-Contrast CT. *IEEE Trans Med Imaging* 2018; **37**(8): 1835-46.
22. Gieraerts C, Dangis A, Janssen L, et al. Prognostic Value and Reproducibility of AI-assisted Analysis of Lung Involvement in COVID-19 on Low-Dose Submillisievert Chest CT: Sample Size Implications for Clinical Trials. *Radiology: Cardiothoracic Imaging* 2020; **2**(5): e200441.
23. Gietema HA, Müller NL, Nasute Fauerbach PV, et al. Quantifying the Extent of Emphysema: Factors Associated with Radiologists' Estimations and Quantitative Indices of Emphysema Severity Using the ECLIPSE Cohort. *Academic Radiology* 2011; **18**(6): 661-71.
24. Altman DG. Categorizing Continuous Variables. *Encyclopedia of Biostatistics*; 2005.
25. Grodecki K, Lin A, Razipour A, et al. Epicardial adipose tissue is associated with extent of pneumonia and adverse outcomes in patients with COVID-19. *Metabolism* 2021; **115**: 154436.
26. Sharifpour M, Rangaraju S, Liu M, et al. C-Reactive protein as a prognostic indicator in hospitalized patients with COVID-19. *PLOS ONE* 2020; **15**(11): e0242400.
27. Henry BM, Aggarwal G, Wong J, et al. Lactate dehydrogenase levels predict coronavirus disease 2019 (COVID-19) severity and mortality: A pooled analysis. *Am J Emerg Med* 2020; **38**(9): 1722-6.
28. Ashari MA, Zainal IA, Zaki FM. Strategies for radiology departments in handling the COVID-19 pandemic. *Diagn Interv Radiol* 2020; **26**(4): 296-300.

Figures

Figure 1. Study flowchart.

Figure 2. Chest CT of a woman with COVID-19 pneumonia who died 7 days later. Lobar severity score was 15 and segmental severity score was 24 (A-B). AI-assisted quantification revealed pneumonia burden of 36.7% (C) involving all five pulmonary lobes (D). Three-dimensional lung renderings depict distribution of disease consisting of both ground-glass opacities (blue) and consolidation (yellow) in (E) coronal and (F) axial planes.

Figure 3. Comparison of pneumonia burden (A), lobar severity score (B) and segmental severity score in patients with and without clinical deterioration or death. Box plots demonstrate the median, interquartile range 25th-75th, and minimum and maximum values.

Figure 4. Performance of different pneumonia scoring systems for prediction of clinical deterioration or death.

Figure 5. Performance of lobar severity scores estimated visually and with AI (A). Agreement chart showed excellent agreement (weighted kappa = 0.862) between visual and AI-assisted evaluation of pneumonia involvement in per-lobe analysis (B).

Table 1. Clinical and laboratory characteristics of patients on admission

	Clinical deterioration or death		P value
	Yes (n = 175)	No (n = 568)	
Clinical characteristics			
Age, years	72±14	62±17	<0.001
Male sex	100 (57.2)	311 (54.8)	0.578
Body mass index, kg/m ²	27.6±6.1	28.8±6.4	0.312
Hypertension	131 (74.9)	326 (57.4)	<0.001
Diabetes mellitus	54 (30.9)	135 (23.8)	0.052
Hyperlipidemia	47 (26.9)	116 (20.4)	0.066
Smoking status			0.250
Former smoker	31 (17.7)	71 (12.5)	
Current smoker	17 (9.7)	51 (9.0)	
History of lung disease	35 (20.0)	83 (14.6)	0.085
History of heart failure	34 (19.4)	74 (13.0)	<0.001
History of coronary artery disease	32 (23.4)	44 (7.7)	<0.001
Chronic kidney disease	34 (19.4)	59 (10.4)	<0.001
Immunodeficiency	41 (23.4)	76 (13.4)	<0.001
Symptoms			
Fever	77 (44.0)	282 (49.6)	0.191
Chills	3 (1.7)	31 (5.5)	0.038
Fatigue	50 (28.6)	219 (38.6)	0.016
Dyspnea	93 (53.1)	226 (39.8)	0.002
Dry cough	64 (36.6)	245 (43.1)	0.124
Sputum production	18 (10.3)	62 (10.9)	0.814
Hemoptysis	0 (0.0)	6 (1.1)	0.172
Sore throat	3 (1.7)	30 (5.3)	0.045
Loss of smell	4 (2.3)	52 (9.2)	0.003
Loss of taste	7 (4.0)	51 (9.0)	0.032
Muscle/joint pain	19 (10.9)	107 (18.8)	0.014
Headache	10 (5.7)	50 (8.8)	0.196
Nausea or vomiting	18 (10.3)	61 (10.7)	0.865
Diarrhea	22 (12.6)	51 (9.0)	0.163
Blood biomarkers			
Lymphocytes (%)	12.5 (7.7 – 18.7), 150	19.3 (13.2 – 27.0), 540	<0.001
Lactate dehydrogenase (U/L)	438±222, 138	256±133, 356	<0.001
C-reactive protein (mg/L)	128.5 (69.1 – 205.8), 136	40.9 (10.8 – 94.4), 434	<0.001
Ferritin (ng/mL)	819 (382 – 1286), 111	421 (221 – 773), 345	<0.001
Prothrombin time (s)	9.2 (8.7 – 10.4), 111	9.0 (8.5 – 9.8), 320	0.022
D-dimer (ng/mL)	2.24 (0.95 – 4.3), 104	0.9 (0.5 – 1.9), 310	<0.001
Troponin (pg /mL)	40.0 (19.0 – 78.5), 101	13.0 (6.0 – 30.0), 299	<0.001
Creatine phosphokinase (U/L)	112.0 (41.2 – 279.0), 90	67.5 (34.8 – 147.5), 274	<0.001

Data are n (%), median (IQR), or mean±SD, n if fewer patients had laboratory results available than the total study population.

Table 2. Classification and quantitative measures of lung lesions on chest CT in COVID-19 pneumonia

	Clinical deterioration or death		P value
	Yes (n = 175)	No (n = 568)	
Lung abnormality			
Only ground-glass opacities	18 (10.3)	132 (23.2)	<0.001
Only consolidation	2 (1.1)	1 (0.2)	0.077
Ground-glass opacities and consolidation	144 (82.3)	323 (56.9)	<0.001
Pleural effusion	42 (24.0)	64 (11.3)	<0.001
Emphysema	19 (10.9)	27 (4.8)	0.003
Fibrosis	13 (7.4)	19 (3.3)	0.020
None	11 (6.3)	112 (19.7)	<0.001
AI-assisted burden (%)			
Total	16.0 (4.5 – 39.3)	3.7 (0.3 – 10.3)	<0.001
Ground-glass opacities	12.1 (3.5 – 34.6)	3.3 (0.2 – 9.0)	<0.001
Consolidation	0.7 (0.1 – 3.7)	0.2 (0.0 – 0.4)	<0.001
Pleural effusion	0.0 (0.0 – 0.7)	0.0 (0.0 – 0.0)	<0.001
Lobar severity score*	10 (6 – 15)	6 (3 – 9)	<0.001
Segmental severity score*	18 (10 – 27)	9 (4 – 15)	<0.001

Data are n (%), median (IQR)

For detailed distribution of individual score components refer to Supplementary Table 1 and 2

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

PROOFS

Table 3. Risk score stratification across quartiles

	Quartile 1	Quartile 2	Quartile 3	Quartile 4	p Value
Range of AI-assisted pneumonia burden	(0-0.6%)	<0.6%-5.1%)	<5.1%-15.4%)	<15.4%-100%)	
Rate of clinical deterioration in AI-assisted pneumonia burden quartiles	10.2% (19/185)	14.5% (27/186)	20.3% (38/187)	49.2% (91/185)	<0.001
Range of lobar severity score	0-3	4-6	7-10	11-25	
Rate of clinical deterioration in lobar severity score quartiles	12.5% (24/191)	13.7% (25/183)	21.5% (42/195)	48.2% (84/174)	<0.001
Range of segmental severity score	0-5	6-11	12-17	18-40	
Rate of clinical deterioration in segmental severity score quartiles	11.8% (23/194)	14.5% (30/206)	19.8% (32/161)	49.5% (90/182)	<0.001

BJR UNCORR

PROOFS

1
2
3
4
5
6
7
8
9
10
11

Table 2. Reclassification table using the artificial intelligence for lobar severity scores.

12
13

Visual	Artificial Intelligence							Reclassification		NRI	P Value for NRI
	0%	<0%;5%)	<5%;25%)	<25%;50%)	<50%;75%)	>75%	Total	Up	Down		
No clinical deterioration								5.1%	4.8%	-0.1%	0.994
0%	768	131	0	0	0	0	899				
<0%;5%)	5	884	3	0	0	0	833				
<5%;25%)	24	46	642	11	0	0	723				
<25%;50%)	7	12	22	208	1	0	250				
<50%;75%)	1	2	2	7	45	0	57				
>75%	1	0	3	2	2	11	19				
Total	806	1075	672	228	48	11	2840				
Clinical deterioration								6.2%	5.8%		
0%	113	25	2	1	1	0	142				
<0%;5%)	2	179	5	1	0	0	187				
<5%;25%)	0	15	192	6	1	0	214				
<25%;50%)	0	1	9	139	9	0	158				
<50%;75%)	0	0	2	19	109	3	133				
>75%	0	0	0	0	3	38	41				
Total	115	220	210	166	123	41	875				

14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37

NRI = net reclassification improvement;

$$NRI = [P(Up|Positive) - P(Down|Positive)] - [P(Up|Negative) - P(Down|Negative)]$$

Patients were reclassified by artificial intelligence and were compared to visual scoring.

38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Bu

ED PROOFS

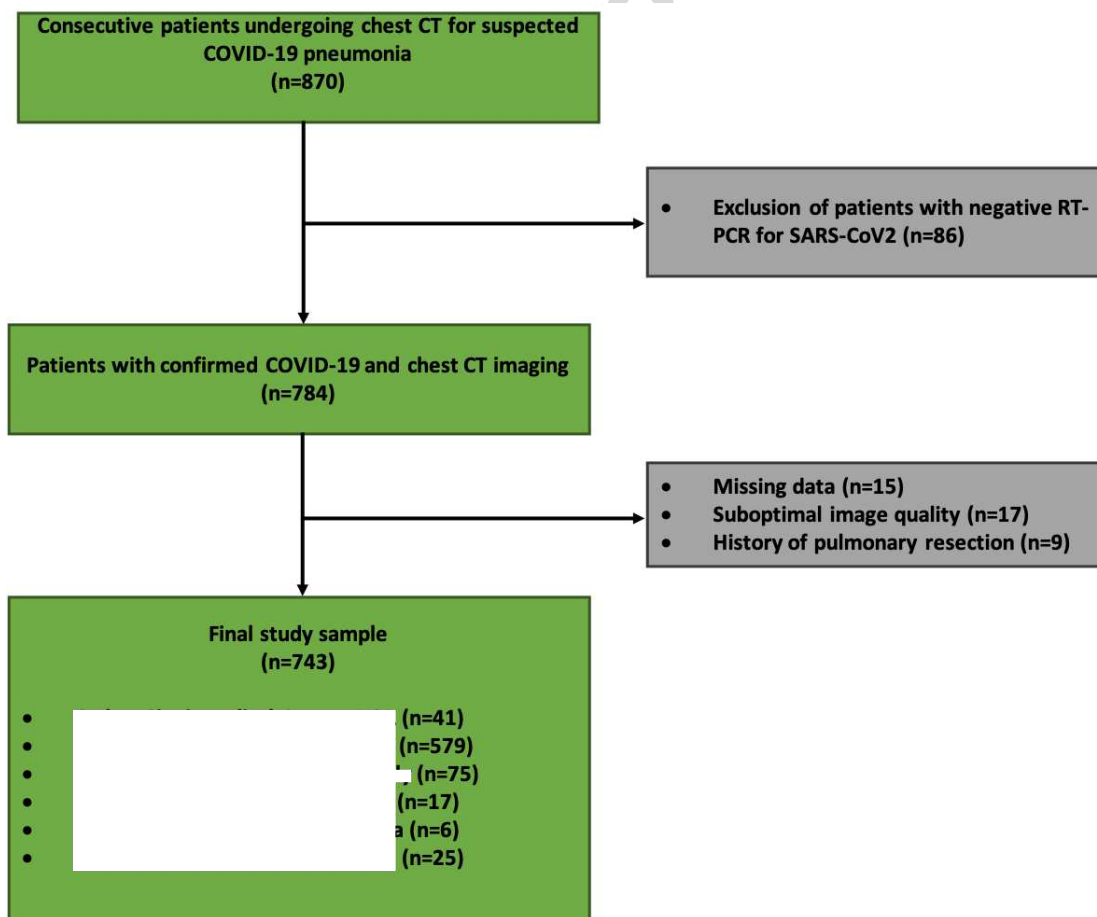
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Table 5. Correlation matrix between investigated variables.

	Al-assisted pneumonia burden	Lobar severity score	Segmental severity score	Lymphocytes	Lactate dehydrogenase	C-reactive protein	Ferritin	Prothrombin time	D-dimer	Troponin
Lobar severity score	0.933**									
Segmental severity score	0.910**	0.972**								
Lymphocytes	-0.299**	-0.306**	-0.302**							
Lactate dehydrogenase	0.567**	0.542**	0.516**	-0.268**						
C-reactive protein	0.507**	0.486**	0.467**	-0.426**	0.460**					
Ferritin	0.399**	0.394**	0.409**	-0.231**	0.454**	0.417**				
Prothrombin time	-0.001	-0.003	0.26	-0.020	-0.011	-0.193**	0.078			
D-dimer	0.205**	0.196**	0.197**	-0.332**	0.264**	0.292**	0.204**	0.125*		
Troponin	-0.025	-0.048	-0.053	-0.282**	0.128**	0.244**	0.055	0.076	0.311**	
Creatine phosphokinase	0.191**	0.172**	0.143**	-0.093*	0.340**	0.163**	0.139**	-0.074	0.151**	0.109*

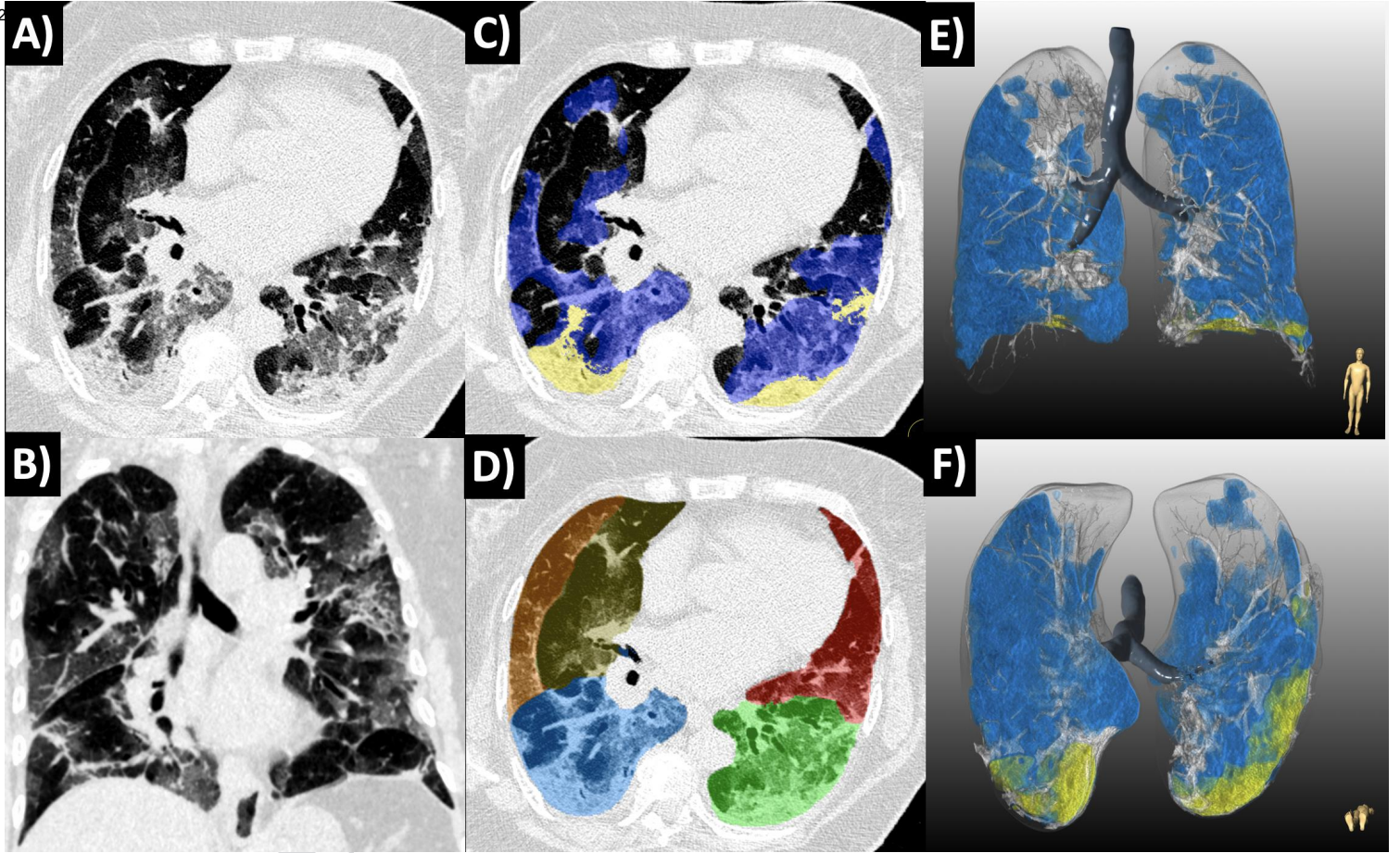
*Correlation is significant at the 0.01 level.
**Correlation is significant at the 0.05 level.

Figure 1



PROOFS

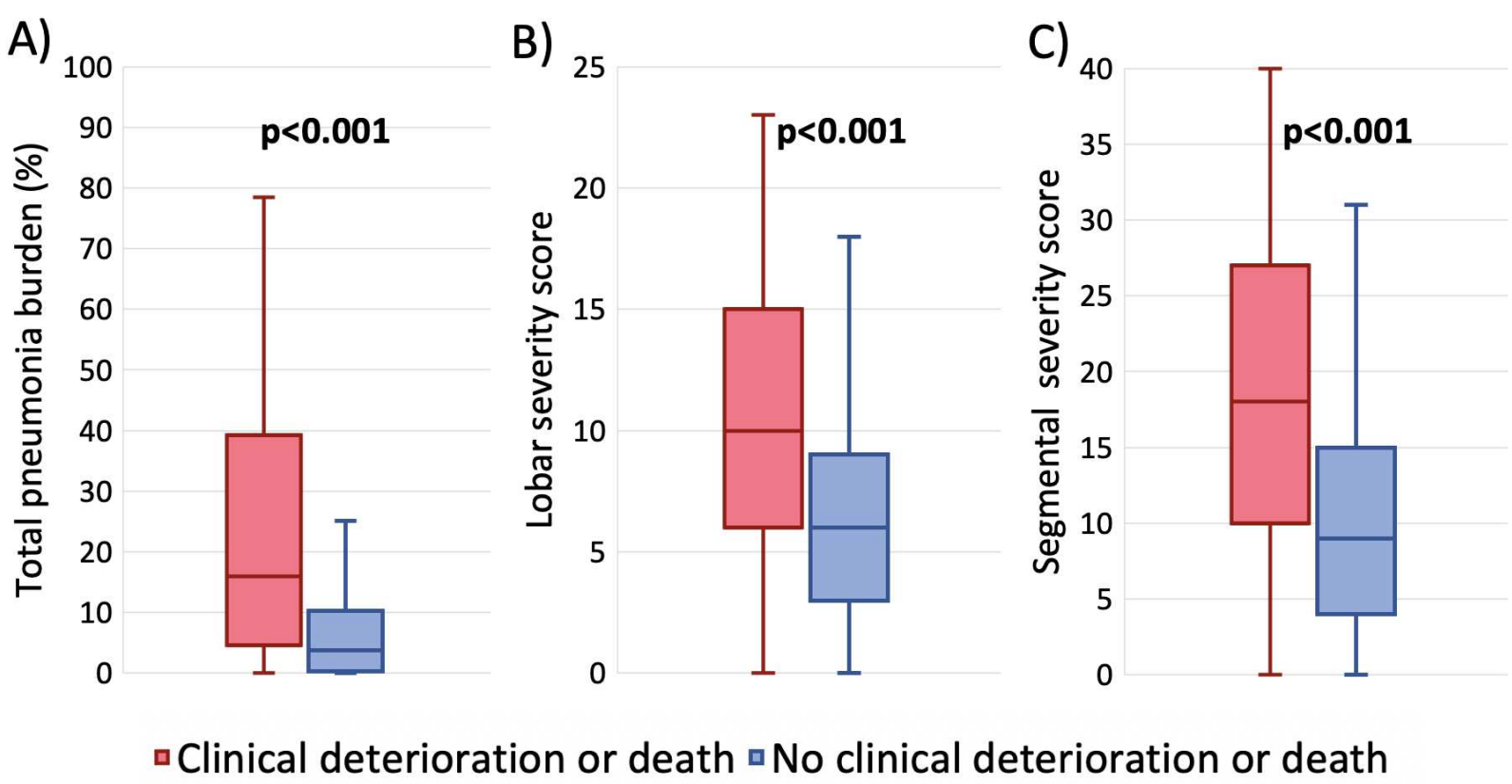
Figure 2



BJR U

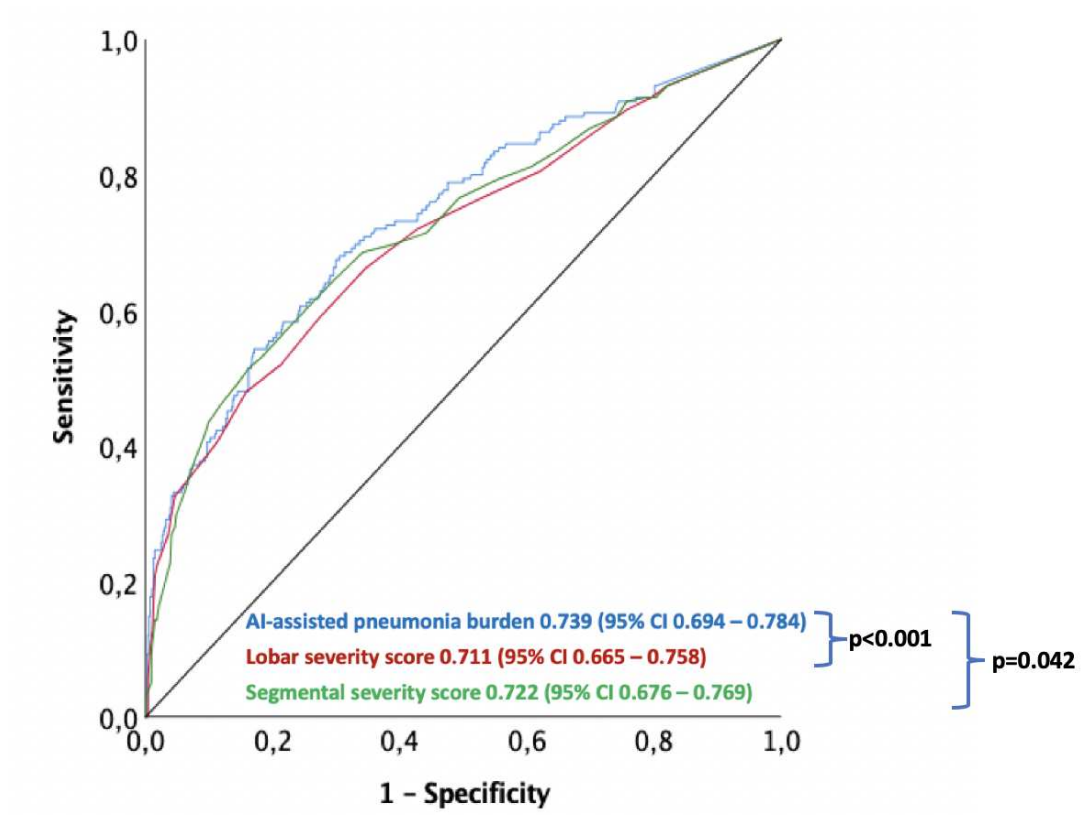
PROOFS

Figure 3



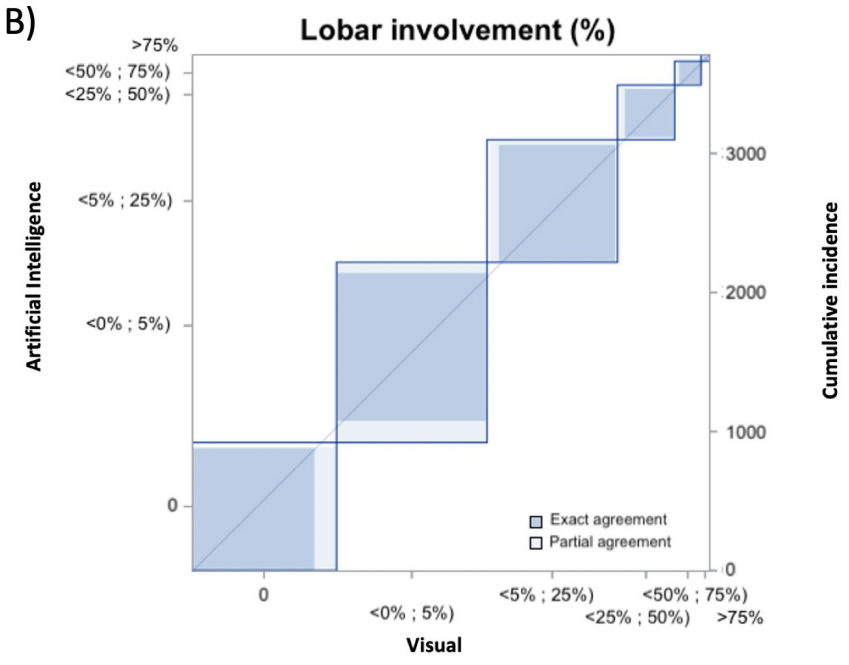
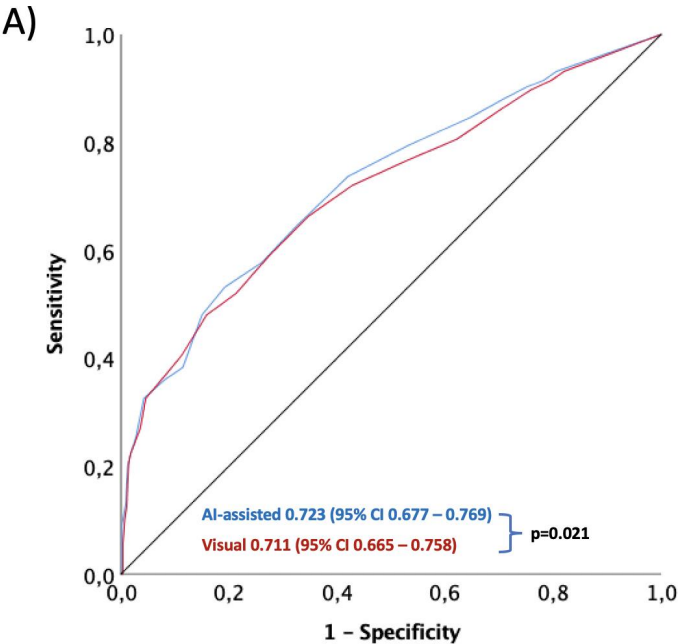
BJR U

Figure 4



PROOFS

Figure 5



BJR U

Table 1. Clinical and laboratory characteristics of patients on admission

	Clinical deterioration or death		P value
	Yes (n = 175)	No (n = 568)	
Clinical characteristics			
Age, years	72±14	62±17	<0.001
Male sex	100	311	0.578
Body mass index, kg/m ²	27.6±6.1	28.8±6.4	0.312
Hypertension	131 (74.9)	326 (57.4)	<0.001
Diabetes mellitus	54 (30.9)	135 (23.8)	0.052
Hyperlipidemia	47 (26.9)	116 (20.4)	0.066
Smoking status			0.250
Former smoker	31 (17.7)	71 (12.5)	
Current smoker	17 (9.7)	51 (9.0)	
History of lung disease	35 (20.0)	83 (14.6)	0.085
History of heart failure	34 (19.4)	74 (13.0)	<0.001
History of coronary artery disease	32 (23.4)	44 (7.7)	<0.001
Chronic kidney disease	34 (19.4)	59 (10.4)	<0.001
Immunodeficiency	41 (23.4)	76 (13.4)	<0.001
Symptoms			
Fever	77 (44.0)	282 (49.6)	0.191
Chills	3 (1.7)	31 (5.5)	0.038
Fatigue	50 (28.6)	219 (38.6)	0.016
Dyspnea	93 (53.1)	226 (39.8)	0.002
Dry cough	64 (36.6)	245 (43.1)	0.124
Sputum production	18 (10.3)	62 (10.9)	0.814
Hemoptysis	0 (0.0)	6 (1.1)	0.172
Sore throat	3 (1.7)	30 (5.3)	0.045
Loss of smell	4 (2.3)	52 (9.2)	0.003
Loss of taste	7 (4.0)	51 (9.0)	0.032
Muscle/joint pain	19 (10.9)	107 (18.8)	0.014
Headache	10 (5.7)	50 (8.8)	0.196
Nausea or vomiting	18 (10.3)	61 (10.7)	0.865
Diarrhea	22 (12.6)	51 (9.0)	0.163
Blood biomarkers			
Lymphocytes (%)	12.5 (7.7 – 18.7), 150	19.3 (13.2 – 27.0), 540	<0.001
Lactate dehydrogenase (U/L)	438±222, 138	256±133, 356	<0.001
C-reactive protein (mg/L)	128.5 (69.1 – 205.8), 136	40.9 (10.8 – 94.4), 434	<0.001
Ferritin (ng/mL)	819 (382 – 1286), 111	421 (221 – 773), 345	<0.001
Prothrombin time (s)	9.2 (8.7 – 10.4), 111	9.0 (8.5 – 9.8), 320	0.022
D-dimer (ng/mL)	2.24 (0.95 – 4.3), 104	0.9 (0.5 – 1.9), 310	<0.001
Troponin (pg /mL)	40.0 (19.0 – 78.5), 101	13.0 (6.0 – 30.0), 299	<0.001
Creatine phosphokinase (U/L)	112.0 (41.2 – 279.0), 90	67.5 (34.8 – 147.5), 274	<0.001

Data are n (%), median (IQR), or mean±SD, n if fewer patients had laboratory results available than the total study population.

Table 2. Classification and quantitative measures of lung lesions on chest CT in COVID-19 pneumonia

	Clinical deterioration or death		<i>P</i> value
	Yes (n = 175)	No (n = 568)	
Lung abnormality			
Only ground-glass opacities	18 (10.3)	132 (23.2)	<0.001
Only consolidation	2 (1.1)	1 (0.2)	0.077
Ground-glass opacities and consolidation	144 (82.3)	323 (56.9)	<0.001
Pleural effusion	42 (24.0)	64 (11.3)	<0.001
Emphysema	19 (10.9)	27 (4.8)	0.003
Fibrosis	13 (7.4)	19 (3.3)	0.020
None	11 (6.3)	112 (19.7)	<0.001
AI-assisted burden (%)			
Total	16.0 (4.5 – 39.3)	3.7 (0.3 – 10.3)	<0.001
Ground-glass opacities	12.1 (3.5 – 34.6)	3.3 (0.2 – 9.0)	<0.001
Consolidation	0.7 (0.1 – 3.7)	0.2 (0.0 – 0.4)	<0.001
Pleural effusion	0.0 (0.0 – 0.7)	0.0 (0.0 – 0.0)	<0.001
Lobar severity score*	10 (6 – 15)	6 (3 – 9)	<0.001
Segmental severity score*	9 (4 – 15)	18 (10 – 27)	<0.001

Data are n (%), median (IQR)

*For detailed distribution of individual score components refer to Supplementary Table 1 and 2

Table 3

Table 3. Risk score stratification across quartiles

	Quartile 1	Quartile 2	Quartile 3	Quartile 4	p Value
Range of AI-assisted pneumonia burden	(0-0.6%)	<0.6%-5.1%)	<5.1%-15.4%)	<15.4%-100%)	
Rate of clinical deterioration in AI-assisted pneumonia burden quartiles	10.2% (19/185)	14.5% (27/186)	20.3% (38/187)	49.2% (91/185)	<0.001
Range of lobar severity score	0-3	4-6	7-10	11-25	
Rate of clinical deterioration in lobar severity score quartiles	12.5% (24/191)	13.7% (25/183)	21.5% (42/195)	48.2% (84/174)	<0.001
Range of segmental severity score	0-5	6-11	12-17	18-40	
Rate of clinical deterioration in segmental severity score quartiles	11.8% (23/194)	14.5% (30/206)	19.8% (32/161)	49.5% (90/182)	<0.001

Table 4

Table 4. Reclassification table using the artificial intelligence for lobar severity scores.

Visual	Artificial Intelligence						Total	Reclassification		NRI	P Value for NRI
	0%	<0%;5%)	<5%;25%)	<25%;50%)	<50%;75%)	>75%		Up	Down		
No clinical deterioration								5.1%	4.8%	-0.1%	0.994
0%	768	131	0	0	0	0	899				
<0%;5%)	5	884	3	0	0	0	833				
<5%;25%)	24	46	642	11	0	0	723				
<25%;50%)	7	12	22	208	1	0	250				
<50%;75%)	1	2	2	7	45	0	57				
>75%	1	0	3	2	2	11	19				
Total	806	1075	672	228	48	11	2840				
Clinical deterioration								6.2%	5.8%		
0%	113	25	2	1	1	0	142				
<0%;5%)	2	179	5	1	0	0	187				
<5%;25%)	0	15	192	6	1	0	214				
<25%;50%)	0	1	9	139	9	0	158				
<50%;75%)	0	0	2	19	109	3	133				
>75%	0	0	0	0	3	38	41				
Total	115	220	210	166	123	41	875				

NRI = net reclassification improvement;

$$\text{NRI} = [P(\text{Up}|\text{Positive}) - P(\text{Down}|\text{Positive})] - [P(\text{Up}|\text{Negative}) - P(\text{Down}|\text{Negative})]$$

Patients were reclassified by artificial intelligence and were compared to visual scoring.

Table 5

Table 5. Correlation matrix between investigated variables.

	AI-assisted pneumonia burden	Lobar severity score	Segmental severity score	Lymphocytes	Lactate dehydrogenase	C-reactive protein	Ferritin	Prothrombin time	D-dimer	Troponin
Lobar severity score	0.933**									
Segmental severity score	0.910**	0.972**								
Lymphocytes	-0.299**	-0.306**	-0.302**							
Lactate dehydrogenase	0.567**	0.542**	0.516**	-0.268**						
C-reactive protein	0.507**	0.486**	0.467**	-0.426**	0.460**					
Ferritin	0.399**	0.394**	0.409**	-0.231**	0.454**	0.417**				
Prothrombin time	-0.001	-0.003	0.26	-0.020	-0.011	-0.193**	0.078			
D-dimer	0.205**	0.196**	0.197**	-0.332**	0.264**	0.292**	0.204**	0.125*		
Troponin	-0.025	-0.048	-0.053	-0.282**	0.128**	0.244**	0.055	0.076	0.311**	
Creatine phosphokinase	0.191**	0.172**	0.143**	-0.093*	0.340**	0.163**	0.139**	-0.074	0.151**	0.109*

*Correlation is significant at the 0.01 level.

**Correlation is significant at the 0.05 level.