

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

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

<sup>1</sup>Engineering Division, NYU Abu Dhabi

<sup>2</sup>Center for Data Science, New York University

<sup>3</sup>Sackler Institute of Graduate Biomedical Sciences, NYU Grossman School of Medicine

<sup>4</sup>Department of Radiology, NYU Langone Health

<sup>5</sup>Center for Advanced Imaging Innovation and Research, NYU Langone Health

<sup>6</sup>Department of Population Health, NYU Langone Health

<sup>7</sup>Department of Medicine, NYU Langone Health

<sup>8</sup>Department of Mathematics, Courant Institute, New York University

\*Equal contribution

✉k.j.geras@nyu.edu

## Abstract

During the 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 3,661 patients, achieves an AUC of 0.786 (95% CI: 0.742-0.827) 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 NYU 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.

## 1 Introduction

In recent months, there has been a surge in patients presenting to the emergency department (ED) with respiratory illnesses associated with SARS CoV-2 infection (COVID-19) [1, 2]. Evaluating the risk of deterioration of these patients to perform triage is crucial for clinical decision-making and resource allocation [3]. While ED triage is difficult under normal circumstances [4, 5], during a pandemic, strained hospital resources increase the challenge [2, 6]. 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.

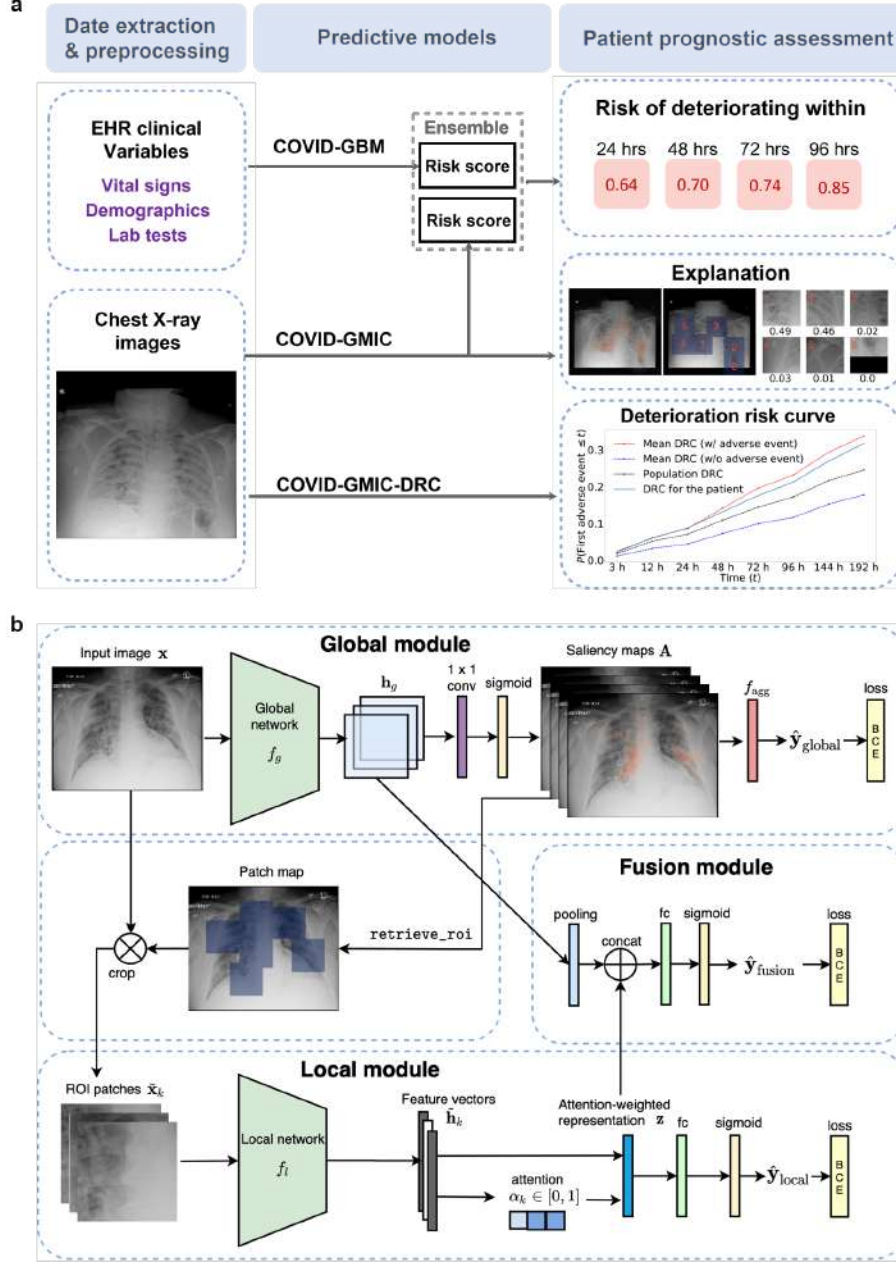
As the primary complication of COVID-19 is pulmonary disease, such as pneumonia [7], chest X-ray imaging is a first-line triage tool for COVID-19 patients. Although other imaging modalities, such as computer tomography (CT), provide higher resolution, chest X-ray images are less costly, inflict a lower radiation dose, and are easier to obtain without incurring the risk of contaminating imaging equipment and disrupting radiologic services [8]. In addition, abnormalities in the chest X-ray images of COVID-19 patients have been found to mirror abnormalities in CT scans [9]. Consequently, chest X-ray imaging is considered a key tool in assessing COVID-19 patients [10]. Unfortunately, although the knowledge of the disease is rapidly evolving, understanding of the correlation between pulmonary parenchymal patterns visible in the chest X-ray images and clinical deterioration is 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 works using imaging data of COVID-19 patients concentrate more on diagnosis than prognosis [11, 12, 13, 14, 15, 16, 17, 18]. Prognostic models 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 [11]. Prior methodology for prognosis of COVID-19 patients via machine learning mainly use routinely-collected clinical variables [2, 19] such as vital signs and laboratory tests, which have long been established as strong predictors of deterioration [20, 21]. Some studies have proposed scoring systems for chest X-ray images to assess the severity and progression of lung involvement using deep learning [22], or more commonly, through manual clinical evaluation [7, 23, 24]. 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.

In this work, we present 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. The goal is to provide support for critical clinical decision-making involving patients arriving at the ED in need of immediate care [2, 25]. We designed our system to satisfy a clinical need of frontline physicians. We were able to build it due to the availability of a large-scale chest X-ray image dataset. The system is based on chest X-ray imaging, which is already being employed as a first-line triage tool in hospitals [7], while also incorporating other routinely collected non-imaging clinical variables that are known to be strong predictors of deterioration.

Our AI system uses deep convolutional neural networks to perform risk evaluation from chest X-ray images. In particular, we base our work on the Globally-Aware Multiple Instance Classifier (GMIC) [26, 27], which is designed to provide interpretability by highlighting the most informative regions of the input images. We call this imaging-based model COVID-GMIC. The system also learns from routinely collected clinical variables using a gradient boosting model (GBM) [28] which we call COVID-GBM. Both models are trained using a dataset of 3,661 patients admitted to NYU Langone Health between March 3, 2020, and May 13, 2020. The outputs of COVID-GMIC and COVID-GBM are combined to predict the risk of deterioration of individual patients over different time horizons, ranging from 24 to 96 hours. In addition, our system includes a model, which we call COVID-GMIC-DRC, that predicts how the risk of deterioration is expected to evolve over time, in the spirit of survival analysis [29].

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: 0.742-0.827), and an area under the precision recall curve (PR AUC) of 0.517 (95% CI: 0.434, 0.605) for prediction of deterioration within 96 hours. 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.



**Figure 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 [26, 27]. The first model, COVID-GMIC, predicts the overall risk of deterioration within 24, 48, 72, and 96 hours, 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 [28] 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 hours 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.

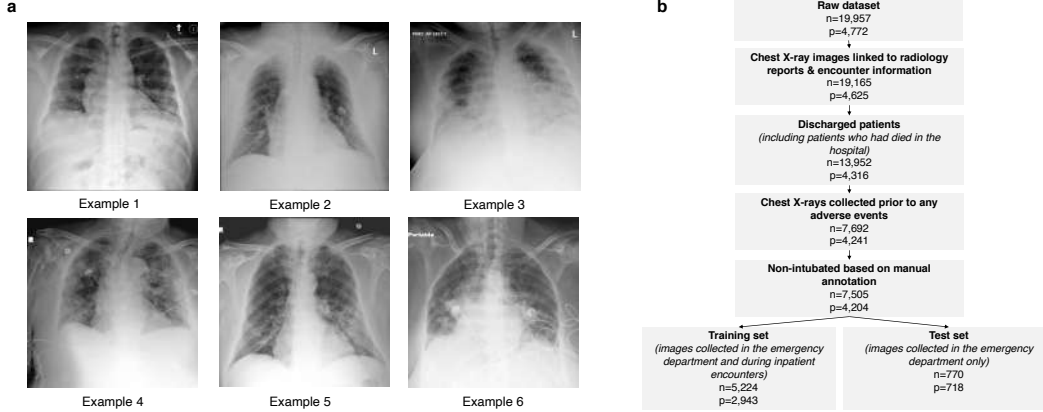


Figure 2: **Illustrations of the dataset.** **a**, Examples of chest X-ray images in our dataset. Example 1: Patient was discharged and experienced no adverse events (44 years old male). Example 2: Patient was transferred to the ICU after 95 hours (71 years old male). Example 3: Patient was intubated after 72 hours (66 years old male). Example 4: Patient was transferred to the ICU after 48 hours (99 years old female). Example 5: Patient was intubated after 24 hours (74 years old male). Example 6: Patient was transferred to the ICU in 30 minutes (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. We excluded chest X-ray images that had missing radiology reports or patient encounter information to ensure data completeness, as well as chest X-ray images that were collected after a patient had experienced an adverse event, since deterioration had already occurred. We included patients who were discharged, and patients who had experienced in-hospital mortality, in order to obtain a full record of adverse events. We also manually checked for images of already intubated patients, and excluded them. In the test set, we only included images collected in the ED, and excluded images collected during inpatient encounters.

## 2 Results

**Dataset.** Our AI system was developed and evaluated using a dataset collected at NYU Langone Health between March 3, 2020 and June 28, 2020.<sup>1</sup> 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). The training set consisting of 5,617 chest X-ray images was used for model development and hyperparameter tuning, while the test set consisting of 832 images was used to report the final results. 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. Supplementary Table 1 summarizes the associated clinical variables included in the dataset.

We define deterioration, the target to be predicted by our models, as the occurrence 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. Figure 2.a shows examples of chest X-ray images collected from different patients. Although the patient in example 5 had less severe parenchymal findings than patients in examples 3 and 4, the patient was intubated within 24 hours compared to 48 and 96 hours in examples 3 and 4. This highlights the difficulty of assessing the risk of deterioration using only chest X-ray images, since the extent of visible parenchymal disease is not fully predictive of the time of deterioration.

**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 hours from the time of the

<sup>1</sup>This study was approved by the Institutional Review Board, with ID# i20-00858.

Table 1: Description of the characteristics of the patient cohort included in the training and test sets used to develop and evaluate our system. The training and test sets are similar in terms of age, BMI, and proportion of females. We note that there is a higher proportion of chest X-ray images associated with deterioration across all time windows in the test set compared to the training set. This implies that there is a higher incidence of adverse events amongst ED patients than inpatients, since the test set only includes chest X-ray images collected from ED patients, while the training set also includes inpatients.

Characteristic	Training set	Test set
Patients, n	2,943	718
Admissions, n	3,175	764
Females, n (%)	1,206 (41.0)	305 (42.5)
Age (years), mean (SD)	62.9 (17.2)	64.9 (17.2)
BMI (kg/m <sup>2</sup> ), mean (SD)	29.4 (7.0)	29.5 (8.6)
Survived	2405	559
Adverse events, n	1,311	594
Intubation, n	386	97
ICU admission, n	387	113
Mortality, n	538	159
Composite outcome, n	730	225
Chest X-ray exams, n	5,224	770
Composite outcome within 24 hours, n (%)	349 (6.7%)	74 (9.6%)
Composite outcome within 48 hours, n (%)	553 (10.6%)	101 (13.1%)
Composite outcome within 72 hours, n (%)	735 (14.1%)	130 (16.9%)
Composite outcome within 96 hours, n (%)	876 (16.8%)	156 (20.3%)
Total number of images, n	5,617	832

chest X-ray exam. The receiver operating characteristic curves and precision-recall curves can be found in Supplementary Figure 4. Our ensemble model consisting of COVID-GMIC and COVID-GBM achieves the best AUC performance across all time windows compared to COVID-GMIC and COVID-GBM individually. 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 Figure 2.b. Similarly, the ensemble of COVID-GMIC and COVID-GBM outperforms all models across all time windows in terms of the PR AUC, except for the 96 hours window.

To illustrate the interpretability of COVID-GMIC, we show in Figure 3 the saliency maps for all time windows (24, 48, 72, and 96 hours) computed for four examples from the test set. Across all four examples, the saliency maps highlight regions that contain visual patterns such as airspace opacities and consolidation, which are correlated with clinical deterioration [22, 24]. 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, while missing the dense consolidation in the periphery of the right upper lobe. This suggests that COVID-GMIC emphasizes only the most informative regions in the image, while human experts can provide a more holistic interpretation covering the entire image. It might, therefore, be useful to enhance GMIC through a classifier agnostic mechanism [31], which finds all the useful evidence in the image, instead of solely the most discriminative part. We leave this for future work.

**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. 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

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. We include 95% confidence intervals estimated by 1,000 iterations of the bootstrap method [30]. 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 Figure 2.b. 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 hours task. In the reader study, our main finding is that COVID-GMIC outperforms radiologists A & B across time windows longer than 24 hours, 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.

Test set (n=832)								
	AUC				PR AUC			
	24 hours	48 hours	72 hours	96 hours	24 hours	48 hours	72 hours	96 hours
COVID-GBM	0.747 (0.692, 0.796)	0.739 (0.683, 0.788)	0.750 (0.701, 0.797)	0.770 (0.727, 0.813)	0.230 (0.164, 0.321)	0.325 (0.254, 0.421)	0.408 (0.337, 0.499)	<b>0.523</b> (0.446, 0.613)
COVID-GMIC	0.695 (0.627, 0.754)	0.716 (0.661, 0.766)	0.717 (0.661, 0.766)	0.738 (0.691, 0.781)	0.200 (0.140, 0.281)	0.302 (0.225, 0.395)	0.374 (0.296, 0.465)	0.439 (0.363, 0.532)
COVID-GBM + COVID-GMIC	<b>0.765</b> (0.713, 0.818)	<b>0.749</b> (0.700, 0.798)	<b>0.769</b> (0.720, 0.814)	<b>0.786</b> (0.742, 0.827)	<b>0.243</b> (0.187, 0.336)	<b>0.332</b> (0.254, 0.427)	<b>0.439</b> (0.351, 0.533)	0.517 (0.434, 0.605)
Reader study dataset (n=200)								
	AUC				PR AUC			
	24 hours	48 hours	72 hours	96 hours	24 hours	48 hours	72 hours	96 hours
Radiologist A	0.613 (0.521, 0.707)	0.645 (0.559, 0.719)	0.691 (0.612, 0.764)	0.740 (0.666, 0.806)	0.346 (0.251, 0.475)	0.490 (0.381, 0.613)	0.640 (0.535, 0.744)	0.742 (0.650, 0.827)
Radiologist B	0.637 (0.544, 0.727)	0.636 (0.556, 0.720)	0.658 (0.578, 0.728)	0.713 (0.640, 0.777)	0.365 (0.268, 0.501)	0.460 (0.360, 0.585)	0.590 (0.479, 0.688)	0.704 (0.603, 0.792)
Radiologist A + Radiologist B	<b>0.642</b> (0.555, 0.729)	0.663 (0.580, 0.737)	0.692 (0.618, 0.763)	0.741 (0.673, 0.804)	<b>0.403</b> (0.286, 0.534)	0.499 (0.385, 0.618)	0.609 (0.507, 0.726)	0.740 (0.649, 0.830)
COVID-GMIC	<b>0.642</b> (0.550, 0.730)	<b>0.701</b> (0.621, 0.775)	<b>0.751</b> (0.681, 0.817)	<b>0.808</b> (0.746, 0.866)	0.381 (0.282, 0.527)	<b>0.546</b> (0.435, 0.671)	<b>0.676</b> (0.572, 0.788)	<b>0.789</b> (0.698, 0.879)
COVID-GBM	0.704 (0.624, 0.776)	0.719 (0.644, 0.790)	0.750 (0.679, 0.816)	0.787 (0.724, 0.847)	0.411 (0.304, 0.563)	0.537 (0.434, 0.680)	0.668 (0.566, 0.778)	0.804 (0.724, 0.870)
COVID-GBM + COVID-GMIC	0.708 (0.617, 0.779)	0.702 (0.629, 0.771)	0.778 (0.705, 0.837)	0.819 (0.753, 0.875)	0.411 (0.305, 0.543)	0.500 (0.399, 0.636)	0.705 (0.604, 0.811)	0.808 (0.718, 0.881)

48, 72, and 96 hours. 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 hours 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 [24]. Supplementary Figure 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 (DRCs) which predict the evaluation of the deterioration risk based on chest X-ray images. Figure 4.a 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 hours. 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 hours, 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 5.

Figure 4.b 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

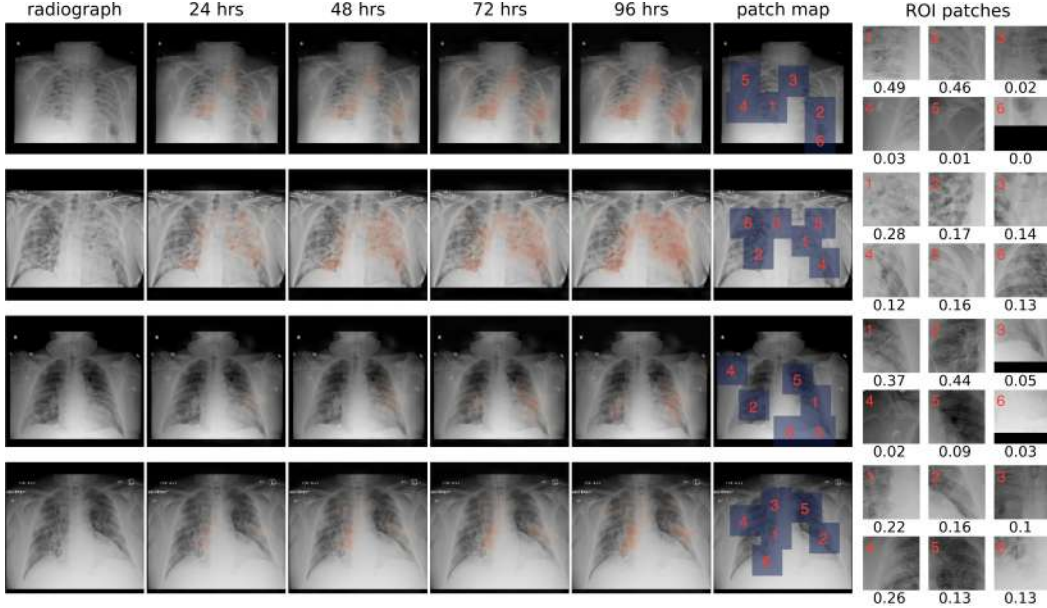


Figure 3: **Explainability of COVID-GMIC.** From left to right: the original X-ray image, saliency maps for clinical deterioration within 24, 48, 72, and 96 hours, locations of region-of-interest (ROI) patches, and ROI patches with their associated attention scores. All four patients were admitted to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

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 [32]. 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 hours using a subset of our data collected prior to deployment from 3,425 patients. The models were serialized and served with TensorFlow Serving components [33] 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 approximately two seconds to extract an image from the DICOM receiver (C-STORE), apply the image preprocessing steps, and get the prediction of a model as a Tensorflow [33] output.

Of the 375 exams collected between May 22, 2020 and June 24, 2020, 38 exams were associated with a positive 96 hour deterioration outcome. 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



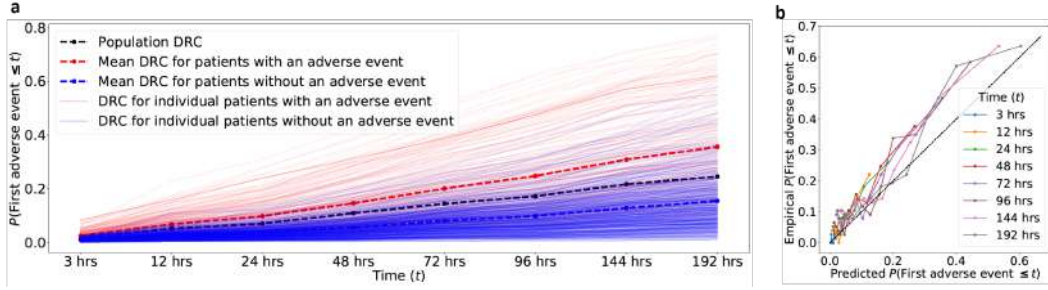


Figure 4: **Deterioration risk curves (DRCs) and reliability plot for COVID-GMIC-DRC.** **a.** DRCs generated by the COVID-GMIC-DRC model for patients in the test set with (faded red lines) and without adverse events (faded blue lines). The mean DRC for patients with adverse events (red dashed line) 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. The diagram shows that these values are similar across the different values of  $t$ , and hence the model’s probability predictions are well-calibrated (for comparison, perfect calibration would correspond to the diagonal black dashed line).

0.365 (95% CI: 0.313-0.465) PR AUC. The decrease in accuracy may indicate changes in the patient population as the pandemic progressed.

### 3 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 [12, 13, 15, 19, 34]. 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 [19, 35]. In Supplementary Table 4, 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 [36]. Existing works on COVID-19 often use external gradient-based algorithms, such as gradCAM [37], to interpret deep neural network classifiers [38, 39, 40]. However, visualizations generated by gradient-based methods are sensitive to minor perturbation in input images, and could yield misleading interpretations [41]. In contrast, COVID-GMIC has an inherently interpretable architecture that better retains localization information of the more informative regions in the input images.



We performed prospective validation of an early version of our system through silent deployment in a hospital which uses the Epic electronic health record system. 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 [42]. To the best of our knowledge, our study is the first to do so for the prognosis of COVID-19 patients.

Our approach has some limitations that will be addressed in future work. 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 performance of this model dropped from an AUC of 0.748 (95% CI: 0.708- 0.790) during retrospective evaluation to 0.717 (95% CI: 0.622-0.801) during prospective validation, suggesting that the model may need to be fine-tuned as additional data is collected. In addition, 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. Incorporating additional data from patient health records may 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 [25]. Although we focus on chest X-ray images because pulmonary disease is the main complication associated to 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 [43, 44, 45]. 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.

## 4 Methods

**Outline.** In this section, we first introduce our data collection and preprocessing pipeline. We 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.

**Dataset collection and preparation.** We extracted a dataset of 19,957 chest X-ray exams collected from 4,772 patients who tested positive for COVID-19 between March 2, 2020, and May 13, 2020. We applied inclusion and exclusion criteria that were defined in collaboration with clinical experts, as shown in Figure 2.b. Specifically, we excluded 783 exams that were not linked to any radiology report, nine exams that were not linked to any encounter information, and 5,213 exams from patients who were still hospitalised by May 13, 2020. To ensure that our system predicts deterioration prior to its occurrence, we excluded 6,260 exams that were collected after an adverse event and 187 exams of already intubated patients. The final dataset consists of 7,502 chest X-ray exams corresponding to 4,204 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. This resulted in 5,224 exams (5,617

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.

The data used to evaluate the models during deployment consist of 375 exams from 217 patients collected between May 22, 2020 and June 24, 2020. The exams were filtered based on the same criteria described above. Among the 375 exams, 25 chest X-ray exams were collected from patients who were admitted to the ICU within 96 hours, and six exams were collected from patients who were intubated within 96 hours.

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 to 45 degrees), and random translation. Supplementary Figure 1 shows the distribution of the size of the images prior to data augmentation, as well as examples of images before and after preprocessing.

In addition to the chest X-ray images, we extracted clinical variables for patients including patient demographics (age, weight, and body mass index), vital signs (heart rate, systolic blood pressure, diastolic blood pressure, temperature, and respiratory rate), and 25 lab test variables listed in Supplementary Table 1. All vital signs were collected prior to the chest X-ray exam.

**Adverse event prediction.** Our main goal is to predict clinical deterioration within four time windows of 24, 48, 72, and 96 hours. We frame this as a 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 [46]. 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 Figure 2.b, 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 laboratory test results were represented as maximum and minimum statistics of all values collected within 12 hours prior to the time of the vital sign measurement. The feature sets of vital signs and laboratory tests were then processed using a gradient boosting model [28] which we refer to as COVID-GBM. 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, 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 [26, 27]. 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 [47, 48], which may distort and blur informative visual patterns in chest X-ray images such as basilar opacities and pulmonary consolidation. In Supplementary Table 2, 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 [49, 50, 51, 52].

The architecture of COVID-GMIC is schematically depicted in Figure 1.b. 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.<sup>2</sup> Each element  $\mathbf{A}_{i,j}^t \in [0, 1]$  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  $f_{\text{agg}} : \mathbb{R}^{h,w} \mapsto [0, 1]$  to transform all saliency maps  $\mathbf{A}^t$  for all time windows  $t$  into classification predictions  $\hat{\mathbf{y}}_{\text{global}}$ :

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

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 Figure 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. 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 [47], which produces a feature vector  $\tilde{\mathbf{h}}_k \in \mathbb{R}^{512}$  from each ROI patch. The predictive value of each ROI patch might vary significantly. Therefore, we utilize the gated attention mechanism [53] to compute an attention score  $\alpha_k \in [0, 1]$  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:

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

The representation  $\mathbf{z}$  is then passed into a fully connected layer with sigmoid activation to generate a prediction  $\hat{\mathbf{y}}_{\text{local}}$ . We refer the readers to Shen et al. [27] 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  $\hat{\mathbf{y}}_{\text{fusion}}$ .

In our experiments, we chose  $H = W = 1024$ . Supplementary Table 2 shows that COVID-GMIC achieves the best validation performance for this resolution. We parameterize  $f_g$  as a ResNet-18 [47] 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:

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

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  $\hat{\mathbf{y}}_{\text{fusion}}$  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

<sup>2</sup>For visualization purposes, we apply nearest neighbor interpolation to upsample the saliency maps to match the resolution of the original image.

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 [54], we parameterize the DRC using a vector of conditional probabilities  $\hat{\mathbf{p}} \in \mathbb{R}^8$ . The  $i^{th}$  entry of this vector,  $\hat{\mathbf{p}}_i$ , 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:<sup>3</sup>

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

Given an estimate of  $\hat{\mathbf{p}}$ , the DRC can be computed applying the chain rule:

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

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  $\hat{\mathbf{p}}_{\text{global}}$ ,  $\hat{\mathbf{p}}_{\text{local}}$ , and  $\hat{\mathbf{p}}_{\text{fusion}}$ . During inference, we use the output of the fusion module,  $\hat{\mathbf{p}}_{\text{fusion}}$ , 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 [47]. 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 Tables 2 and 5. During training, we minimize the following loss function defined on a single example:

$$l(T, \hat{\mathbf{p}}_{\text{global}}, \hat{\mathbf{p}}_{\text{local}}, \hat{\mathbf{p}}_{\text{fusion}}) = l_s(T, \hat{\mathbf{p}}_{\text{global}}) + l_s(T, \hat{\mathbf{p}}_{\text{local}}) + l_s(T, \hat{\mathbf{p}}_{\text{fusion}}) + \sum_{m=0}^8 \beta |\mathbf{A}^m|, \quad (7)$$

where  $l_s$  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

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

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

---

<sup>3</sup>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 hours.

information after  $t_i$  and no adverse events before  $t_i$ , equals

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

Note that each  $\hat{\mathbf{p}}_i$  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 [55] require using the whole dataset to perform model updates [56, 57, 58], 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-GBM. The chest X-ray image models were implemented in PyTorch [59] and trained using NVIDIA Tesla V100 GPUs. The clinical variables models were implemented using the Python library LightGBM [28].

We initialized the weights of COVID-GMIC and COVID-GMIC-DRC by pretraining them on the ChestX-ray14 dataset [60] (Supplementary Table 3 compares the performance of different initialization strategies). We used Adam [61] with a minibatch size of eight to train the models on our data. We applied data augmentation during training and testing, but not during validation. During testing, we augmented each image ten times and averaged the corresponding outputs to produce the final prediction.

We optimized the hyperparameters using random search [62]. 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 [63] (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.

**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 hours). 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 hours. 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.

## Acknowledgements

The authors would like to thank Mario Videna, Abdul Khaja and Michael Constantino 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) and the National Science Foundation (HDR-1922658, HDR-1940097).

## Author Contributions

FES, YS, NW, AK, JP and TM designed and conducted the experiments with neural networks. FES, NW, JP, SJ and TM built the data preprocessing pipeline. FES, NR and BZ designed the clinical variables model. SJ conducted the reader study and analyzed the data. SD and MC conducted literature search. YL, DW, BZ and YA collected the data. DK, LA and WM analyzed the results from a clinical perspective. YA, CFG and KJG 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.

## References

- [1] Baugh, J. J. *et al.* Creating a COVID-19 surge clinic to offload the emergency department. *Am. J. Emerg. Med.* **38**, 1535–1537 (2020).
- [2] Debnath, S. *et al.* Machine learning to assist clinical decision-making during the COVID-19 pandemic. *Bioelectron. Med.* **6**, 1–8 (2020).
- [3] 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).
- [4] Dorsett, M. Point of no return: COVID-19 and the us health care system: An emergency physician’s perspective. *Sci. Adv.* **6** (2020).
- [5] McKenna, P. *et al.* Emergency department and hospital crowding: causes, consequences, and cures. *Clin. Exp. Emerg. Med.* **6**, 189 (2019).
- [6] Warner, M. A. Stop doing needless things! Saving healthcare resources during COVID-19 and beyond. *J. Gen. Intern. Med.* **35**, 2186–2188 (2020).
- [7] 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).
- [8] American College of Radiology. ACR recommendations for the use of chest radiography and computed tomography (CT) for suspected COVID-19 infection. <https://www.acr.org/Advocacy-and-Economics/ACR-Position-Statements/Recommendations-for-Chest-Radiography-and-CT-for-Suspected-COVID19-Infection> (2020).
- [9] 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).
- [10] 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).
- [11] 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).
- [12] 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).
- [13] 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).

- [14] 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).
- [15] Ozturk, T. *et al.* Automated detection of COVID-19 cases using deep neural networks with X-ray images. *Comput. Biol. Med.* **121**, 103792 (2020).
- [16] 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).
- [17] 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).
- [18] 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).
- [19] Wynants, L. *et al.* Prediction models for diagnosis and prognosis of COVID-19 infection: systematic review and critical appraisal. *BMJ* **369**, m1328 (2020).
- [20] 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).
- [21] 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).
- [22] 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).
- [23] 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).
- [24] 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).
- [25] 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).
- [26] 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).
- [27] Shen, Y. *et al.* An interpretable classifier for high-resolution breast cancer screening images utilizing weakly supervised localization. Preprint at <https://arxiv.org/abs/2002.07613> (2020).
- [28] Ke, G. *et al.* Lightgbm: A highly efficient gradient boosting decision tree. In *Adv. Neur. In.*, 3146–3154 (2017).
- [29] Miller Jr, R. G. *Survival Analysis*, vol. 66 (John Wiley & Sons, New York, 2011).
- [30] Efron, B. & Tibshirani, R. J. *An introduction to the bootstrap* (CRC press, 1994).
- [31] Żołna, 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] Martín, 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).



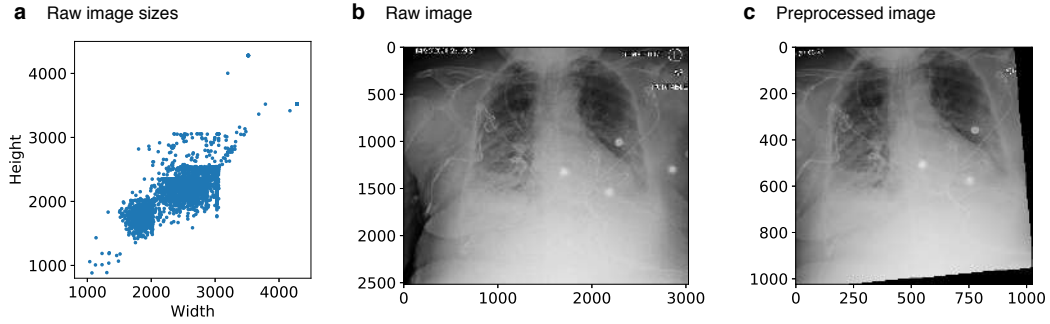
- [36] Ahmad, M. A., Eckert, C. & Teredesai, A. Interpretable machine learning in healthcare. In *Proceedings of the 2018 ACM International Conference on Bioinformatics, Computational Biology, and Health Informatics*, 559–560 (2018).
- [37] Selvaraju, R. R. *et al.* Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE International Conference on Computer Vision*, 618–626 (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 *Adv. Neur. In.*, 9505–9515 (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.* (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** (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 (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 (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 (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] Allaouzi, 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. Preprint at <https://arxiv.org/abs/1802.04712> (2018).
- [54] Gensheimer, M. F. & Narasimhan, B. A scalable discrete-time survival model for neural networks. Preprint at <https://arxiv.org/abs/1805.00917> (2018).
- [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] Paszke, A. *et al.* PyTorch: An imperative style, high-performance deep learning library. In *Adv. Neur. In.*, 8026–8037 (2019).
- [60] 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* (2017).
- [61] Kingma, D. P. & Ba, J. Adam: A method for stochastic optimization. In *Proceedings of the 3rd International Conference on Learning Representations* (2015).
- [62] Bergstra, J. & Bengio, Y. Random search for hyper-parameter optimization. *J. Mach. Learn. Res.* **13**, 281–305 (2012).
- [63] Xu, Q. & Liang, Y. Monte Carlo cross validation. *Chemometr. Intell. Lab.* **56**, 1–11 (2001).
- [64] Liu, K., Chen, Y., Lin, R. & Han, K. Clinical features of COVID-19 in elderly patients: A comparison with young and middle-aged patients. *J. Infection.* **80**, e14–e18 (2020).
- [65] Krizhevsky, A. Learning multiple layers of features from tiny images. Tech. Rep., University of Toronto (2009).
- [66] Deng, J. *et al.* ImageNet: A Large-Scale Hierarchical Image Database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, 248–255 (2009).
- [67] Geras, K. J. *et al.* High-resolution breast cancer screening with multi-view deep convolutional neural networks. Preprint at <https://arxiv.org/abs/1703.07047> (2017).
- [68] Pan, S. J. & Yang, Q. A survey on transfer learning. *IEEE Trans. Knowl. Data Eng.* **22**, 1345–1359 (2009).
- [69] Yosinski, J., Clune, J., Bengio, Y. & Lipson, H. How transferable are features in deep neural networks? In *Adv. Neur. In.*, 3320–3328 (2014).
- [70] He, K., Zhang, X., Ren, S. & Sun, J. Delving deep into rectifiers: surpassing human-level performance on ImageNet classification. In *Proceedings of the IEEE International Conference on Computer Vision*, 1026–1034 (2015).

## Supplementary Information

### Supplementary Note 1: Image preprocessing

In Supplementary Figure 1.a, we show the heights and widths of the images prior to data augmentation. In Supplementary Figure 1.b, we show an example of a raw image and the final image after applying the preprocessing steps in Figure 1.c.



Supplementary Figure 1: (a) Heights and widths (in pixels) of images prior to data augmentation. (b) An example raw image. (c) To ensure that the inputs to the model have a consistent size, we perform center cropping and rescaling. In addition, we apply random horizontal flipping, rotation, and translation to augment the training dataset.

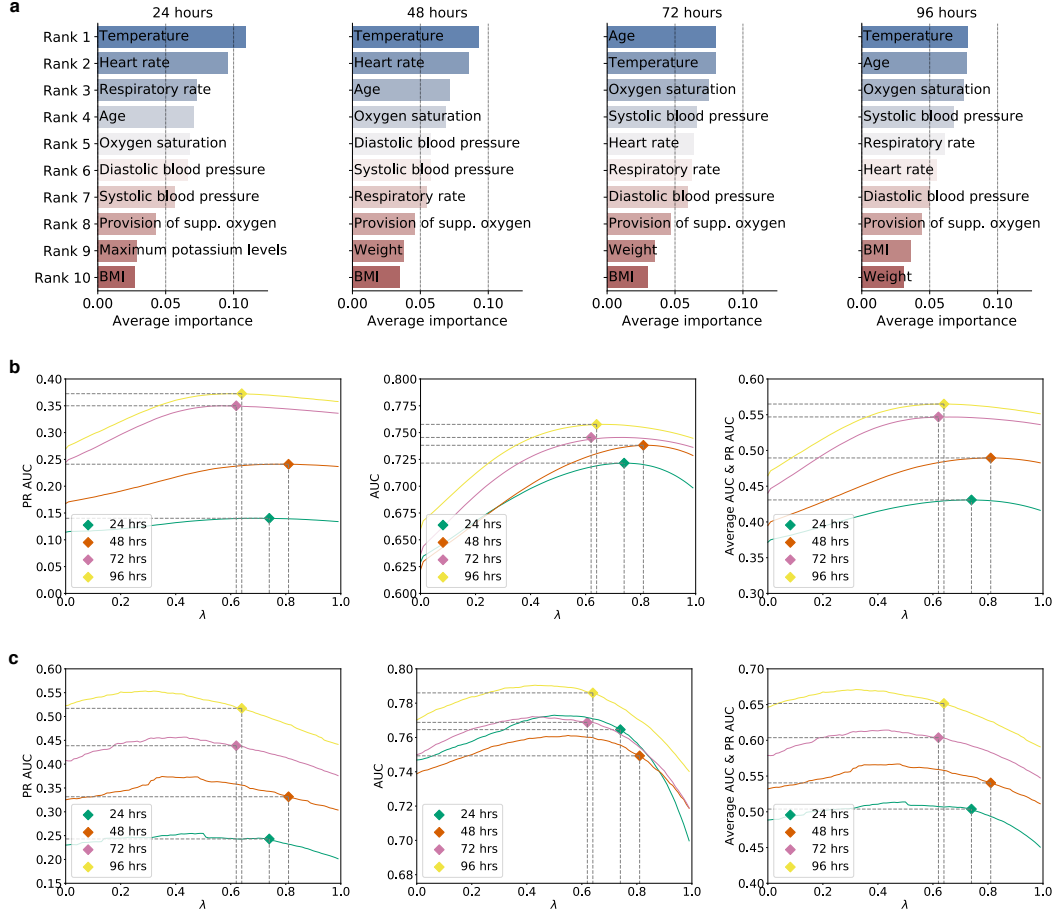
## Supplementary Note 2: Clinical variables modeling

The statistics of the clinical variables that were used to develop the COVID-GBM models are listed in Table 1. The raw laboratory test variables were further processed to extract the minimum and maximum statistics.

Supplementary Table 1: Mean and interquartile range statistics of the raw vital signs and laboratory test results, corresponding to the patients included in the training and test sets for COVID-GBM. Note that  $n$  represents a counting unit.

Variable, <i>unit</i>	Training Set	Test Set
Vital signs		
Heart rate, <i>beats per minute</i>	93.7 (25.0)	93.5 (27.0)
Respiratory rate, <i>breaths per minute</i>	22.4 (7.0)	23.4 (7.0)
Temperature, $^{\circ}F$	99.4 (1.9)	99.4 (1.9)
Systolic blood pressure, <i>mmHg</i>	130.7 (30.0)	129.8 (29.3)
Diastolic blood pressure, <i>mmHg</i>	75.9 (17.0)	76.0 (18.0)
Oxygen saturation, %	94.1 (4.0)	93.8 (5.0)
Laboratory tests		
Albumin, <i>g/dL</i>	3.5 (0.9)	3.5 (0.9)
ALT, <i>U/L</i>	49.8 (32.0)	52.2 (36.0)
AST, <i>U/L</i>	67.3 (37.0)	69.7 (43.0)
Total bilirubin, <i>mg/dL</i>	0.7 (0.4)	0.7 (0.4)
Blood urea nitrogen, <i>mg/dL</i>	25.9 (17.0)	26.4 (18.0)
Calcium, <i>mg/dL</i>	8.7 (0.8)	8.7 (0.8)
Chloride, <i>mEq/L</i>	101.1 (7.0)	101.6 (7.0)
Creatinine, <i>mg/dL</i>	1.6 (0.7)	1.6 (0.7)
D-dimer, <i>ng/mL</i>	1,321.6 (535.5)	1,146.3 (618.5)
Eosinophils, %	0.4 (0.0)	0.4 (0.0)
Eosinophils, $n$	0.03 (0.00)	0.03 (0.00)
Hematocrit, %	38.9 (7.3)	38.9 (7.5)
LDH, <i>U/L</i>	412.8 (207.0)	404.0 (213.0)
Lymphocytes, %	14.1 (10.0)	14.9 (11.0)
Lymphocytes, $n$	1.0 (0.7)	1.0 (0.7)
Platelet volume, <i>fL</i>	10.6 (1.4)	10.6 (1.4)
Neutrophils, $n$	6.4 (4.0)	6.3 (3.8)
Neutrophils, %	76.6 (14.0)	75.9 (13.0)
Platelet, $n$	226.1 (114.0)	223.7 (103.0)
Potassium, <i>mmol/L</i>	4.2 (0.8)	4.2 (0.8)
Procalcitonin, <i>ng/mL</i>	1.9 (0.3)	1.9 (0.4)
Total protein, <i>g/dL</i>	7.1 (1.1)	7.2 (1.0)
Sodium, <i>mmol/L</i>	136.2 (6.0)	136.6 (7.0)
Troponin, <i>ng/mL</i>	0.2 (0.1)	0.2 (0.1)

The average importance of the top ten features computed by the COVID-GBM models are shown in Supplementary Figure 2.a. The importance of a feature is measured by the numbers of times the feature is used to split the data across all trees in a single COVID-GBM model. Age is amongst the top ten features across all time windows, which is consistent with existing findings that mortality is more common amongst elderly COVID-19 patients than younger patients [64]. The inclusion of the vital sign variables, amongst the top ten features across all models, is also aligned with existing research suggesting that they are strong indicators of deterioration [20].



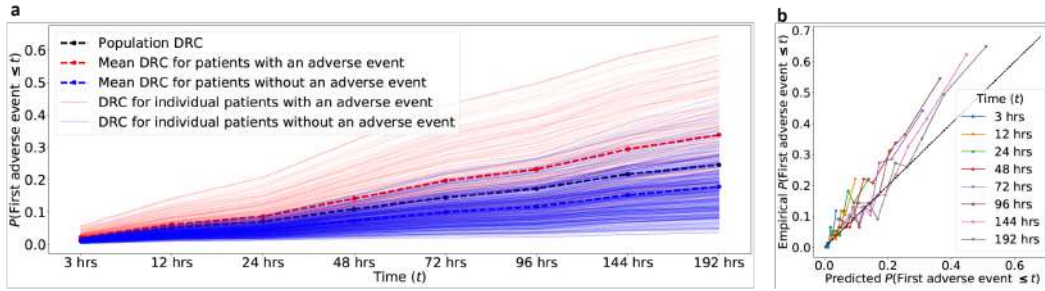
Supplementary Figure 2: **Additional results for COVID-GBM and the ensemble of COVID-GBM and COVID-GMIC.** **a**, The average importance of the top ten features computed by the nine COVID-GBM ensemble models for 24, 48, 72, and 96 hours. The importance of a feature is measured by the numbers of times the feature is used to split the data across all trees in a model. **b**, The effect of varying  $\lambda$ , the weight on the COVID-GMIC prediction, in combining the predictions of COVID-GMIC and COVID-GBM when using AUC, PR AUC and the average AUC and PR AUC on the validation set. For the average AUC and PR AUC, the optimal  $\lambda$  was 0.74 for 24 hours, 0.81 for 48 hours, 0.62 for 72 hours, and 0.64 for 96 hours. **c**, the optimal values of  $\lambda$  selected through the validation set in **b** are shown for the test set.

### Supplementary Note 3: Ablation studies

**DenseNet-121-based models.** DenseNet [48] is a deep neural network architecture which consists of dense blocks in which layers are directly connected to every other layer in a feed-forward fashion. It achieves strong performance on benchmark natural images dataset, such as CIFAR10/100 [65] and ILSVRC 2012 (ImageNet) dataset [66] while being computationally efficient. Here we compare COVID-GMIC to a specific variant of DenseNet, DenseNet-121, which has been applied to process chest X-ray images in the literature [49, 50, 51, 52].

The model assumes an input size of  $224 \times 224$ . We applied DenseNet-121-based models to predict deterioration and also to compute deterioration risk curves. We initialized the models with weights pretrained on the ChestX-ray14 dataset [60], provided at <https://github.com/arnoweng/CheXNet>. We used weight decay in the optimizer. To perform hyperparameter search, we sampled the learning rate and the rate of weight decay per step uniformly on a logarithmic scale between  $10^{[-6, -1]}$  and  $10^{[-6, -3]}$ .

For adverse event prediction, the DenseNet-121 based model yielded test AUCs of 0.687 (95% CI: 0.621 - 0.749), 0.709 (95% CI: 0.653 - 0.757), 0.710 (95% CI: 0.660 - 0.763), and 0.736 (95% CI: 0.691 - 0.782), and PRAUCs of 0.216 (95% CI: 0.155 - 0.317), 0.315 (95% CI: 0.239 - 0.419), 0.373 (95% CI: 0.300 - 0.464), and 0.454 (95% CI: 0.384 - 0.542) for 24, 48, 72, and 96 hours. The deterioration risk curves produced by the DenseNet-121 based models and the corresponding reliability plot are presented in Figure 3.



Supplementary Figure 3: **Deterioration risk curves (DRCs) and reliability plot for DenseNet-121.** Compare to Figure 4, which shows analogous graphs for COVID-GMIC-DRC. **a**, DRCs generated by DenseNet-121 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 DenseNet-121 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. The diagram shows that these values are similar across the different values of  $t$ , and hence the model is well-calibrated (for comparison, perfect calibration would correspond to the diagonal black dashed line).

**Impact of input image resolution.** Prior work on deep learning for medical images [67] report that using high resolution input images can improve performance. In this section, we analyze the impact of image resolution on our tasks of interest. We consider the following image sizes:  $128 \times 128$ ,  $256 \times 256$ ,  $512 \times 512$ , and  $1024 \times 1024$ . We pretrain all models on the ChestX-ray14 dataset [60] and then fine-tune them on our dataset. Results on the test set are reported in Supplementary Table 2.

The DenseNet-121 based model achieves the best AUCs when using an image size of  $256 \times 256$ , and the best concordance index for  $512 \times 512$ . Further increasing the resolution does not improve performance. COVID-GMIC achieves the best performance for the highest input image resolution of  $1024 \times 1024$ , while achieving the best concordance index for  $512 \times 512$ . While a further increase in performance may be possible, we did not consider any larger image sizes resolutions because the computational cost would become prohibitively high.

Supplementary Table 2: Model performance with 95% confidence intervals when using input images of sizes of  $128 \times 128$ ,  $256 \times 256$ ,  $512 \times 512$ , and  $1024 \times 1024$ . For COVID-GMIC, we started with a size of  $256 \times 256$  since an image with resolution of  $128 \times 128$  pixels results in saliency maps that are too small to generate meaningful ROI patches. We report AUCs for predicting the risk of deterioration within 24, 48, 72, and 96 hours. When evaluating the deterioration risk curves, we report the concordance index with a reference time of 96 hours, as well as the average of the index over all possible reference times (3, 12, 24, 48, 72, 96, 144, and 192 hours).

		AUC / PR AUC				Concordance index	
		24 hours	48 hours	72 hours	96 hours	96 hours	Average
DenseNet-121	$128 \times 128$	0.663 (0.593, 0.724) / 0.214 (0.144, 0.309)	0.688 (0.627, 0.743) / 0.300 (0.224, 0.402)	0.700 (0.647, 0.751) / 0.370 (0.292, 0.461)	0.728 (0.675, 0.771) / 0.453 (0.373, 0.542)	0.700 (0.666, 0.733)	0.700 (0.664, 0.728)
	$256 \times 256$	<b>0.698</b> (0.632, 0.763) / <b>0.218</b> (0.153, 0.310)	<b>0.721</b> (0.668, 0.778) / 0.310 (0.238, 0.413)	<b>0.719</b> (0.670, 0.773) / <b>0.390</b> (0.318, 0.486)	<b>0.748</b> (0.701, 0.795) / <b>0.469</b> (0.392, 0.562)	0.701 (0.664, 0.736)	0.698 (0.662, 0.733)
	$512 \times 512$	0.682 (0.615, 0.747) / 0.208 (0.149, 0.305)	0.710 (0.656, 0.762) / <b>0.318</b> (0.238, 0.422)	0.709 (0.654, 0.762) / 0.383 (0.307, 0.480)	0.732 (0.684, 0.778) / 0.441 (0.366, 0.529)	<b>0.705</b> (0.673, 0.739)	<b>0.701</b> (0.669, 0.735)
	$1024 \times 1024$	0.680 (0.618, 0.741) / 0.180 (0.130, 0.259)	0.709 (0.655, 0.761) / 0.278 (0.212, 0.371)	0.716 (0.666, 0.766) / 0.369 (0.296, 0.469)	0.739 (0.691, 0.784) / 0.441 (0.366, 0.529)	0.701 (0.668, 0.734)	0.696 (0.663, 0.728)
COVID-GMIC	$256 \times 256$	0.664 (0.594, 0.735) / 0.202 (0.144, 0.303)	0.688 (0.629, 0.746) / 0.263 (0.200, 0.354)	0.699 (0.648, 0.747) / 0.342 (0.270, 0.431)	0.728 (0.682, 0.772) / 0.424 (0.356, 0.505)	0.712 (0.680, 0.745)	0.707 (0.673, 0.739)
	$512 \times 512$	<b>0.700</b> (0.635, 0.765) / <b>0.210</b> (0.154, 0.298)	0.714 (0.659, 0.767) / 0.300 (0.230, 0.395)	0.714 (0.662, 0.757) / <b>0.389</b> (0.314, 0.481)	0.733 (0.686, 0.776) / 0.443 (0.371, 0.532)	<b>0.713</b> (0.679, 0.748)	<b>0.708</b> (0.675, 0.742)
	$1024 \times 1024$	0.695 (0.627, 0.760) / 0.200 (0.142, 0.279)	<b>0.716</b> (0.661, 0.767) / <b>0.302</b> (0.230, 0.394)	<b>0.717</b> (0.663, 0.764) / 0.374 (0.301, 0.459)	<b>0.738</b> (0.692, 0.780) / <b>0.439</b> (0.368, 0.522)	0.686 (0.652, 0.722)	0.685 (0.653, 0.722)

**Impact of different transfer learning strategies.** In data-scarce applications, it is crucial to pretrain deep neural networks on a related task for which a large dataset is available, prior to fine-tuning on the task of interest [68, 69]. Given the relatively small number of COVID-19 positive cases in our dataset, we investigate the impact of different weight initialization strategies on our results. Specifically, we compare three strategies: 1) initialization by He et al. [70], 2) initialization with weights from models trained on natural images (ImageNet [66]), and 3) initialization with weights from models trained on chest X-ray images (ChestX-ray14 dataset [60]). We apply the initialization procedure to all layers except the last fully connected layer, which is always initialized randomly. We then fine-tune the entire network on our COVID-19 task.

Based on results shown in Supplementary Table 3, fine-tuning the network from weights pretrained on the ChestX-ray14 dataset is the most effective strategy for COVID-GMIC. This dataset contains over 100,000 chest X-ray images from more than 30,000 patients, including many with advanced lung disease. The images are paired with labels representing fourteen common thoracic observations: atelectasis, cardiomegaly, effusion, infiltration, mass, nodule, pneumonia, pneumothorax, consolidation, edema, emphysema, fibrosis, pleural thickening, and hernia. By pretraining a model to detect these conditions, we hypothesize that the model learns a representation that is useful for our downstream task of COVID-19 prognosis.

Supplementary Table 3: Model performance with 95% confidence intervals across three different initialization strategies: random initialization, initialization with the weights of the model pretrained on ImageNet [66] and initialization with the weights of the model pretrained model on the ChestX-ray14 dataset [60]. We report AUCs for each time window in the outcome classification task. When evaluating the deterioration risk curves, we report the concordance index with a reference time of 96 hours, as well as the average of the index over all discretized times (3, 12, 24, 48, 72, 96, 144, and 192 hours).

		AUC / PR AUC				Concordance index	
		24 hours	48 hours	72 hours	96 hours	96 hours	Average
DenseNet-121	Random	0.687 (0.621, 0.749) / 0.178 (0.134, 0.251)	0.699 (0.644, 0.750) / 0.258 (0.201, 0.339)	0.693 (0.639, 0.744) / 0.326 (0.264, 0.416)	0.705 (0.658, 0.750) / 0.386 (0.323, 0.474)	0.649 (0.612, 0.684)	0.648 (0.611, 0.683)
	ImageNet	<b>0.701</b> (0.639, 0.761) / 0.206 (0.152, 0.295)	<b>0.722</b> (0.668, 0.776) / 0.299 (0.232, 0.401)	<b>0.719</b> (0.670, 0.772) / 0.365 (0.294, 0.466)	<b>0.745</b> (0.701, 0.789) / 0.444 (0.375, 0.539)	0.686 (0.652, 0.720)	0.683 (0.651, 0.715)
	ChestX-ray14	0.687 (0.619, 0.758) / <b>0.216</b> (0.155, 0.317)	0.709 (0.653, 0.767) / <b>0.315</b> (0.239, 0.419)	0.710 (0.660, 0.763) / <b>0.373</b> (0.300, 0.464)	0.736 (0.691, 0.782) / <b>0.454</b> (0.384, 0.542)	<b>0.705</b> (0.673, 0.739)	<b>0.701</b> (0.669, 0.735)
COVID-GMIC	Random	0.675 (0.607, 0.741) / 0.174 (0.125, 0.247)	0.671 (0.617, 0.728) / 0.227 (0.177, 0.308)	0.686 (0.640, 0.732) / 0.290 (0.235, 0.366)	0.708 (0.664, 0.748) / 0.352 (0.294, 0.428)	0.643 (0.608, 0.680)	0.640 (0.607, 0.676)
	ImageNet	0.694 (0.631, 0.753) / 0.195 (0.138, 0.280)	0.709 (0.657, 0.761) / 0.258 (0.197, 0.351)	<b>0.724</b> (0.673, 0.769) / 0.347 (0.278, 0.431)	0.737 (0.692, 0.778) / 0.433 (0.360, 0.512)	0.684 (0.651, 0.716)	0.680 (0.649, 0.711)
	ChestX-ray14	<b>0.695</b> (0.626, 0.757) / <b>0.200</b> (0.142, 0.283)	<b>0.716</b> (0.659, 0.768) / <b>0.302</b> (0.228, 0.400)	0.717 (0.665, 0.762) / 0.374 (0.302, 0.463)	<b>0.738</b> (0.690, 0.783) / <b>0.439</b> (0.368, 0.532)	<b>0.713</b> (0.679, 0.748)	<b>0.708</b> (0.675, 0.742)

**Impact of training set size.** We evaluated the impact of the sample size used for training our machine learning models. Specifically, we evaluated our models on a subset of the training data,



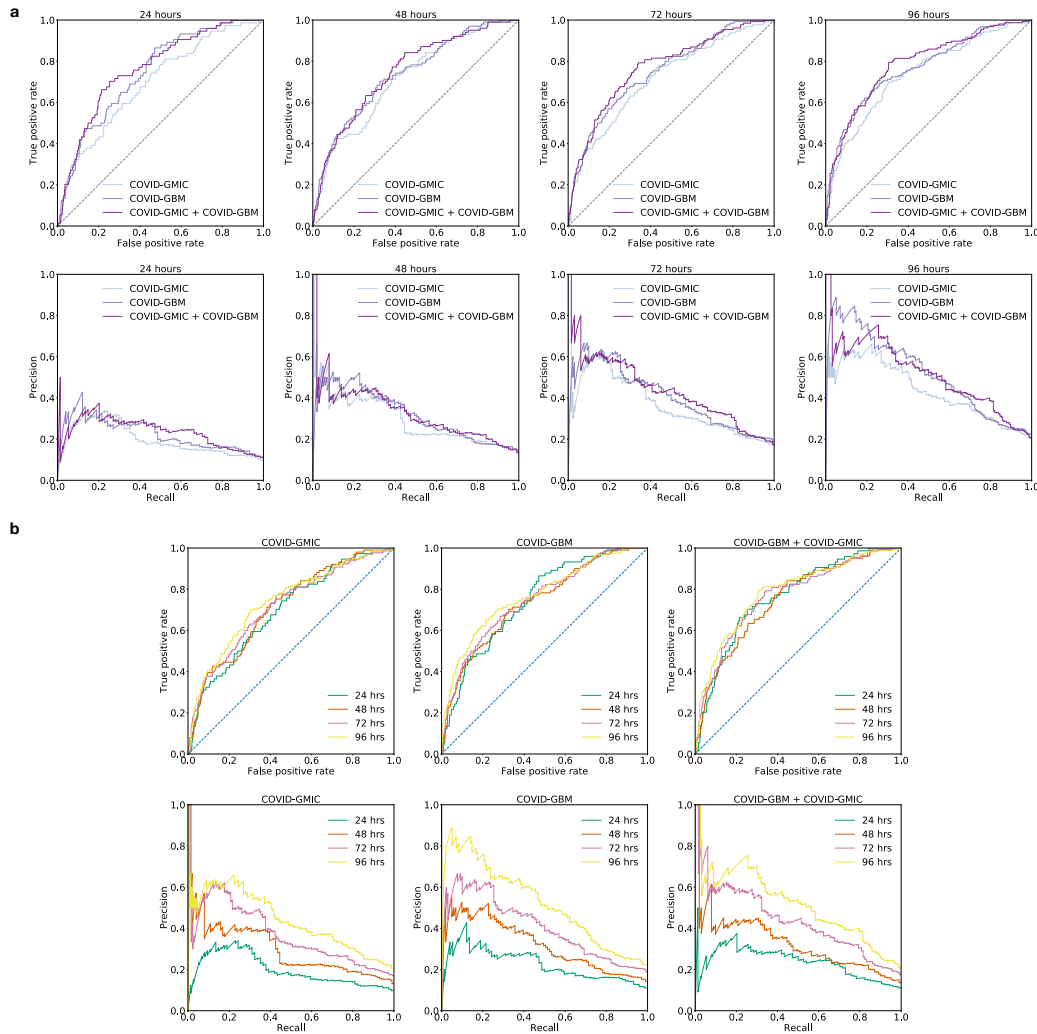
obtained by randomly sampling 12.5%, 25%, and 50% of the exams. Table 4 presents the AUCs and PR AUCs and the concordance indices achieved on the test set. It is evident that the performance of COVID-GMIC and COVID-GBM improve when increasing the number of images and clinical variables used for training, which highlights the importance of using a large dataset.

Supplementary Table 4: Model performance with 95% confidence intervals when using 12.5%, 25%, 50%, and 100% of the training data. We report AUCs for each time window in the adverse event prediction task. When evaluating the deterioration risk curves, we report the concordance index with a reference time of 96 hours, as well as the average of the index over all discretized times (3, 12, 24, 48, 72, 96, 144, and 192 hours).

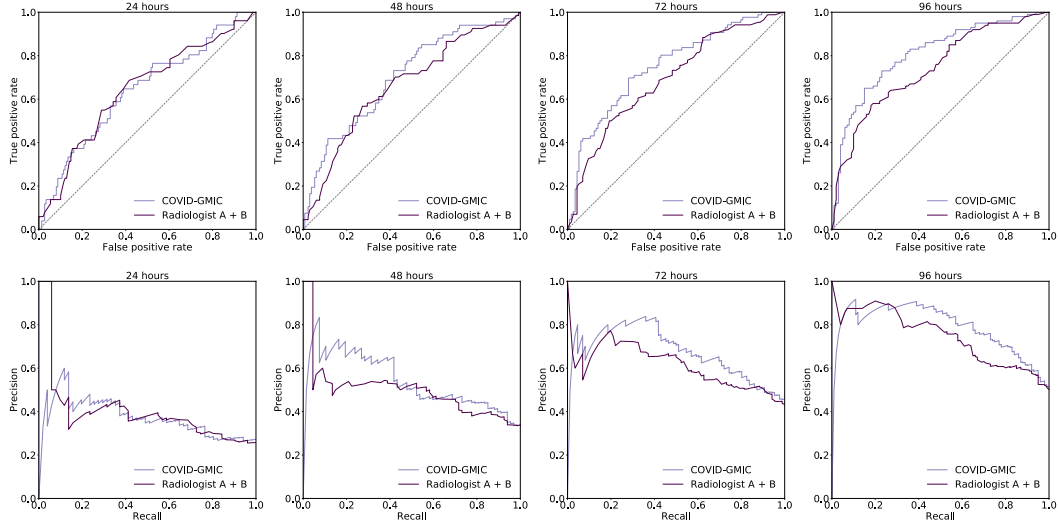
		AUC / PR AUC				Concordance index	
		24 hours	48 hours	72 hours	96 hours	96 hours	Average
DenseNet-121	12.5%	0.608 (0.538, 0.686) / 0.182 (0.123, 0.270)	0.653 (0.595, 0.712) / 0.265 (0.198, 0.353)	0.672 (0.622, 0.727) / 0.336 (0.271, 0.424)	0.703 (0.657, 0.752) / 0.415 (0.344, 0.500)	0.675 (0.642, 0.710)	0.670 (0.637, 0.704)
	25%	0.638 (0.568, 0.706) / 0.174 (0.121, 0.258)	0.678 (0.619, 0.735) / 0.266 (0.205, 0.362)	0.682 (0.630, 0.736) / 0.327 (0.261, 0.415)	0.711 (0.664, 0.760) / 0.408 (0.341, 0.495)	0.676 (0.643, 0.711)	0.671 (0.638, 0.705)
	50%	0.672 (0.607, 0.739) / 0.214 (0.150, 0.319)	0.699 (0.646, 0.754) / 0.303 (0.233, 0.397)	0.698 (0.649, 0.750) / 0.351 (0.285, 0.437)	0.725 (0.681, 0.771) / 0.433 (0.365, 0.517)	0.694 (0.660, 0.728)	0.691 (0.657, 0.725)
	100%	<b>0.687</b> (0.621, 0.753) / <b>0.216</b> (0.154, 0.317)	<b>0.709</b> (0.654, 0.763) / <b>0.315</b> (0.239, 0.417)	<b>0.710</b> (0.658, 0.761) / <b>0.373</b> (0.298, 0.475)	<b>0.736</b> (0.689, 0.781) / <b>0.454</b> (0.377, 0.552)	<b>0.705</b> (0.673, 0.739)	<b>0.701</b> (0.669, 0.735)
COVID-GMIC	12.5%	0.640 (0.577, 0.703) / 0.145 (0.110, 0.206)	0.672 (0.618, 0.723) / 0.231 (0.179, 0.316)	0.677 (0.626, 0.723) / 0.318 (0.249, 0.406)	0.695 (0.652, 0.738) / 0.384 (0.319, 0.474)	0.673 (0.640, 0.706)	0.668 (0.635, 0.701)
	25%	0.661 (0.598, 0.724) / 0.177 (0.125, 0.263)	0.672 (0.618, 0.728) / 0.254 (0.196, 0.346)	0.677 (0.631, 0.727) / 0.327 (0.266, 0.416)	0.693 (0.648, 0.737) / 0.395 (0.329, 0.477)	0.689 (0.655, 0.723)	0.680 (0.646, 0.714)
	50%	0.646 (0.577, 0.716) / 0.164 (0.116, 0.238)	0.681 (0.622, 0.738) / 0.266 (0.199, 0.360)	0.687 (0.632, 0.739) / 0.351 (0.274, 0.445)	0.716 (0.668, 0.763) / 0.424 (0.346, 0.516)	0.699 (0.665, 0.734)	0.690 (0.658, 0.723)
	100%	<b>0.695</b> (0.626, 0.753) / <b>0.200</b> (0.142, 0.276)	<b>0.716</b> (0.663, 0.769) / <b>0.302</b> (0.230, 0.395)	<b>0.717</b> (0.667, 0.767) / <b>0.374</b> (0.297, 0.461)	<b>0.738</b> (0.693, 0.782) / <b>0.439</b> (0.363, 0.521)	<b>0.713</b> (0.679, 0.748)	<b>0.708</b> (0.675, 0.742)
COVID-GBM	12.5%	0.674 (0.612, 0.739) / 0.262 (0.180, 0.371)	0.699 (0.645, 0.751) / 0.297 (0.228, 0.395)	0.710 (0.659, 0.754) / 0.395 (0.318, 0.480)	0.708 (0.661, 0.753) / 0.439 (0.362, 0.517)		
	25%	0.688 (0.636, 0.748) / 0.175 (0.130, 0.248)	0.716 (0.667, 0.766) / 0.319 (0.237, 0.411)	0.733 (0.688, 0.777) / 0.385 (0.309, 0.466)	0.739 (0.694, 0.783) / 0.476 (0.407, 0.550)		
	50%	0.743 (0.690, 0.787) / 0.210 (0.157, 0.301)	<b>0.752</b> (0.702, 0.797) / <b>0.325</b> (0.252, 0.425)	0.749 (0.703, 0.792) / 0.418 (0.341, 0.510)	0.751 (0.706, 0.791) / 0.482 (0.407, 0.568)		
	100%	<b>0.747</b> (0.692, 0.798) / <b>0.230</b> (0.167, 0.322)	0.739 (0.685, 0.791) / <b>0.325</b> (0.253, 0.425)	<b>0.750</b> (0.704, 0.794) / <b>0.408</b> (0.334, 0.502)	<b>0.770</b> (0.728, 0.811) / <b>0.523</b> (0.439, 0.611)		

## Supplementary Note 4: Additional results on the test sets

We visualize the receiver operating characteristic (ROC) and precision-recall (PR) curves on the test set in Supplementary Figure 4. In **a**, we group the results based on the predictive models (COVID-GMIC, COVID-GBM, and the ensemble of both), while in **b**, we group the performances based on the time window of the task (i.e., 24, 48, 72, and 96 hours). In Supplementary Figure 5, we visualize the ROC and PR curves on the test set considered in the reader study.



Supplementary Figure 4: **Receiver operating characteristic (ROC) and Precision-Recall (PR) curves for predicting the onset of adverse events within 24, 48, 72, and 96 hours evaluated on the test set.** **a**, ROC and PR curves are grouped by predictive models. Ensembling COVID-GMIC and COVID-GBM improves performance in almost all cases. **b**, ROC and PR curves are grouped by time window of the task. The AUC and PR AUC improve as the length of the time window increases, which is consistency across models. Numerical values of AUCs and PR AUCs can be found in Table 2.



Supplementary Figure 5: Test set ROC (top) and PR (bottom) curves of COVID-GMIC and the radiologists for predicting the risk of deterioration over 24, 48, 72, and 96 hours. These results suggest that COVID-GMIC performs comparably to the radiologists. Numerical values of AUCs and PR AUCs can be found in Table 2.

In Supplementary Table 5, we show the concordance index results across all time intervals for the best DenseNet-121 and COVID-GMIC-DRC models.

Supplementary Table 5: Concordance index (with 95% confidence intervals) of the DRC curves generated by the best DenseNet-121 and COVID-GMIC-DRC models. Both models use input images of size  $512 \times 512$  and are pretrained on the ChestX-ray14 dataset [60]. The results shows that the concordance index does not change much with the choice of time reference.

Time (in hours)	Concordance index								
	3	12	24	48	72	96	144	192	Ave.
DenseNet-121	0.681 (0.648,0.715)	0.694 (0.661,0.730)	0.701 (0.667,0.736)	0.702 (0.669,0.737)	0.703 (0.672,0.738)	0.705 (0.673,0.739)	0.706 (0.673,0.740)	0.705 (0.673,0.740)	0.701 (0.669,0.735)
COVID-GMIC-DRC	0.692 (0.650,0.723)	0.698 (0.660,0.732)	0.706 (0.672,0.740)	0.710 (0.674,0.743)	0.712 (0.679,0.748)	0.713 (0.679,0.748)	0.716 (0.684,0.751)	0.715 (0.682,0.750)	<b>0.708</b> (0.675,0.742)

## Supplementary Note 5: Model selection

We describe our model selection procedure used throughout the paper in Algorithm 1. For the ablation study in Table 4, we control the size of the dataset by setting the parameter  $u$  to 12.5, 25 and 50. Specifically, in that experiment, we randomly sampled  $u\%$  of the training set  $\mathcal{D}_t$  as the “universe”  $\mathcal{U}$  that our model used for training and validation.

---

**Algorithm 1** Model selection

---

**Input:** training set  $\mathcal{D}_t$ , test set  $\mathcal{D}_s$ , universe fraction  $u \in [0, 100]$ , and a predictive model  $\mathcal{M}$

**Output:**  $a^*$  performance of  $\mathcal{M}$  evaluated on  $\mathcal{D}_s$

- 1:  $\mathcal{U}$  = randomly sample  $u\%$  of data from  $\mathcal{D}_t$
  - 2:  $\Phi$  = 30 randomly sampled configuration of hyperparameters of the  $\mathcal{M}$
  - 3: **for** each hyperparameter configuration  $\phi_i \in \Phi$  **do**
  - 4:   **for**  $j \in \{1, 2, 3\}$  **do**
  - 5:     draw a random seed  $r_j$
  - 6:      $\mathcal{U}_t^j, \mathcal{U}_v^j$  = universe  $\mathcal{U}$  split into training and validation subset using the random seed  $r_j$
  - 7:      $\mathcal{M}_{ij}$  = trained  $\mathcal{M}$  using hyperparameter configuration  $\phi_i$  on  $\mathcal{U}_t^j$
  - 8:      $a_{ij}$  = performance of  $\mathcal{M}_{ij}$  evaluated on  $\mathcal{U}_v^j$
  - 9:   **end for**
  - 10:    $a_i = \frac{1}{3} \sum_{j=1}^3 a_{ij}$
  - 11: **end for**
  - 12:  $\mathcal{A} = \{a_i \mid \forall \phi_i \in \Phi\}$
  - 13:  $\mathcal{B} = \{\mathcal{M}_{ij} \mid \forall a_i \in \text{top-3}(\mathcal{A})\}$
  - 14:  $\mathcal{M}^*$  = an equally weighted ensemble of all models in  $\mathcal{B}$
  - 15:  $a^*$  = performance of  $\mathcal{M}^*$  on  $\mathcal{D}_s$
  - 16: **return**  $a^*$
-