

## ARTICLE OPEN



# An artificial intelligence system for predicting the deterioration of COVID-19 patients in the emergency department

Farah E. Shamout<sup>1,9</sup>, Yiqiu Shen<sup>2,9</sup>, Nan Wu<sup>2,9</sup>, Aakash Kaku<sup>2,9</sup>, Jungkyu Park<sup>3,4,9</sup>, Taro Makino<sup>5,2,3,9</sup>, Stanisław Jastrzębski<sup>2,3,5</sup>, Jan Witowski<sup>6,3,5</sup>, Duo Wang<sup>6</sup>, Ben Zhang<sup>6</sup>, Siddhant Dogra<sup>6,3</sup>, Meng Cao<sup>7</sup>, Narges Razavian<sup>2,3,6</sup>, David Kudlowitz<sup>7</sup>, Lea Azour<sup>6,3</sup>, William Moore<sup>3</sup>, Yvonne W. Lui<sup>6,3,5</sup>, Yindalon Aphinyanaphongs<sup>6,3</sup>, Carlos Fernandez-Granda<sup>6,2,8</sup> and Krzysztof J. Geras<sup>6,2,3,5</sup>✉

During the coronavirus disease 2019 (COVID-19) pandemic, rapid and accurate triage of patients at the emergency department is critical to inform decision-making. We propose a data-driven approach for automatic prediction of deterioration risk using a deep neural network that learns from chest X-ray images and a gradient boosting model that learns from routine clinical variables. Our AI prognosis system, trained using data from 3661 patients, achieves an area under the receiver operating characteristic curve (AUC) of 0.786 (95% CI: 0.745–0.830) when predicting deterioration within 96 hours. The deep neural network extracts informative areas of chest X-ray images to assist clinicians in interpreting the predictions and performs comparably to two radiologists in a reader study. In order to verify performance in a real clinical setting, we silently deployed a preliminary version of the deep neural network at New York University Langone Health during the first wave of the pandemic, which produced accurate predictions in real-time. In summary, our findings demonstrate the potential of the proposed system for assisting front-line physicians in the triage of COVID-19 patients.

*npj Digital Medicine* (2021)4:80; <https://doi.org/10.1038/s41746-021-00453-0>

## INTRODUCTION

In recent months, there has been a surge in patients presenting to the emergency department (ED) with respiratory illnesses associated with the coronavirus disease 2019 (COVID-19)<sup>1,2</sup>. Evaluating the risk of deterioration of these patients to perform triage is crucial for clinical decision-making and resource allocation<sup>3</sup>. While ED triage is difficult under normal circumstances<sup>4,5</sup>, during a pandemic, strained hospital resources increase the challenge<sup>2,6</sup>. This is compounded by our incomplete understanding of COVID-19. Data-driven risk evaluation based on artificial intelligence (AI) could, therefore, play an important role in streamlining ED triage.

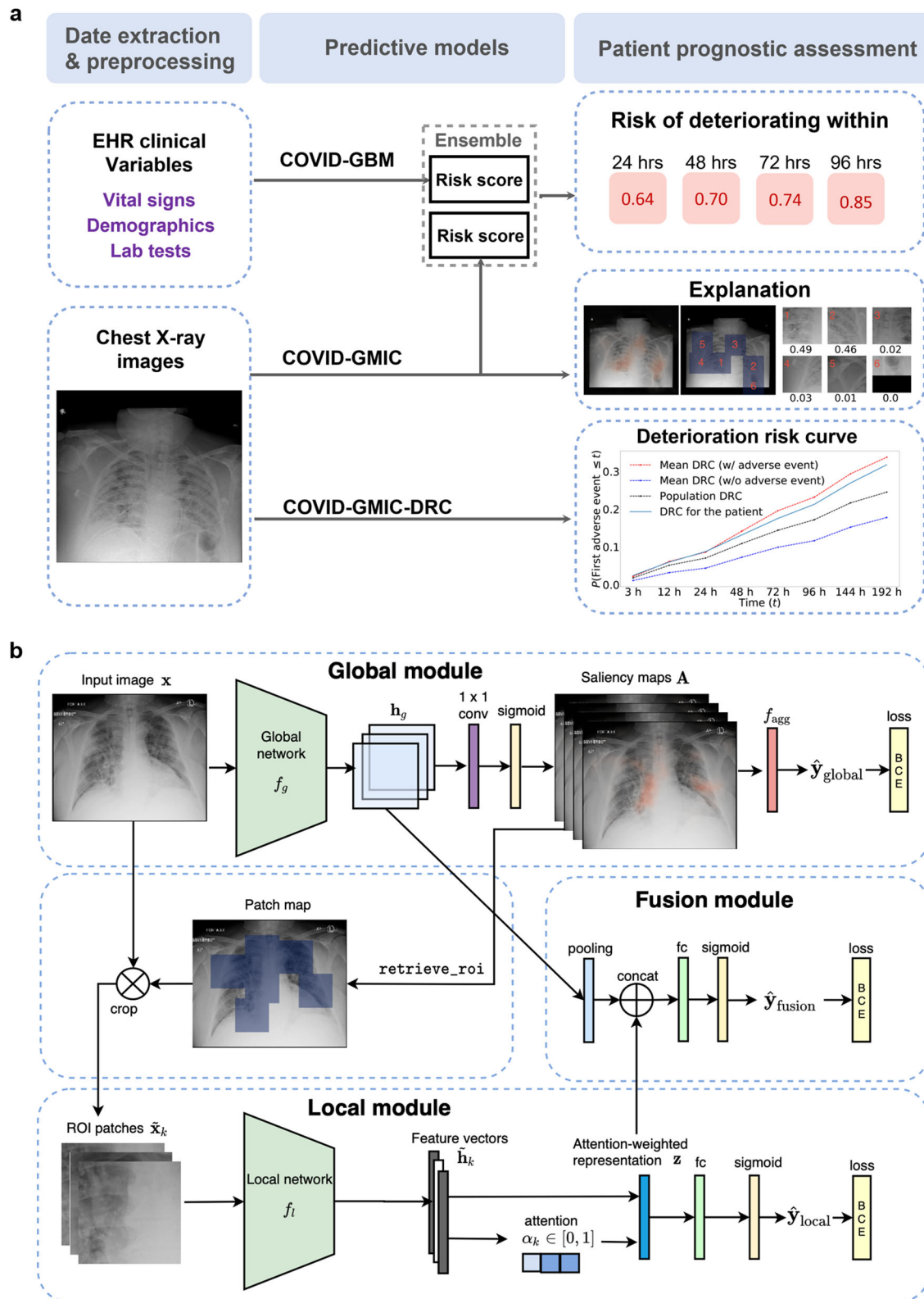
As the primary complication of COVID-19 is pulmonary disease, such as pneumonia<sup>7</sup>, chest X-ray imaging is a first-line triage tool for COVID-19 patients<sup>8</sup>. Although other imaging modalities, such as computed tomography (CT), provide higher resolution, chest X-ray imaging is less costly, inflicts a lower radiation dose, and is easier to obtain without incurring the risk of contaminating imaging equipment and disrupting radiologic services<sup>9</sup>. In addition, abnormalities in the chest X-ray images of COVID-19 patients have been found to mirror abnormalities in CT scans<sup>10</sup>. Although the knowledge of the disease is rapidly evolving, the understanding of the correlation between pulmonary parenchymal patterns visible in the chest X-ray images and clinical deterioration remains limited. This motivates the use of machine learning approaches for risk stratification using chest X-ray imaging, which may be able to learn such correlations automatically from data.

The majority of related previous work using imaging data of COVID-19 patients focus more on diagnosis than prognosis<sup>11–18</sup>.

Prognostic models used for predicting mortality, morbidity and other outcomes related to the disease course have a number of potential real-life applications, such as: consistently defining and triaging sick patients, alerting bed management teams on expected demands, providing situational awareness across teams of individual patients, and more general resource allocation<sup>11</sup>. Prior methodology for prognosis of COVID-19 patients via machine learning mainly use routinely collected clinical variables<sup>2,19</sup> such as vital signs and laboratory tests, which have long been established as strong predictors of deterioration<sup>20,21</sup>. Some studies have proposed scoring systems for chest X-ray images to assess the severity and progression of lung involvement using deep learning<sup>22</sup>, or more commonly, through manual clinical evaluation<sup>7,23,24</sup>. In general, the role of deep learning for the prognosis of COVID-19 patients using chest X-ray imaging has not yet been fully established. Using both the images and the clinical variables in a single AI system also has not been studied before. We show that they both contain complimentary information, which opens a new perspective on building prognostic AI systems for COVID-19.

In this retrospective study, we develop an AI system that performs an automatic evaluation of deterioration risk, based on chest X-ray imaging, combined with other routinely collected non-imaging clinical variables. An overview of the system is shown in Fig. 1a. The goal is to provide support for critical clinical decision-making involving patients arriving at the ED in need of immediate care<sup>2,25</sup>, based on the need for efficient patient triage. The system is based on chest X-ray imaging, while also incorporating other routinely collected non-imaging clinical variables that are known to be strong predictors of deterioration.

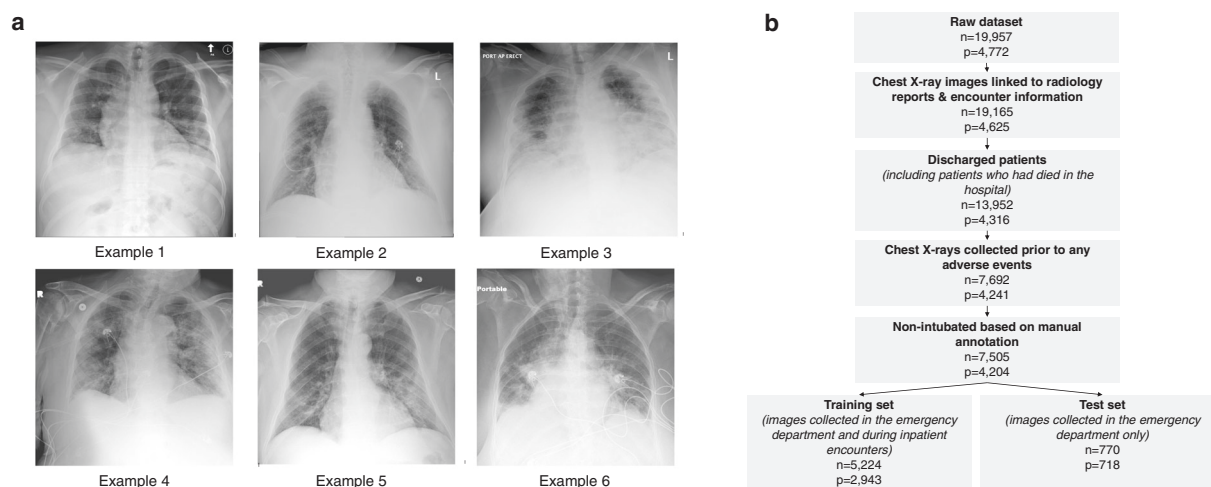
<sup>1</sup>Engineering Division, NYU Abu Dhabi, Abu Dhabi, UAE. <sup>2</sup>Center for Data Science, New York University, New York, NY, USA. <sup>3</sup>Department of Radiology, NYU Langone Health, New York, NY, USA. <sup>4</sup>Vilcek Institute of Graduate Biomedical Sciences, NYU Grossman School of Medicine, New York, NY, USA. <sup>5</sup>Center for Advanced Imaging Innovation and Research, NYU Langone Health, New York, NY, USA. <sup>6</sup>Department of Population Health, NYU Langone Health, New York, NY, USA. <sup>7</sup>Department of Medicine, NYU Langone Health, New York, NY, USA. <sup>8</sup>Department of Mathematics, Courant Institute, New York University, New York, NY, USA. <sup>9</sup>These authors contributed equally: Farah E. Shamout, Yiqiu Shen, Nan Wu, Aakash Kaku, Jungkyu Park, Taro Makino ✉email: [k.j.geras@nyu.edu](mailto:k.j.geras@nyu.edu)



Our AI system uses deep convolutional neural networks to perform risk evaluation from chest X-ray images. In particular, we designed our imaging-based classifier based on the Globally-Aware Multiple Instance Classifier (GMIC)<sup>26,27</sup>, denoted as COVID-GMIC, aiming for accurate performance and interpretability (see Fig. 1b). The system also learns from routinely collected clinical

variables using a gradient boosting model (GBM)<sup>28</sup>, denoted as COVID-GBM. Both models were trained using a dataset of 3661 patients admitted to NYU Langone Health between March 3, 2020, and May 13, 2020. To learn from both modalities, we combined the output predictions of COVID-GMIC and COVID-GBM to predict each patient's overall risk of deterioration over different time

**Fig. 1 Overview of the AI system and the architecture of its deep learning component.** **a** Overview of the AI system that assesses the patient's risk of deterioration every time a chest X-ray image is collected in the ED. We design two different models to process the chest X-ray images, both based on the GMIC neural network architecture<sup>26,27</sup>. The first model, COVID-GMIC, predicts the overall risk of deterioration within 24, 48, 72, and 96 h, and computes saliency maps that highlight the regions of the image that most informed its predictions. The predictions of COVID-GMIC are combined with predictions of a gradient boosting model<sup>28</sup> that learns from routinely collected clinical variables, referred to as COVID-GBM. The second model, COVID-GMIC-DRC, predicts how the patient's risk of deterioration evolves over time in the form of deterioration risk curves. **b** Architecture of COVID-GMIC. First, COVID-GMIC utilizes the global network to generate four saliency maps that highlight the regions on the X-ray image that are predictive of the onset of adverse events within 24, 48, 72, and 96 h, respectively. COVID-GMIC then applies a local network to extract fine-grained visual details from these regions. Finally, it employs a fusion module that aggregates information from both the global context and local details to make a holistic diagnosis.



**Fig. 2 Illustrations of the dataset and the dataset flowchart.** **a** Examples of chest X-ray images in our dataset. Example 1: Patient was discharged and experienced no adverse events (44 years old male). Example 2: Patient was transferred to the ICU after 95 h (71 years old male). Example 3: Patient was intubated after 72 h (66 years old male). Example 4: Patient was transferred to the ICU after 48 h (99 years old female). Example 5: Patient was intubated after 24 h (74 years old male). Example 6: Patient was transferred to the ICU in 30 min (73 years old female). It is important to note that the extent of parenchymal disease does not necessarily have a direct correlation with deterioration time. For example, Example 5 has less severe parenchymal findings than Examples 3 and 4, but deteriorated faster. **b** Flowchart showing how the inclusion and exclusion criteria were applied to obtain the final training and test sets, where  $n$  represents the number of chest X-ray exams, and  $p$  represents the number of unique patients. Specifically, we excluded 783 exams that were not linked to any radiology report, nine exams that had missing encounter information, and 5213 exams from patients who were still hospitalized by May 13, 2020. To ensure that our system predicts deterioration prior to its occurrence, we excluded 6260 exams that were collected after an adverse event and 187 exams of already intubated patients. The final dataset consisted of 7502 chest X-ray exams corresponding to 4204 unique patients. We split the dataset at the patient level such that exams from the same patient exclusively appear either in the training or test set. In the training set, we included exams that were collected both in the ED and during inpatient encounters. Since the intended clinical use of our model is in the ED, the test set only includes exams collected in the ED and hence we excluded 543 patients who did not have exams collected in the ED. This resulted in 5224 exams (5617 images) in the training set and 770 exams (832 images) in the test set. We included both frontal and lateral images, however there were less than 50 lateral images in the entire dataset.

horizons, ranging from 24 to 96 h. In addition, the system includes a model that predicts how the risk of deterioration is expected to evolve over time by computing deterioration risk curves (DRC), in the spirit of survival analysis<sup>29</sup>, denoted as COVID-GMIC-DRC.

Our system is able to accurately predict the deterioration risk on a test set of new patients. It achieves an area under the receiver operating characteristic curve (AUC) of 0.786 (95% CI: 0.745–0.830), and an area under the precision-recall curve (PR AUC) of 0.517 (95% CI: 0.429–0.600) for prediction of deterioration within 96 h. Additionally, its estimated probability of the temporal risk evolution discriminates effectively between patients, and is well-calibrated. The imaging-based model achieves a comparable AUC to two experienced chest radiologists in a reader study, highlighting the potential of our data-driven approach. In order to verify our system's performance in a real clinical setting, we silently deployed a preliminary version of it at NYU Langone Health during the first wave of the pandemic, demonstrating that it can produce accurate predictions in real-time. Overall, these results strongly suggest that our system is a viable and valuable tool to inform triage of COVID-19 patients. For reproducibility,

we published our code and the trained models at [https://github.com/nyukat/COVID-19\\_prognosis](https://github.com/nyukat/COVID-19_prognosis).

## RESULTS

### Dataset

Our AI system was developed and evaluated using a dataset of 19,957 chest X-ray exams collected from 4,722 patients at NYU Langone Health between March 3, 2020 and May 13, 2020. The dataset consists of chest X-ray images collected from patients who tested positive for COVID-19 using the polymerase chain reaction (PCR) test, along with the clinical variables recorded closest to the time of image acquisition (e.g., vital signs, laboratory test results, and patient characteristics). Figure 2a shows examples of chest X-ray images collected from different patients. We applied inclusion and exclusion criteria that were defined in collaboration with clinical experts, as shown in Fig. 2b. The training set consisting of 2943 patients and 5617 chest X-ray images was used for model development and hyperparameter tuning using Monte Carlo cross-validation, where 20% of the training set was used for

model validation for each hyperparameter configuration. The test set consisting of 718 patients and 832 images was used to report the final results and was not used during training. The training and the test sets were disjoint, with no patient overlap. Table 1 summarizes the overall demographics and characteristics of the patient cohort in the training and test sets, including distributions of the included clinical variables. The raw laboratory test variables were further processed to extract the minimum and maximum statistics.

### Ground-truth labels

We define deterioration, the target to be predicted by our models, as the composite outcome of one of three adverse events: intubation, admission to the intensive care unit (ICU), and in-hospital mortality. If multiple adverse events occurred, we only consider the time of the first event.

### Evaluation metrics

Throughout the paper we used AUC (area under the receiver operating characteristic curve) and PR AUC (area under the precision-recall curve), which offer a complimentary view on the performance of our models. We additionally computed positive predictive values (PPV) and negative predictive values (NPV) for each of 24, 48, 72, 96 h tasks. We dichotomized the probabilistic predictions to reflect the class distribution in the training set. These metrics integrate the performance of the evaluated models over all possible thresholds for predictions to be considered positive. As there are no available guidelines on how to select the threshold, we prefer these metrics to metrics that are computed for a fixed threshold (i.e., F1 score or classification accuracy). We also computed 95% confidence intervals estimated by 1000 iterations of the bootstrap method<sup>30</sup>.

### Model performance

Table 2 summarizes the performance of all the models in terms of the AUC and PR AUC for the prediction of deterioration within 24, 48, 72, and 96 h from the time of the chest X-ray exam. The receiver operating characteristic curves and precision-recall curves can be found in Supplementary Fig. 1. The clinical variables only model (COVID-GBM) achieves a better performance than a logistic regression baseline across all time windows. We trained a logistic regression model utilizing only clinical variables achieved 0.698, 0.699, 0.712, and 0.728 AUC and 0.214, 0.266, 0.339, and 0.436 PR AUC across the 24, 48, 72, and 96 h windows, respectively. It is not possible to directly compare the performance of COVID-GBM and COVID-GMIC in the current setting since they model different training and test sets, although they rely on the same patient-level data splits. We reported the NPVs and PPVs of COVID-GMIC, COVID-GBM, and the ensemble of the two in Supplementary Table 1. We also show examples that were incorrectly classified (false positives and false negatives) in Supplementary Fig. 2.

However, the performance of the ensemble model consisting of COVID-GMIC and COVID-GBM achieves an improved AUC and PR AUC across all time windows compared to the COVID-GMIC baseline. This highlights the complementary role of chest X-ray images and routine clinical variables in predicting deterioration. The weighting of the predictions of COVID-GMIC and COVID-GBM was optimized on the validation set, as shown in Supplementary Fig. 3b. The consistent advantage of the ensemble model in our results is especially encouraging. Investigating more complex strategies for fusion of information from these two modalities could further improve the results and this will be a subject of our future research. Sample learning curves of COVID-GMIC are shown in Supplementary Fig. 4 for reference.

To illustrate the interpretability of COVID-GMIC, we show in Fig. 3 the saliency maps for all time windows (24, 48, 72, and 96 h) computed for four examples from the test set. Across all four

**Table 1.** Description of the characteristics of the patient cohort included in the training and test sets and the mean and interquartile range statistics of the raw vital signs and laboratory test results used for COVID-GBM.

Characteristic	Training set	Test set
Patients, <i>n</i>	2943	718
Admissions, <i>n</i>	3175	764
Sex (females), <i>n</i> (%)	1206 (41.0)	305 (42.5)
Age (years)*, mean (SD)	61.9 (17.6)	64.7 (17.4)
BMI (kg/m <sup>2</sup> ), mean (SD)	29.6 (6.7)	29.4 (7.3)
Weight (kg), mean (SD)	83.1 (22.2)	82.2 (23.1)
Survived	2405	559
Adverse events, <i>n</i>	1311	369
Intubation, <i>n</i> (%)	386 (29.4)	97 (26.3)
ICU admission, <i>n</i> (%)	387 (29.5)	113 (30.6)
Mortality, <i>n</i> (%)	538 (41.0)	159 (43.1)
Chest X-ray exams, <i>n</i>	5224	770
Composite outcome within 24 h, <i>n</i> (%)	349 (6.7)	74 (9.6)
Composite outcome within 48 h, <i>n</i> (%)	553 (10.6)	101 (13.1)
Composite outcome within 72 h, <i>n</i> (%)	735 (14.1)	130 (16.9)
Composite outcome within 96 h, <i>n</i> (%)	876 (16.8)	156 (20.3)
Total number of images, <i>n</i>	5617	832
Vital signs feature sets, <i>n</i> units	10,640	2776
Heart rate, beats per minute	93.7 (25.0)	93.5 (27.0)
Respiratory rate*, breaths per minute	22.4 (7.0)	23.4 (7.0)
Temperature, °F	99.4 (1.9)	99.4 (1.9)
Systolic blood pressure, mmHg	130.7 (30.0)	129.8 (29.3)
Diastolic blood pressure, mmHg	75.9 (17.0)	76.0 (18.0)
Oxygen saturation*, %	94.1 (4.0)	93.8 (5.0)
Provision of supplemental oxygen*, <i>n</i> (%)	3970 (37.3)	1166 (42.0)
Raw laboratory test results, <i>units</i>		
Albumin, g/dL	3.5 (0.9)	3.5 (0.9)
Alanine transaminase, U/L	49.8 (32.0)	52.2 (36.0)
Aspartate aminotransferase, U/L	67.3 (37.0)	69.7 (43.0)
Total bilirubin, mg/dL	0.7 (0.4)	0.7 (0.4)
Blood urea nitrogen, mg/dL	25.9 (17.0)	26.4 (18.0)
Calcium, mg/dL	8.7 (0.8)	8.7 (0.8)
Chloride, mEq/L	101.1 (7.0)	101.6 (7.0)
Creatinine, mg/dL	1.6 (0.7)	1.6 (0.7)
D-dimer, ng/mL	1321.6 (535.5)	1146.3 (618.5)
Eosinophils, %	0.4 (0.0)	0.4 (0.0)
Eosinophils, <i>n</i>	0.03 (0.00)	0.03 (0.00)
Hematocrit, %	38.9 (7.3)	38.9 (7.5)
Lactate dehydrogenase, U/L	412.8 (207.0)	404.0 (213.0)
Lymphocytes, %	14.1 (10.0)	14.9 (11.0)
Lymphocytes, <i>n</i>	1.0 (0.7)	1.0 (0.7)
Platelet volume, fL	10.6 (1.4)	10.6 (1.4)
Neutrophils, <i>n</i>	6.4 (4.0)	6.3 (3.8)
Neutrophils, %	76.6 (14.0)	75.9 (13.0)
Platelet, <i>n</i>	226.1 (114.0)	223.7 (103.0)
Potassium, mmol/L	4.2 (0.8)	4.2 (0.8)
Procalcitonin, ng/mL	1.9 (0.3)	1.9 (0.4)
Total protein, g/dL	7.1 (1.1)	7.2 (1.0)
Sodium, mmol/L	136.2 (6.0)	136.6 (7.0)
Troponin, ng/mL	0.2 (0.1)	0.2 (0.1)

Note that *n* represents a counting unit. The asterisk (\*) denotes statistically significant difference with  $p < 0.01$  between the training and test sets. We used the two-sided *t*-test to compare continuous variables (age, BMI, weight, vital signs, and laboratory tests) and a 2-sample *z*-test to compare the proportions of categorical variables (sex and provision of supplemental oxygen).



**Table 2.** Performance of the outcome classification task on the held-out test set, and on the subset of the test set used in the reader study (*n* represents the number of images).

Test set ( <i>n</i> = 832)								
	AUC				PR AUC			
	24 h	48 h	72 h	96 h	24 h	48 h	72 h	96 h
COVID-GBM	0.747 (0.698, 0.802)	0.739 (0.69, 0.795)	0.750 (0.703, 0.799)	0.770 (0.727, 0.813)	0.230 (0.139, 0.296)	0.325 (0.229, 0.396)	0.408 (0.317, 0.479)	<b>0.523</b> (0.433, 0.6)
COVID-GMIC	0.695 (0.636, 0.763)	0.716 (0.666, 0.771)	0.717 (0.668, 0.773)	0.738 (0.695, 0.785)	0.200 (0.119, 0.260)	0.302 (0.209, 0.379)	0.374 (0.283, 0.452)	0.439 (0.346, 0.515)
COVID-GBM + COVID-GMIC	<b>0.765</b> (0.712, 0.817)	<b>0.749</b> (0.700, 0.798)	<b>0.769</b> (0.724, 0.818)	<b>0.786</b> (0.745, 0.830)	<b>0.243</b> (0.150, 0.299)	<b>0.332</b> (0.237, 0.41)	<b>0.439</b> (0.345, 0.527)	0.517 (0.429, 0.600)
Reader study dataset ( <i>n</i> = 200)								
	AUC				PR AUC			
	24 h	48 h	72 h	96 h	24 h	48 h	72 h	96 h
Radiologist A	0.613 (0.519, 0.705)	0.645 (0.571, 0.731)	0.691 (0.618, 0.77)	0.740 (0.674, 0.814)	0.346 (0.217, 0.441)	0.490 (0.367, 0.599)	0.640 (0.536, 0.745)	0.742 (0.657, 0.834)
Radiologist B	0.637 (0.547, 0.73)	0.636 (0.552, 0.716)	0.658 (0.588, 0.738)	0.713 (0.649, 0.786)	0.365 (0.229, 0.462)	0.460 (0.335, 0.56)	0.590 (0.492, 0.701)	0.704 (0.616, 0.805)
Radiologist A + Radiologist B	<b>0.642</b> (0.555, 0.729)	0.663 (0.589, 0.746)	0.692 (0.621, 0.766)	0.741 (0.678, 0.809)	<b>0.403</b> (0.272, 0.52)	0.499 (0.380, 0.613)	0.609 (0.492, 0.711)	0.740 (0.650, 0.831)
COVID-GMIC	<b>0.642</b> (0.554, 0.734)	<b>0.701</b> (0.627, 0.781)	<b>0.751</b> (0.685, 0.821)	<b>0.808</b> (0.75, 0.87)	0.381 (0.235, 0.480)	<b>0.546</b> (0.421, 0.657)	<b>0.676</b> (0.564, 0.780)	<b>0.789</b> (0.699, 0.880)
COVID-GBM	0.704 (0.632, 0.784)	0.719 (0.648, 0.794)	0.750 (0.684, 0.821)	0.787 (0.727, 0.850)	0.411 (0.259, 0.518)	0.537 (0.394, 0.64)	0.668 (0.558, 0.77)	0.804 (0.738, 0.884)
COVID-GBM + COVID-GMIC	0.708 (0.637, 0.799)	0.702 (0.633, 0.775)	0.778 (0.719, 0.851)	0.819 (0.763, 0.885)	0.411 (0.279, 0.517)	0.500 (0.364, 0.601)	0.705 (0.599, 0.806)	0.808 (0.735, 0.898)

We include 95% confidence intervals estimated by 1000 iterations of the bootstrap method<sup>30</sup>. The optimal weights assigned to the COVID-GMIC prediction in the COVID-GMIC and COVID-GBM ensemble were derived through optimizing the AUC on the validation set as described in Supplementary Fig. 3b. The ensemble of COVID-GMIC and COVID-GBM, denoted as ‘COVID-GMIC + COVID-GBM’, achieves the best performance across all time windows in terms of the AUC and PRAUC, except for the PR AUC in the 96 h task. In the reader study, our main finding is that COVID-GMIC outperforms radiologists A & B across time windows longer than 24 h, with 3 and 17 years of experience, respectively. Note that the radiologists did not have access to clinical variables and as such their performance is not directly comparable to the COVID-GBM model; we include it only for reference. The area under the precision-recall curve is sensitive to class distribution, which explains the large differences between the scores on the test set and the reader study subset. Best performance per metric is shown in bold.

examples, the saliency maps highlight regions that contain visual patterns such as airspace opacities and consolidation, which are correlated with clinical deterioration<sup>22,24</sup>. These saliency maps are utilized to guide the extraction of six regions of interest (ROI) patches cropped from the entire image, which are then associated with a score that indicates its relevance to the prediction task. We also note that in the last example, the saliency maps highlight right mid to lower paramediastinal and left mid-lung periphery. The dense consolidation in the periphery of the right upper lobe is highlighted by ROI patch 4. It might also be useful to enhance GMIC through a classifier agnostic mechanism<sup>31</sup>, which finds all the useful evidence in the image, instead of solely the most discriminative part. We leave this for future work.

The most predictive features (top 10) of COVID-GBM are shown in Supplementary Fig. 3a. Temperature was ranked among the top two predictive features and age was ranked among the top four predictive features across all time windows.

### Comparison to radiologists

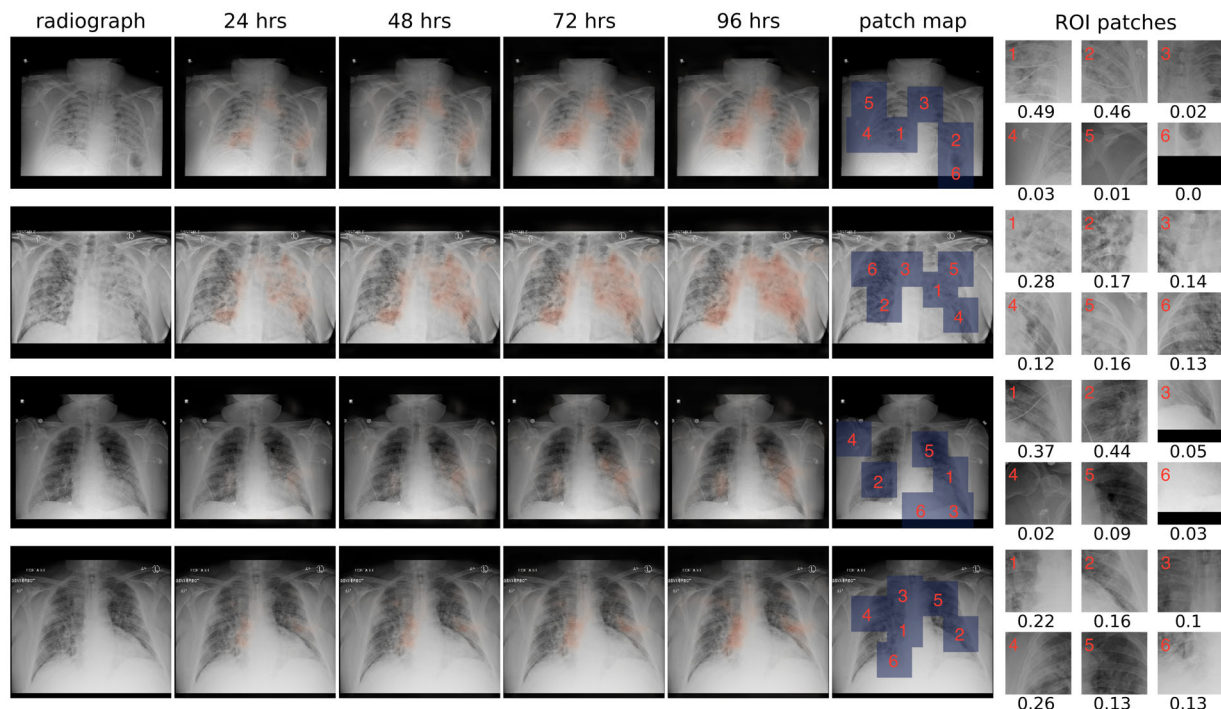
We compared the performance of COVID-GMIC with two chest radiologists from NYU Langone Health (with 3 and 17 years of experience) in a reader study with a sample of 200 frontal chest

X-ray exams from the test set. We used stratified sampling to improve the representation of patients with a negative outcome in the reader study dataset. Specifically, we randomly sampled the first 100 exams from patients that had an adverse event in the next 96 h from the time the exam was taken. The remaining 100 exams came from the complement of the test set. We describe the design of the reader study in more detail in the Methods section.

As shown in Table 2, our main finding is that COVID-GMIC achieves a comparable performance to radiologists across all time windows in terms of AUC and PR AUC, and outperforms radiologists for 48, 72, and 96 h. For example, COVID-GMIC achieves AUC of 0.808 (95% CI, 0.746–0.866) compared to AUC of 0.741 average AUC of both radiologists in the 96 h prediction task. We hypothesize that COVID-GMIC outperforms radiologists on this task due to the currently limited clinical understanding of which pulmonary parenchymal patterns predict clinical deterioration, rather than the severity of lung involvement<sup>24</sup>. Supplementary Fig. 5 shows AUC and PR AUC curves across all time windows.

### Deterioration risk curves

We use a modified version of COVID-GMIC, referred to hereafter as COVID-GMIC-DRC, to generate discretized deterioration risk curves



**Fig. 3 Explainability of COVID-GMIC.** From left to right: the original X-ray image, saliency maps for clinical deterioration within 24, 48, 72, and 96 h, locations of region-of-interest (ROI) patches, and ROI patches with their associated attention scores. All four patients were admitted to the intensive care unit and were intubated within 48 h. In the first example, there are diffuse airspace opacities, though the saliency maps primarily highlight the medial right basal and peripheral left basal opacities. Similarly, the two ROI patches (1 and 2) on the basal region demonstrate comparable attention values, 0.49 and 0.46, respectively. In the second example, the extensive left mid to upper-lung abnormality (ROI patch 1) is highlighted, which correlates with the most extensive area of parenchymal consolidation. In the third example, the saliency maps highlight the left mid lung (ROI patch 1) and right hilar/infrahilar regions (ROI patch 2) which show groundglass opacities. In the last example, the saliency maps highlight the right infrahilar region (ROI patch 1) and the left mid lung periphery (ROI patch 2). The ROI patch 4 is also assigned the highest attention score as a predictive region of clinical deterioration, which corresponds to dense peripheral right upper lobe consolidation.

(DRCs) which predict the evaluation of the deterioration risk based on chest X-ray images. Figure 4a shows the DRCs for all the patients in the test set. The DRC represents the probability that the first adverse event occurs before time  $t$ , where  $t$  is equal to 3, 12, 24, 48, 72, 96, 144, and 192 h. The mean DRCs of patients who deteriorate (red bold line) is significantly higher than the mean DRCs of patients who are discharged without experiencing any adverse events (blue bold line). We evaluate the performance of the model using the concordance index, which is computed on patients in the test set who experienced adverse events. For a fixed time  $t$  the index equals the fraction of pairs of patients in the test data for which the patient with the higher DRC value at  $t$  experiences an adverse event earlier. For  $t$  equal to 96 h, the concordance index is 0.713 (95% CI: 0.682–0.747), which demonstrates that COVID-GMIC-DRC can effectively discriminate between patients. Other values of  $t$  yield similar results, as reported in Supplementary Table 2. Sample learning curves of COVID-GMIC-DRC are shown in Supplementary Fig. 6 for reference.

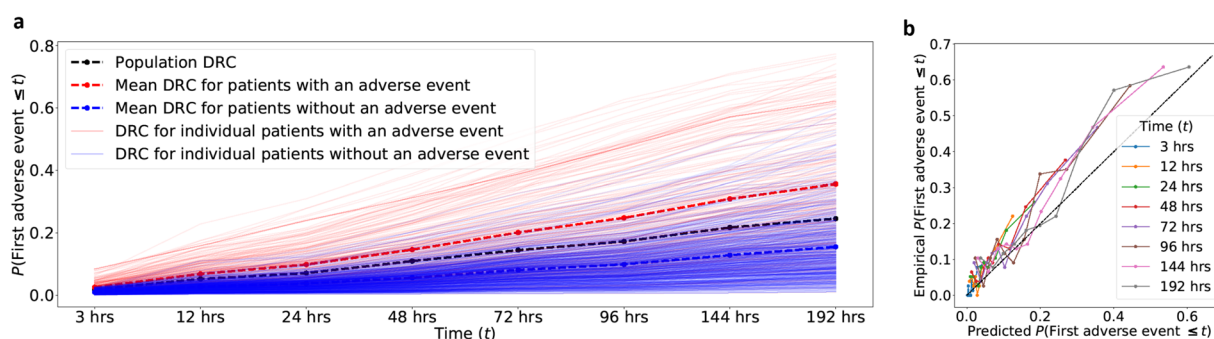
Figure 4b shows a reliability plot, which evaluates the calibration of the probabilities encoded in the DRCs. The diagram compares the values of the estimated DRCs for the patients in the test set with empirical probabilities that represent the true frequency of adverse events. To compute the empirical probabilities, we divided the patients into deciles according to the value of the DRC at each time  $t$ . We then computed the fraction of patients in each decile that suffered adverse events up to  $t$ . The fraction is plotted against the mean DRC of the patients in the decile. The diagram shows that these values are similar across the different values of  $t$ , meaning the model is well-calibrated (for

comparison, perfect calibration would correspond to the diagonal black dashed line).

### Prospective silent validation in a clinical setting

Our long-term goal is to deploy our system in existing clinical workflows to assist clinicians. The clinical implementation of machine learning models is a very challenging process, both from technical and organizational standpoints<sup>32</sup>. To test the feasibility of deploying the AI system in the hospital, we silently deployed a preliminary version of our AI system in the hospital system and let it operate in real-time beginning on May 22, 2020. The deployed version includes 15 models that are based on DenseNet-121 architectures, and use only chest X-ray images. The models were developed to predict deterioration within 96 h using a subset of our data collected prior to deployment from 3425 patients. The models were serialized and served with TensorFlow Serving components<sup>33</sup> on an Intel(R) Xeon(R) Gold 6154 CPU @ 3.00GHz; no GPUs were used. Images are preprocessed as explained in the Methods section. Our system produces predictions essentially in real-time - it takes ~2 s to extract an image from the PACS system, apply the image preprocessing steps, and get the prediction of a model as a TensorFlow<sup>33</sup> output.

A total of 375 exams were collected between May 22, 2020 and June 24, 2020. Of the 375 exams collected between May 22, 2020 and June 24, 2020, 38 exams (10.1%) were associated with a positive 96 h deterioration outcome, compared to 20.3% in the retrospective test set. When we compare the composition of types of events between the two cohorts, we can observe a significant difference. Mortality within 96 h is the most prevalent event in the



**Fig. 4** Deterioration risk curves (DRCs) and reliability plot for COVID-GMIC-DRC. **a** DRCs generated by the COVID-GMIC-DRC model for patients in the test set with (faded red lines) and without adverse events (faded blue lines). The mean DRC for patients with adverse events (red dashed line) is higher than the DRC for patients without adverse events (blue dashed line) at all times. The graph also includes the ground-truth population DRC (black dashed line) computed from the test data. **b** Reliability plot of the DRCs generated by the COVID-GMIC-DRC model for patients in the test set. The empirical probabilities are computed by dividing the patients into deciles according to the value of the DRC at each time  $t$ . The empirical probability equals the fraction of patients in each decile that suffered adverse events up to  $t$ . This is plotted against the predicted probability, which equals the mean DRC of the patients in the decile. Perfect calibration is indicated by the diagonal black dashed line.

retrospective cohort (43.1%), while ICU admission within 96 h is the most prevalent event in the silent evaluation cohort (65.8%). Additionally, mortality constitutes 18.4% of adverse events and intubation constitutes 15.8% of events in the silent cohort. An ensemble of the deployed models, obtained by averaging their predictions, achieved an AUC of 0.717 (95% CI: 0.622–0.801) and a PR AUC of 0.289 (95% CI: 0.181–0.465). These results are comparable to those obtained on a retrospective test set used for evaluation before deployment, which are 0.748 (95% CI: 0.708–0.790) AUC and 0.365 (95% CI: 0.313–0.465) PR AUC. The decrease in accuracy is expected and may indicate changes in the patient population, increased class imbalance, and treatment guidelines as the pandemic progressed. When practically deployed, our system would still need periodical retraining with the latest data.

## DISCUSSION

In this work, we present an AI system that is able to predict deterioration of COVID-19 patients presenting to the ED, where deterioration is defined as the composite outcome of mortality, intubation, or ICU admission. The system aims to provide clinicians with a quantitative estimate of the risk of deterioration, and how it is expected to evolve over time, in order to enable efficient triage and prioritization of patients at the high risk of deterioration. The tool may be of particular interest for pandemic hotspots where triage at admission is critical to allocate limited resources such as hospital beds.

Recent studies have shown that chest X-ray images are useful for the diagnosis of COVID-19<sup>12,13,15,19,34</sup>. Our work supplements those studies by demonstrating the significance of this imaging modality for COVID-19 prognosis. Additionally, our results suggest that chest X-ray images and routinely collected clinical variables contain complementary information, and that it is best to use both to predict clinical deterioration. This builds upon existing prognostic research, which typically focuses on developing risk prediction models using non-imaging variables extracted from electronic health records<sup>19,35</sup>. In Supplementary Table 3, we demonstrate that our models' performance can be improved by increasing the dataset size. The current dearth of prognosis models that use both imaging and clinical variables may partly be due to the limited availability of large-scale datasets including both data types and outcome labels, which is a key strength of our study. In order to assess the clinical benefits of our approach, we conducted a reader study, and the results indicate that the proposed system can perform comparably to radiologists.

This highlights the potential of data-driven tools for assisting the interpretation of X-ray images.

The proposed deep learning model, COVID-GMIC, provides visually intuitive saliency maps to help clinicians interpret the model predictions<sup>36</sup>. Existing work on COVID-19 often use external gradient-based algorithms, such as gradCAM<sup>37</sup>, to interpret deep neural network classifiers<sup>38–40</sup>. However, visualizations generated by gradient-based methods are sensitive to minor perturbation in input images, and could yield misleading interpretations<sup>41</sup>. In contrast, COVID-GMIC has an inherently interpretable architecture that better retains localization information of the more informative regions in the input images. We assessed the model's interpretability qualitatively due to the difficulty in obtaining ground-truth segmentation labels from radiologists during the pandemic. Assessing the interpretability of GMIC quantitatively, by measuring its ability to indicate the same areas in the images as the radiologists indicate, is an area of future work.

We performed prospective validation of an early version of our system through silent deployment in an NYU Langone Health hospital. The results suggest that the implementation of our AI system in the existing clinical workflows is feasible. Our model does not incur any overhead operational costs on data collection, since chest X-ray images are routinely collected from COVID-19 patients. Additionally, the model can process the image efficiently in real-time, without requiring extensive computational resources such as GPUs. This is an important strength of our study, since very few studies have implemented and prospectively validated risk prediction models in general<sup>42</sup>. Our approach has some limitations that will be addressed in future work. Our deep neural network considers a single chest X-ray image as an input and does not consider longitudinal changes in consecutive images. This is primarily due to our focus on emergency department triage, where the patient typically gets only one or a few scans. Another limitation is that the silent deployment was based only on the model that processes chest X-ray exams, and did not include routine clinical variables, nor any interventions. The COVID-GMIC-DRC model also did not incorporate any clinical variables. This is because the computation of the deterioration risk curves heavily depends on model calibration. Gradient boosting models are generally not as well calibrated as neural networks. Therefore, incorporating the clinical variables within the DRC model requires more extensive calibration analysis and/or the design of an additional neural network for clinical variables. We developed the DRC model to meet a secondary objective of our study, which is to assess whether chest X-rays contain useful information for survival analysis. Our future work will focus on calibration in the context of



multi-modal learning for survival analysis. In addition, an inherent limitation of the type of our study, an internal retrospective and prospective validation, is that the system's performance measures may be affected when COVID-19 outcomes have a different prevalence compared to at the height of the pandemic or when different imaging protocols are used. Therefore, further validation is required to assess whether the system can improve key performance measures, such as patient outcomes, through prospective and external validation across different hospitals and electronic health records systems.

Our system currently considers two data types, which are chest X-ray images and clinical variables. The multi-modal system adopts a late fusion strategy to combine the predictions of COVID-GMIC and COVID-GBM. Future work should investigate more sophisticated multi-modal learning strategies that cross-transfer information between the two modalities to improve performance and understanding of the utility of clinical data. Incorporating additional data from patient health records may also further improve its performance. For example, the inclusion of presenting symptoms using natural language processing has been shown to improve the performance of a risk prediction model in the ED<sup>25</sup>. Although we focus on chest X-ray images because pulmonary disease is the main complication associated with COVID-19, COVID-19 patients may also suffer poor outcomes due to non-pulmonary complications such as: non-pulmonary thromboembolic events, stroke, and pediatric inflammatory syndromes<sup>43–45</sup>. This could explain some of the false negatives incurred by our system; therefore, incorporating other types of data that reflect non-pulmonary complications may also improve prognostic accuracy.

Our system was developed and evaluated using data collected from the NYU Langone Health in New York, USA. Therefore, it is possible that our models overfit to the patient demographics and specific configurations in the imaging acquisition devices of our dataset.

Our findings show the promise of data-driven AI systems in predicting the risk of deterioration for COVID-19 patients, and highlights the importance of designing multi-modal AI systems capable of processing different types of data. We anticipate that such tools will play an increasingly important role in supporting clinical decision-making in the future.

## METHODS

### Ethics

This study was approved by the NYU Langone Health Institutional Review Board (IRB), with ID# i20-00858. A waiver for informed consent was granted by the IRB, since the study presents no more than minimal risk.

### Outline

In this section, we first introduce our image preprocessing pipeline then formulate the adverse event prediction task and present our multi-modal approach which utilizes both chest X-ray images and clinical variables. Next, we formally define deterioration risk curve (DRC) and introduce our X-ray image-based approach to estimate DRC. Subsequently, we summarize the technical details of model training and implementation. Lastly, we describe the design of the reader study.

### Image preprocessing

After extracting the images from DICOM files, we applied the following preprocessing procedure. We first thresholded and normalized pixel values, and then cropped the images to remove any zero-valued pixels surrounding the image. Then, we unified the dimensions of all images by cropping the images outside the center and rescaling. We performed data augmentation by applying random horizontal flipping ( $p = 0.5$ ), random rotation ( $-45^\circ$  to  $45^\circ$ ), and random translation. Supplementary Fig. 7 shows the distribution of the size of the images prior to data augmentation, as well as examples of images before and after preprocessing.

### Adverse event prediction

Our main goal is to predict clinical deterioration within four time windows of 24, 48, 72, and 96 h. We frame this as a multi-label classification task with binary labels  $\mathbf{y} = [y^{24}, y^{48}, y^{72}, y^{96}]$  indicating clinical deterioration of a patient within the four time windows. The probability of deterioration is estimated using two types of data associated with the patient: a chest X-ray image, and routine clinical variables. We use two different machine learning models for this task: COVID-GMIC to process chest X-ray images, and COVID-GBM to process clinical variables. For each time window  $t \in \mathbb{T}_a = \{24, 48, 72, 96\}$ , both models produce probability estimates of clinical deterioration,  $\hat{\mathbf{y}}_{\text{COVID-GMIC}}^t, \hat{\mathbf{y}}_{\text{COVID-GBM}}^t \in [0, 1]$ .

In order to combine the predictions from COVID-GMIC and COVID-GBM, we employ the technique of model ensembling<sup>46</sup>. Specifically, for each example, we compute a multi-modal prediction  $\hat{\mathbf{y}}_{\text{ENSEMBLE}}$  as a linear combination of  $\hat{\mathbf{y}}_{\text{COVID-GMIC}}$  and  $\hat{\mathbf{y}}_{\text{COVID-GBM}}$ :

$$\hat{\mathbf{y}}_{\text{ENSEMBLE}} = \lambda \hat{\mathbf{y}}_{\text{COVID-GMIC}} + (1 - \lambda) \hat{\mathbf{y}}_{\text{COVID-GBM}}, \quad (1)$$

where  $\lambda \in [0, 1]$  is a hyperparameter. We selected the best  $\lambda$  by optimizing the average of the AUC and PR AUC on the validation set. In Supplementary Fig. 3b, we show the validation performance of  $\hat{\mathbf{y}}_{\text{ENSEMBLE}}$  for varying  $\lambda$ .

### Clinical variables model

The goal of the clinical variables model is to predict the risk of deterioration when the patient's vital signs are measured. Thus, each prediction was computed using a set of vital sign measurements, in addition to the patient's most recent laboratory test results, age, weight, and body mass index (BMI). The vital signs (7 in total) were heart rate, respiratory rate, temperature, systolic blood pressure, diastolic blood pressure, oxygen saturation, and the provision of supplemental oxygen. The laboratory test measurements (24 in total) were albumin, alanine transaminase, aspartate aminotransferase, total bilirubin, blood urea nitrogen, calcium, chloride, creatinine, d-dimer, eosinophils count, eosinophils percentage, hematocrit, lactate dehydrogenase, lymphocytes count, lymphocytes percentage, platelet volume, neutrophils count, neutrophils percentage, platelet, potassium, procalcitonin, total protein, sodium, and troponin. The laboratory test features were further represented as maximum and minimum statistics of any results collected within 12 h prior to the time of the vital sign measurement, leading to 48 processed features in total. The feature sets of age, weight, BMI, vital signs and processed laboratory tests (58 input features in total) were then processed using a gradient boosting model<sup>28</sup> which we refer to as COVID-GBM. In cases where a patient had a missing vital sign or laboratory test measurement, we carried forward the most recently recorded measurement. If there were no recent measurements, then the value was left as missing since GBM can handle missing values. For the final ensemble prediction,  $\hat{\mathbf{y}}_{\text{ENSEMBLE}}$ , we combined the COVID-GMIC prediction with the COVID-GBM prediction computed using the most recently collected clinical variables prior to the chest X-ray exam. In cases where there were no clinical variables collected prior to the chest X-ray (i.e., missing clinical variables), we performed a mean imputation of the predictions assigned to the validation set.

### Chest X-ray image model

We process chest X-ray images using a deep convolutional neural network model, which we call COVID-GMIC, based on the GMIC model<sup>26,27</sup>. COVID-GMIC has two desirable properties. First, COVID-GMIC generates interpretable saliency maps that highlight regions in the X-ray images that correlate with clinical deterioration. Second, it possesses a local module that is able to utilize high-resolution information in a memory-efficient manner. This avoids aggressive downsampling of the input image, a technique that is commonly used on natural images<sup>47,48</sup>, which may distort and blur informative visual patterns in chest X-ray images such as basilar opacities and pulmonary consolidation. In Supplementary Table 4, we demonstrate that COVID-GMIC achieves comparable results to DenseNet-121, a neural network model that is not interpretable by design, but is commonly used for chest X-ray analysis<sup>49–52</sup>.

The architecture of COVID-GMIC is schematically depicted in Fig. 1b. COVID-GMIC processes an X-ray image  $\mathbf{x} \in \mathbb{R}^{h,w}$  ( $h$  and  $w$  denote the height and width) in three steps. First, the global module helps COVID-GMIC learn an overall view of the X-ray image. Within this module, COVID-GMIC utilizes a global network  $f_g$  to extract feature maps  $\mathbf{h}_g \in \mathbb{R}^{h,w,n}$ , where  $h$ ,  $w$ , and  $n$  denote the height, width, and number of channels of the



feature maps. The resolution of the feature maps is chosen to be coarser than the resolution of the input image. For each time window  $t \in \mathbb{T}_a$ , we apply a  $1 \times 1$  convolution layer with sigmoid activation to transform  $\mathbf{h}_g$  into a saliency map  $\mathbf{A}^t \in \mathbb{R}^{h,w}$  that highlights regions on the X-ray image which correlate with clinical deterioration. For visualization purposes, we apply nearest neighbor interpolation to upsample the saliency maps to match the resolution of the original image. Each element  $\mathbf{A}_{ij}^t \in [0, 1]$  represents the contribution of the spatial location  $(i, j)$  in predicting the onset of adverse events within time window  $t$ . In order to train  $f_g$ , we use an aggregation function  $f_{agg} : \mathbb{R}^{h,w} \mapsto [0, 1]$  to transform all saliency maps  $\mathbf{A}^t$  for all time windows  $t$  into classification predictions  $\hat{\mathbf{y}}_{global}$ :

$$f_{agg}(\mathbf{A}^t) = \frac{1}{|H^+|} \sum_{(i,j) \in H^+} \mathbf{A}_{ij}^t, \quad (2)$$

where  $H^+$  denotes the set containing the locations of the  $r\%$  largest values in  $\mathbf{A}^t$ , and  $r$  is a hyperparameter.

The local module enables COVID-GMIC to selectively focus on a small set of informative regions. As shown in Fig. 1, COVID-GMIC utilizes the saliency maps, which contain the approximate locations of informative regions, to retrieve six image patches from the input X-ray image, which we call region-of-interest (ROI) patches. We refer the readers to Supplementary Note 5 for more details about the ROI retrieval algorithm. Figure 3 shows some examples of ROI patches. To utilize high-resolution information within each ROI patch  $\tilde{\mathbf{x}} \in \mathbb{R}^{224,224}$ , COVID-GMIC applies a local network  $f_l$ , parameterized as a ResNet-18<sup>47</sup>, which produces a feature vector  $\mathbf{h}_k \in \mathbb{R}^{512}$  from each ROI patch. The predictive value of each ROI patch might vary significantly. Therefore, we utilize the gated attention mechanism<sup>53</sup> to compute an attention score  $a_k \in [0, 1]$  that indicates the relevance of each ROI patch  $\tilde{\mathbf{x}}$  for the prediction task. To aggregate information from all ROI patches, we compute an attention-weighted representation:

$$\mathbf{z} = \sum_{k=1}^6 a_k \tilde{\mathbf{h}}_k. \quad (3)$$

The representation  $\mathbf{z}$  is then passed into a fully connected layer with sigmoid activation to generate a prediction  $\hat{\mathbf{y}}_{local}$ . We refer the readers to Shen et al.<sup>27</sup> for further details.

The fusion module combines both global and local information to compute a final prediction. We apply global max pooling to  $\mathbf{h}_g$ , and concatenate it with  $\mathbf{z}$  to combine information from both saliency maps and ROI patches. The concatenated representation is then fed into a fully connected layer with sigmoid activation to produce the final prediction  $\hat{\mathbf{y}}_{fusion}$ .

In our experiments, we chose  $H=W=1024$ . Supplementary Table 4 shows that COVID-GMIC achieves the best validation performance for this resolution. We parameterize  $f_g$  as a ResNet-18<sup>47</sup> which yields feature maps  $\mathbf{h}^g$  with resolution  $h=w=32$ , and number of channels  $n=512$ . During training, we optimize the loss function:

$$\begin{aligned} \mathcal{L}(\mathbf{y}; \hat{\mathbf{y}}_{global}, \hat{\mathbf{y}}_{local}, \hat{\mathbf{y}}_{fusion}) = & \frac{1}{|\mathbb{T}_a|} \sum_{t \in \mathbb{T}_a} \text{BCE}(\mathbf{y}^t, \hat{\mathbf{y}}_{global}^t) + \text{BCE}(\mathbf{y}^t, \hat{\mathbf{y}}_{local}^t) \\ & + \text{BCE}(\mathbf{y}^t, \hat{\mathbf{y}}_{fusion}^t) + \beta |\mathbf{A}^t|, \end{aligned} \quad (4)$$

where BCE denotes binary cross-entropy and  $\beta$  is a hyperparameter representing the relative weights on an  $\ell_1$ -norm regularization term that promotes sparsity of the saliency maps. During inference, we use  $\hat{\mathbf{y}}_{fusion}$  as the final prediction generated by the model.

### Estimation of deterioration risk curves

The deterioration risk curve (DRC) represents the evolution of the deterioration risk over time for each patient. Let  $T$  denote the time of the first adverse event. The DRC is defined as a discretized curve that equals the probability  $P(T \leq t_i)$  of the first adverse event occurring before time  $t_i \in \{t_i | 1 \leq i \leq 8\}$ , where  $t_1 = 3$ ,  $t_2 = 12$ ,  $t_3 = 24$ ,  $t_4 = 48$ ,  $t_5 = 72$ ,  $t_6 = 96$ ,  $t_7 = 144$ ,  $t_8 = 192$  (all times are in hours).

Following recent work on survival analysis via deep learning<sup>54</sup>, we parameterize the DRC using a vector of conditional probabilities  $\hat{\mathbf{p}} \in \mathbb{R}^8$ . The  $i$ th entry of this vector,  $\hat{p}_i$ , is equal to the conditional probability of the adverse event happening before time  $t_i$  given that no adverse event occurred before time  $t_{i-1}$ , that is:

$$\hat{p}_i = \begin{cases} P(T \leq t_1), & i = 1, \\ P(T \leq t_i | T > t_{i-1}), & 2 \leq i \leq 8. \end{cases} \quad (5)$$

The parameters in our implementation are the complementary

probabilities  $\hat{\mathbf{q}} = 1 - \hat{\mathbf{p}}$ , which is a mathematically equivalent parameterization. We also include an additional parameter to account for patients whose first adverse event occurs after 192 h. Given an estimate of  $\hat{\mathbf{p}}$ , the DRC can be computed applying the chain rule:

$$\begin{aligned} \text{DRC}(t_i) &= P(T \leq t_i) \\ &= 1 - P(T > t_i) \\ &= 1 - \prod_{j=1}^i P(T > t_j | T > t_{j-1}) \\ &= 1 - \prod_{j=1}^i (1 - \hat{p}_j). \end{aligned} \quad (6)$$

We use the GMIC model to estimate the conditional probabilities  $\hat{\mathbf{p}}$  from chest X-ray images. We refer to this model as COVID-GMIC-DRC. As explained in the previous section, the GMIC model has three different outputs corresponding to the global module, local module and fusion module. When estimating conditional probabilities for the eight time intervals, we denote these outputs by  $\hat{\mathbf{p}}_{global}$ ,  $\hat{\mathbf{p}}_{local}$ , and  $\hat{\mathbf{p}}_{fusion}$ . During inference, we use the output of the fusion module,  $\hat{\mathbf{p}}_{fusion}$ , as the final prediction of the conditional-probability vector  $\hat{\mathbf{p}}$ . We use an input resolution of  $H=W=512$  and parameterize  $f_g$  as ResNet-34<sup>47</sup>. The resulting feature maps  $\mathbf{h}_g$  have resolution  $h=w=16$  and number of channels  $n=512$ . The results of an ablation study that evaluates the impact of input resolution and compares COVID-GMIC-DRC to a model based on the Densenet-121 architecture, are shown in the Supplementary Fig. 8 and Supplementary Tables 2 and 4. During training, we minimize the following loss function defined on a single example:

$$\mathcal{L}(T; \hat{\mathbf{p}}_{global}, \hat{\mathbf{p}}_{local}, \hat{\mathbf{p}}_{fusion}) = \mathcal{L}_s(T; \hat{\mathbf{p}}_{global}) + \mathcal{L}_s(T; \hat{\mathbf{p}}_{local}) + \mathcal{L}_s(T; \hat{\mathbf{p}}_{fusion}) + \sum_{m=0}^8 \beta |\mathbf{A}^m|, \quad (7)$$

where  $\mathcal{L}_s$  is the negative log-likelihood of the conditional probabilities. For a patient who had an adverse event between  $t_{i-1}$  and  $t_i$  (where  $t_0 = 0$ ), this negative log-likelihood is given by

$$\begin{aligned} \mathcal{L}_s(T; \hat{\mathbf{p}}) &= -\ln P(t_{i-1} \leq T \leq t_i) \\ &= -\ln \prod_{j=1}^{i-1} P(T > t_j | T > t_{j-1}) P(T \leq t_i | T > t_{i-1}) \\ &= -\sum_{j=1}^{i-1} \ln(1 - \hat{p}_j) - \ln \hat{p}_i. \end{aligned} \quad (8)$$

The framework can easily incorporate censored data corresponding to patients whose information is not available after a certain point. The negative log-likelihood corresponding to a patient, who has no information after  $t_i$  and no adverse events before  $t_i$ , equals

$$\begin{aligned} \mathcal{L}_s(T; \hat{\mathbf{p}}) &= -\ln P(T > t_i) \\ &= -\ln \prod_{j=1}^i P(T > t_j | T > t_{j-1}) \\ &= -\sum_{j=1}^i \ln(1 - \hat{p}_j). \end{aligned} \quad (9)$$

Note that each  $\hat{p}_i$  is estimated only using patients that have data available up to  $t_i$ . The total negative log-likelihood of the training set is equal to the sum of the individual negative log-likelihoods corresponding to each patient, which makes it possible to perform minimization efficiently via stochastic gradient descent. In contrast, deep learning models for survival analysis based on Cox proportional hazards regression<sup>55</sup> require using the whole dataset to perform model updates<sup>56–58</sup>, which is computationally infeasible when processing large image datasets.

### Model training and selection

In this section, we discuss the experimental setup used for COVID-GMIC, COVID-GMIC-DRC, and COVID-GMIC-BM. We initialized the weights of COVID-GMIC and COVID-GMIC-DRC by pretraining them on the ChestX-ray14 dataset<sup>59</sup> (Supplementary Table 5 compares the performance of different initialization strategies). We used Adam<sup>60</sup> with a minibatch size of eight to train the models on our data. During the training and test stages, we applied a set of data transformations to the inputs in order to make the model more robust to rotation and spatial translation. During the test stage, we applied ten different augmentations to each image and used the average of their predictions in order to further improve performance.

We did not apply any data augmentation during the validation stage since it introduces randomness, which can be confounding when determining whether or not validation performance is improving.

We optimized the hyperparameters using random search<sup>61</sup>. For COVID-GMIC, we searched for the learning rate  $\eta \in 10^{[-6, -4]}$  on a logarithmic scale, the regularization hyperparameter  $\beta \in 4 \times 10^{[-6, -3]}$  on a logarithmic scale, and the pooling threshold  $r \in [0.2, 0.8]$  on a linear scale. For COVID-GMIC-DRC, based on the preliminary experiments, we fixed the learning rate to  $1.25 \times 10^{-4}$ . We searched for the regularization hyperparameter,  $\beta \in 10^{[-6, -4]}$  on a logarithmic scale, and the pooling threshold  $r \in \{0.2, 0.5, 0.8\}$ . For COVID-GBM, we searched for the learning rate  $\eta \in 10^{[-2, -1]}$  on a logarithmic scale, the number of estimators  $e \in 10^{[2, 3]}$  on a logarithmic scale, and the number of leaves  $l \in [5, 15]$  on a linear scale. For each hyperparameter configuration, we performed Monte Carlo cross-validation<sup>62</sup> (we sampled 80% of the data for training and 20% of the data was used for validation). We performed cross-validation using three different random splits for each hyperparameter configuration. We then selected the top three hyperparameter configurations based on the average validation performance across the three splits. Finally, we combined the nine models from the top three hyperparameter configurations by averaging their predictions on the held-out test set to evaluate the performance. This procedure is formally described in Supplementary Algorithm 1.

## Software

The chest X-ray image models were implemented in PyTorch<sup>63</sup> and trained using NVIDIA Tesla V100 GPUs. The clinical variables models were implemented using the Python library LightGBM<sup>28</sup>.

## Design of the reader study

The reader study consists of 200 frontal chest X-ray exams from the test set. We selected one exam per patient to increase the diversity of exams. We used stratified sampling to ensure that a sufficient number of exams in the study corresponded to the least common outcome (patients with adverse outcomes in the next 24 h). In more detail, we oversampled exams of patients who developed an adverse event by sampling the first 100 exams only from patients from the test set that had an adverse outcome within the first 96 h. The remaining 100 exams came from the remaining patients in the test set. The radiologists were asked to assign the overall probability of deterioration to each scan across all time windows of evaluation.

## Reporting summary

Further information on research design is available in the Nature Research Reporting Summary linked to this article.

## DATA AVAILABILITY

The ImageNet dataset is available at <http://www.image-net.org/>. The ChestX-ray8 dataset is available at <https://nihcc.app.box.com/v/ChestXray-NIHCC>. The COVID-19 X-ray images and associated clinical variables from NYU Langone Health are not publicly available, but we provide sample patients in our source code repository.

## CODE AVAILABILITY

The code of the models in this study, along with their trained weights, are available at [https://github.com/nyukat/COVID-19\\_prognosis](https://github.com/nyukat/COVID-19_prognosis).

Received: 4 November 2020; Accepted: 19 March 2021;

Published online: 12 May 2021

## REFERENCES

- Baugh, J. J. et al. Creating a COVID-19 surge clinic to offload the emergency department. *Am. J. Emerg. Med.* **38**, 1535–1537 (2020).
- Debnath, S. et al. Machine learning to assist clinical decision-making during the COVID-19 pandemic. *Bioelectron. Med.* **6**, 1–8 (2020).
- Whiteside, T., Kane, E., Aljohani, B., Alsamman, M. & Pourmand, A. Redesigning emergency department operations amidst a viral pandemic. *Am. J. Emerg. Med.* **38**, 1448–1453 (2020).
- Dorsett, M. Point of no return: COVID-19 and the us health care system: an emergency physician's perspective. *Sci. Adv.* **6**, eabc5354 (2020).
- McKenna, P. et al. Emergency department and hospital crowding: causes, consequences, and cures. *Clin. Exp. Emerg. Med.* **6**, 189 (2019).
- Warner, M. A. Stop doing needless things! Saving healthcare resources during COVID-19 and beyond. *J. Gen. Intern. Med.* **35**, 2186–2188 (2020).
- Cozzi, D. et al. Chest X-ray in new coronavirus disease 2019 (COVID-19) infection: findings and correlation with clinical outcome. *Radiol. Med.* <https://doi.org/10.1007/s11547-020-01232-9> (2020).
- Rubin, G. D. et al. The role of chest imaging in patient management during the COVID-19 pandemic: a multinational consensus statement from the Fleischner society. *Chest* **158**, 106–116 (2020).
- American College of Radiology. ACR recommendations for the use of chest radiography and computed tomography (CT) for suspected COVID-19 infection. (2020). <https://www.acr.org/Advocacy-and-Economics/ACR-Position-Statements/Recommendations-for-Chest-Radiography-and-CT-for-Suspected-COVID19-Infection>.
- Wong, H.Y.F. et al. Frequency and distribution of chest radiographic findings in COVID-19 positive patients. *Radiology*. <https://doi.org/10.1148/radiol.2020201160> (2020).
- Kundu, S., Elhalawani, H., Gichoya, J. W. & Kahn Jr, C. E. How might ai and chest imaging help unravel COVID-19's mysteries? *Radiol. Artif. Intell.* **2**, e200053 (2020).
- Khan, A. I., Shah, J. L. & Bhat, M. M. CoroNet: a deep neural network for detection and diagnosis of COVID-19 from chest X-ray images. *Comput. Meth. Prog. Bio.* **196**, 105581 (2020).
- Ucar, F. & Korkmaz, D. COVIDiagnosis-net: deep bayes-squeezeNet based diagnostic of the coronavirus disease 2019 (COVID-19) from X-ray images. *Med. Hypotheses* **140**, 109761 (2020).
- Li, L. et al. Artificial intelligence distinguishes COVID-19 from community acquired pneumonia on chest ct. *Radiology*. <https://doi.org/10.1148/radiol.2020200905> (2020).
- Ozturk, T. et al. Automated detection of COVID-19 cases using deep neural networks with X-ray images. *Comput. Biol. Med.* **121**, 103792 (2020).
- Wang, S. et al. A fully automatic deep learning system for COVID-19 diagnostic and prognostic analysis. *Eur. Respir. J.* <https://doi.org/10.1183/13993003.00775-2020> (2020).
- Zhang, K. et al. Clinically applicable AI system for accurate diagnosis, quantitative measurements, and prognosis of COVID-19 pneumonia using computed tomography. *Cell* **181**, 1423–1433.e11 (2020).
- Singh, D., Kumar, V. & Kaur, M. Classification of COVID-19 patients from chest ct images using multi-objective differential evolution-based convolutional neural networks. *Eur. J. Clin. Microbiol.* **39**, 1379–1389 (2020).
- Wynants, L. et al. Prediction models for diagnosis and prognosis of COVID-19 infection: systematic review and critical appraisal. *BMJ* **369**, m1328 (2020).
- Royal College of Physicians. National early warning score (news) 2: Standardising the assessment of acute-illness severity in the nhs. report of a working party. <https://www.rcplondon.ac.uk/projects/outputs/national-early-warning-score-news-2> (2017).
- Shamout, F. E., Zhu, T., Sharma, P., Watkinson, P. J. & Clifton, D. A. Deep interpretable early warning system for the detection of clinical deterioration. *IEEE J. Biomed. Health* **24**, 437–446 (2019).
- Li, M.D. et al. Automated assessment of COVID-19 pulmonary disease severity on chest radiographs using convolutional siamese neural networks. Preprint at <https://www.medrxiv.org/content/10.1101/2020.05.20.20108159v1> (2020).
- Borghesi, A. & Maroldi, R. COVID-19 outbreak in Italy: experimental chest X-ray scoring system for quantifying and monitoring disease progression. *Radiol. Med.* **125**, 509–513 (2020).
- Toussie, D. et al. Clinical and chest radiography features determine patient outcomes in young and middle age adults with COVID-19. *Radiology*. <https://doi.org/10.1148/radiol.2020201754> (2020).
- Fernandes, M. et al. Clinical decision support systems for triage in the emergency department using intelligent systems: a review. *Artif. Intell. Med.* **102**, 101762 (2020).
- Shen, Y. et al. Globally-aware multiple instance classifier for breast cancer screening. In *International Workshop on Machine Learning in Medical Imaging*, 18–26 (2019).
- Shen, Y. et al. An interpretable classifier for high-resolution breast cancer screening images utilizing weakly supervised localization. *Medical image analysis*. **68**, 101908 (2020).
- Ke, G. et al. Lightgbm: A highly efficient gradient boosting decision tree. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, 3146–3154 (ACM, 2017).
- Miller Jr, R.G. *Survival analysis*, vol. 66 (John Wiley & Sons, New York, 2011).
- Efron, B. & Tibshirani, R.J. *An introduction to the bootstrap* (CRC press, 1994).

31. Żoła, K., Geras, K. J. & Cho, K. Classifier-agnostic saliency map extraction. *Comput. Vis. Image Und.* **196**, 102969 (2020).
32. Baier, L., Jöhren, F. & Seebacher, S. Challenges in the deployment and operation of machine learning in practice. In *Proceedings of the 27th European Conference on Information Systems (ECIS)*, 2019).
33. Martin, A. et al. TensorFlow: Large-scale machine learning on heterogeneous distributed systems. Preprint at <https://arxiv.org/abs/1603.04467> (2015).
34. Narin, A., Kaya, C. & Pamuk, Z. Automatic detection of coronavirus disease (COVID-19) using X-ray images and deep convolutional neural networks. Preprint at <https://arxiv.org/abs/2003.10849> (2020).
35. Shamout, F.E., Zhu, T. & Clifton, D.A. Machine learning for clinical outcome prediction. *IEEE Rev. Biomed. Eng.* <https://doi.org/10.1109/RBME.2020.3007816> (2020).
36. Ahmad, M.A., Eckert, C. & Teredesai, A. Interpretable machine learning in healthcare. In *Proceedings of the 2018 ACM International Conference on Bioinformatics, Computational Biology, and Health Informatics*, 559–560 (ACM, 2018).
37. Selvaraju, R.R. et al. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE International Conference on Computer Vision*, 618–626 (IEEE, 2017).
38. Song, L. et al. Exploring the active mechanism of berberine against hcc by systematic pharmacology and experimental validation. *Mol. Med. Rep.* **20**, 4654–4664 (2019).
39. Brunese, L., Mercaldo, F., Reginelli, A. & Santone, A. Explainable deep learning for pulmonary disease and coronavirus COVID-19 detection from X-rays. *Comput. Meth. Prog. Bio.* **196**, 105608 (2020).
40. Paul, H. Y., Kim, T. K. & Lin, C. T. Generalizability of deep learning tuberculosis classifier to COVID-19 chest radiographs: new tricks for an old algorithm? *J. Thorac. Imag.* **35**, W102–W104 (2020).
41. Adebayo, J. et al. Sanity checks for saliency maps. In *NeurIPS Proceedings*, 9505–9515 (NeurIPS, 2018).
42. Brajer, N. et al. Prospective and external evaluation of a machine learning model to predict in-hospital mortality of adults at time of admission. *JAMA Netw. Open* **3**, e1920733–e1920733 (2020).
43. Lodigiani, C. et al. Venous and arterial thromboembolic complications in COVID-19 patients admitted to an academic hospital in Milan, Italy. *Thromb. Res.* **191**, 9–14 (2020).
44. Oxley, T.J. et al. Large-vessel stroke as a presenting feature of COVID-19 in the young. *New Engl. J. Med.* **382**, e60 (2020).
45. Viner, R. M. & Whittaker, E. Kawasaki-like disease: emerging complication during the COVID-19 pandemic. *Lancet* **395**, 1741–1743 (2020).
46. Dietterich, T.G. Ensemble methods in machine learning. In *International Workshop on Multiple Classifier Systems*, 1–15 (Multiple Classifier Systems, 2000).
47. He, K., Zhang, X., Ren, S. & Sun, J. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 770–778 (IEEE, 2016).
48. Huang, G., Liu, Z., Van Der Maaten, L. & Weinberger, K.Q. Densely connected convolutional networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 4700–4708 (IEEE, 2017).
49. Rajpurkar, P. et al. CheXNet: radiologist-level pneumonia detection on chest X-rays with deep learning. Preprint at <https://arxiv.org/abs/1711.05225> (2017).
50. Allauzi, I. & Ahmed, M. B. A novel approach for multi-label chest X-ray classification of common thorax diseases. *IEEE Access* **7**, 64279–64288 (2019).
51. Liu, H. et al. Sdfn: segmentation-based deep fusion network for thoracic disease classification in chest X-ray images. *Comput. Med. Imag. Grap.* **75**, 66–73 (2019).
52. Guan, Q. & Huang, Y. Multi-label chest X-ray image classification via category-wise residual attention learning. *Pattern Recogn. Lett.* **130**, 259–266 (2020).
53. Ilse, M., Tomczak, J.M. & Welling, M. Attention-based deep multiple instance learning. In *International Conference on Machine Learning*, 2127–2136 (PMLR, 2018).
54. Gensheimer, M.F. & Narasimhan, B. A scalable discrete-time survival model for neural networks. *PeerJ* **7**, e6257 (2019).
55. Cox, D.R. & Oakes, D. *Analysis of survival data*, vol. 21 (CRC Press, Boca Raton, 1984).
56. Ching, T., Zhu, X. & Garmire, L. X. Cox-nnet: an artificial neural network method for prognosis prediction of high-throughput omics data. *PLoS Comput. Biol.* **14**, e1006076 (2018).
57. Katzman, J. L. et al. DeepSurv: personalized treatment recommender system using a cox proportional hazards deep neural network. *BMC Med. Res. Methodol.* **18**, 24 (2018).
58. Liang, W. et al. Early triage of critically ill COVID-19 patients using deep learning. *Nat. Commun.* **11**, 1–7 (2020).
59. Wang, X. et al. ChestX-ray8: Hospital-scale chest X-ray database and benchmarks on weakly-supervised classification and localization of common thorax diseases. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (IEEE)*, 2017).
60. Kingma, D.P. & Ba, J. Adam: a method for stochastic optimization. In *Proceedings of the 3rd International Conference on Learning Representations (ICLR)*, 2015).
61. Bergstra, J. & Bengio, Y. Random search for hyper-parameter optimization. *J. Mach. Learn. Res.* **13**, 281–305 (2012).
62. Xu, Q. & Liang, Y. Monte Carlo cross validation. *Chemometr. Intell. Lab.* **56**, 1–11 (2001).
63. Paszke, A. et al. PyTorch: an imperative style, high-performance deep learning library. In *NeurIPS Proceedings*, 8026–8037 (NeurIPS, 2019).

## ACKNOWLEDGEMENTS

The authors would like to thank Mario Videna, Abdul Khaja and Michael Costantino for supporting our computing environment, Philip P. Rodenbough (the NYUAD Writing Center) and Catriona C. Geras for revising the manuscript, and Boyang Yu, Jimin Tan, Kyunghyun Cho and Matthew Muckley for useful discussions. We also gratefully acknowledge the support of Nvidia Corporation with the donation of some of the GPUs used in this research. This work was supported in part by grants from the National Institutes of Health (P41EB017183, R01LM013316), the National Science Foundation (HDR-1922658, HDR-1940097), and NYU Abu Dhabi.

## AUTHOR CONTRIBUTIONS

F.E.S., Y.S., N.W., A.K., J.P. and T.M. are the co-first authors of this paper. F.E.S., Y.S., N.W., A.K., J.P. and T.M. designed and conducted the experiments with neural networks. F.E.S., N.W., J.P., S.J., T.M. and J.W. built the data preprocessing pipeline. F.E.S., N.R. and B.Z. designed the clinical variables model. S.J. conducted the reader study and analyzed the data. S.D. and M.C. conducted literature search. Y.L., D.W., B.Z. and Y.A. collected the data. D.K., L.A. and W.M. analyzed the results from a clinical perspective. Y.A., C.F.G. and K.J.G. supervised the execution of all elements of the project. All authors provided critical feedback and helped shape the manuscript.

## 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/s41746-021-00453-0>.

**Correspondence** and requests for materials should be addressed to K.J.G.

**Reprints and permission information** is available at <http://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 license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license 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 license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2021