

---

# Accelerating COVID-19 Differential Diagnosis with Explainable Ultrasound Image Analysis

---

Jannis Born<sup>1,\*</sup>Nina Wiedemann<sup>2,\*</sup>Gabriel Brändle<sup>3</sup>Charlotte Buhre<sup>4</sup>Bastian Rieck<sup>1</sup>Karsten Borgwardt<sup>1</sup>

\*Shared first-authorship. {jborn,wnina}@ethz.ch

<sup>1</sup>DEPARTMENT OF BIOSYSTEMS SCIENCE AND ENGINEERING, ETH ZURICH, SWITZERLAND<sup>2</sup>DEPARTMENT OF COMPUTER SCIENCE, ETH ZURICH, SWITZERLAND<sup>3</sup>HIRSBLANDEN CLINIQUE DES GRANGETTES, GENEVA, SWITZERLAND<sup>4</sup>MEDIZINISCHE HOCHSCHULE BRANDENBURG THEODOR FONTANE, GERMANY

## Abstract

Controlling the COVID-19 pandemic largely hinges upon the existence of fast, safe, and highly-available diagnostic tools. Ultrasound, in contrast to CT or X-Ray, has many practical advantages and can serve as a globally-applicable first-line examination technique. We provide the largest publicly available lung ultrasound (US) dataset for COVID-19 consisting of 106 videos from three classes (COVID-19, bacterial pneumonia, and healthy controls); curated and approved by medical experts. On this dataset, we perform an in-depth study of the value of deep learning methods for differential diagnosis of COVID-19. We propose a frame-based convolutional neural network that correctly classifies COVID-19 US videos with a sensitivity of  $0.98 \pm 0.04$  and a specificity of  $0.91 \pm 0.08$  (frame-based sensitivity  $0.93 \pm 0.05$ , specificity  $0.87 \pm 0.07$ ). We further employ class activation maps for the spatio-temporal localization of pulmonary biomarkers, which we subsequently validate for human-in-the-loop scenarios in a blindfolded study with medical experts. Aiming for scalability and robustness, we perform ablation studies comparing mobile-friendly, frame- and video-based architectures and show reliability of the best model by aleatoric and epistemic uncertainty estimates. We hope to pave the road for a community effort toward an accessible, efficient and interpretable screening method and we have started to work on a clinical validation of the proposed method. Data and code are publicly available.

## 1 Introduction

To date, SARS-CoV-2 has infected several millions and COVID-19 has killed hundreds of thousands around the globe. Its long incubation time calls for fast, accurate, and reliable techniques for early disease diagnosis to successfully fight the spread [29]. The standard genetic test (RT-PCR), suffers from a processing time of up to 2 days [33], several publications reported sensitivity as low as 70% [2, 25] and a recent meta-analysis estimated the *false negative* rate to be at least 20% over the course of the infection [26]. Medical imaging has great potential to complement the diagnostic process as a fast assessment tool that guides further PCR-testing, especially in triage situations [16]. Currently, CT scans are the gold standard for pneumonia [8] and are considered relatively reliable for COVID-19 diagnosis [2, 5, 18], although a significant amount of patients exhibit normal CT scans [52]. However, performing CT is expensive and highly irradiating, posing risks of infection

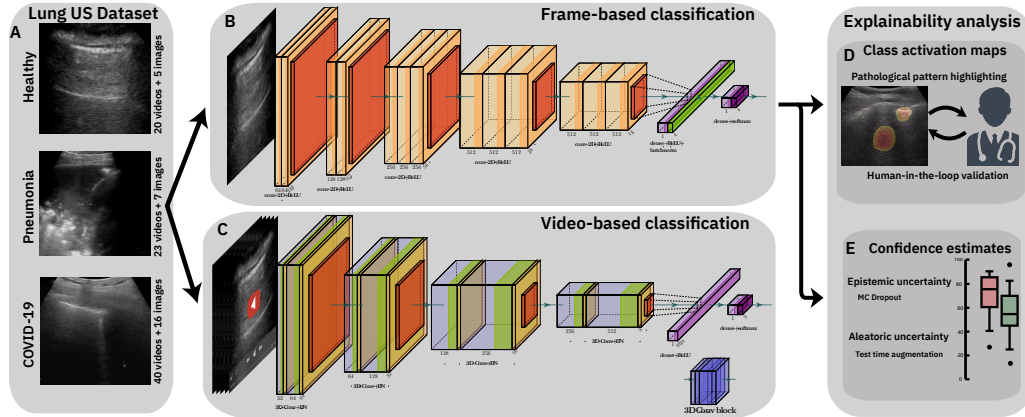


Figure 1: **Flowchart of our contribution.** **A:** 3 samples from our public COVID-19 lung US dataset. *Top:* Healthy lung with horizontal A-lines, *Middle:* pneumonia infected lung with alveolar consolidations, *Bottom:* SARS-CoV-2 infected lung with subpleural consolidation and a focal B-line. **B,C:** We fine-tune and compare frame- and video-based CNNs on this new dataset and demonstrate the feasibility of differential diagnosis from ultrasound. **D:** Class activation maps highlight patterns that drove the model’s decision and are reviewed and evaluated for diagnostic value by medical experts. **E:** Uncertainty techniques are employed and shown to equip the model with the ability to recognize samples with high error probability.

for patients and staff, requires extensive sterilization [34], and is reserved for developed countries; there are only  $\sim 30k$  CT scanner globally [13]. Chest X-ray (CXR) scans are still the first line examination, despite some reports of low specificity and sensitivity for COVID-19 (for example, Weinstock et al. [51] found 89% normal CXR in 493 COVID-19 patients). Ultrasound (US), by contrast, is a cheap, safe, non-invasive and repeatable technique that can be performed with portable devices at patient bedside and is ubiquitously available around the globe. Over the last two decades, ultrasound became an established tool to diagnose pulmonary diseases [14, 30, 36], has been forecast to replace radiographic techniques [9], was demonstrated to be superior to CXR for pulmonary diseases [15, 38], and started to replace X-Ray as first-line examination [1, 10].

In the COVID-19 outbreak, a growing body of evidence for disease-specific patterns in US (e.g. B-lines and subpleural consolidations) has led to advocacy for an amplified role of US from the research community [11, 27, 44, 46] and radiologists reported great agreement between US and CT findings for COVID-19 infections [19, 37]. Moreover, in triage situations or in third-world countries, where CT, PCR and CXR are not available, US was demonstrated to be a valuable patient stratification technique for pneumonia [3, 17]. This gives US, in contrast to other imaging techniques, the potential to become a globally applicable first-line examination method [45]. However, the relevant pattern are hard to discern for humans [35, 47], calling into play medical image analysis based on machine learning technique as a decision support tool for physicians. Here, we provide the first study of automatic lung ultrasound analysis for differential diagnosis of bacterial and viral pneumonia; aiming to develop a medical decision support tool.

**Related work.** Literature on exploiting medical image analysis and computer vision techniques to classify or segment CT or CXR data of COVID-19 patients recently exploded (for reviews, see Shi et al. [42], Ulhaq et al. [48], for a list of public data sources see Kalkreuth and Kaufmann [24]). For example, in an early study, Butt et al. [12] reported a sensitivity of 98% (specificity 92%) in a binary classification on CT scans from 110 COVID-19 patients, while Mei et al. [33] very recently achieved equal sensitivity (but lower specificity) compared to senior radiologists in detecting COVID-19 from CT and clinical information of 279 patients. US instead has been neglected heavily by the ML community [6]; only the Italian COVID-19 Lung Ultrasound (ICLUS) project has proposed a deep learning approach for a severity assessment of COVID-19 from ultrasound data [39]. The work convincingly predicts disease severity and segments COVID-19 specific patterns, building up on their previous work on localizing B-lines [49]. The paper claims to release a dataset of annotated COVID-19 cases, but to date, no annotations are available. While this effort is highly relevant for disease monitoring, it is not directly applicable for first-line diagnosis, where the main problem lies

in *distinguishing* COVID-19 from other pneumonia. We aim to close this gap with our approach to classify COVID-19, healthy, and pneumonia point-of-care ultrasound (POCUS) images.

**Our contributions.** Figure 1 depicts a graphical overview of our contributions. We provide the largest publicly-available dataset of lung US recordings consisting of 106 videos. This dataset is heterogeneous and mostly from public sources, but was curated manually and approved by a medical doctor. We further take a first step towards a tool for differential diagnosis of pulmonary diseases, here especially focused on bacterial and viral pneumonia such as COVID-19. An earlier version of our dataset alongside some preliminary results, we already presented in [7]. Without deprivation of novelty, we here demonstrate that competitive performance can be achieved from raw US recordings, thereby challenging the current focus on irradiating imaging techniques. Moreover, we employ explainability techniques such as class activation maps or uncertainty estimates and present a roadmap towards an automatic detection system that can *segment* and *highlight* relevant spatio-temporal patterns. Such a system could not only lead to superior diagnostic performance, as was partially shown for CT [33], but can also reduce the time doctors require to make a diagnosis [41]. Our approach is of evident need because physicians must be trained thoroughly to reliably differentiate COVID-19 from pneumonia [35], making it necessary to use powerful deep learning to develop a system that can complement the work of physicians in a timely manner.

## 2 A lung ultrasound dataset for COVID-19 detection

We provide the to-date largest pre-processed and publicly available lung POCUS dataset<sup>1</sup>, comprising samples of COVID-19 patients, pneumonia-infected lungs and healthy patients. As shown in Table 1, we collected and gathered 139 recordings (106 videos + 33 images) recorded with either convex or linear probes, where the latter is a higher frequency probe yielding more superficial images. Our sources comprise community platforms, open medical repositories, health-tech companies, other scientific literature, and data recorded by healthy volunteers from our team. Main sources of data were [grepmed.com](https://www.grepmed.com), [thepocusatlas.com](https://www.thepocusatlas.com), [butterflynetwork.com](https://www.butterflynetwork.com) and [radiopaedia.org](https://www.radiopaedia.org). All samples of our database were annotated and approved by a medical doctor; moreover notes on the visible patterns in each video (e.g. B-Lines or consolidations) were added. In all collected videos of COVID-19 and pneumonia, disease-specific patterns are visible. The dataset is heterogeneous in terms of resolution, frame rates, the conducted lung US protocol, the devices used, and little to no meta data about the patients is available. For more details about the dataset and metadata, see the release on GitHub.

	Convex		Linear		Sum
	Vid.	Img.	Vid.	Img.	
<b>COVID</b>	40	16	4	3	<b>63</b>
<b>BP</b>	23	7	2	2	<b>34</b>
<b>VP</b>	3	–	4	–	<b>7</b>
<b>Healthy</b>	20	5	10	–	<b>35</b>
<b>Sum</b>	<b>86</b>	<b>28</b>	<b>20</b>	<b>5</b>	<b>139</b>

Table 1: Number of videos and images in our dataset, per class and probe. BP is bacterial pneumonia, VP is viral pneumonia.

## 3 Differential diagnosis of COVID-19 with lung ultrasound

### 3.1 Experimental setup

**Data processing.** All experiments are conducted on data recorded with convex ultrasound probes, the standard probe for lung assessment that allows to see deeply into the lung [31]. We manually processed all convex ultrasound recordings and split them into images at a frame rate of 3Hz (with maximal 30 frames per video), leading to a database of 693 COVID-19, 377 bacterial pneumonia, and 295 healthy control images. For examples see Figure 1A. All images were cropped to a quadratic window excluding measure bars and texts and artifacts on the borders before they were resized to  $224 \times 224$  pixels. Apart from the independent test data, all reported results were obtained in a 5-fold stratified cross validation. It was ensured that the frames of a single video are present within a single fold only, and that the number of samples per class is similar in all folds. All models were trained to classify images as COVID-19, pneumonia, healthy, or uninformative. The latter consists of ImageNet pictures as well as neck ultrasound data; we added these picture for the purpose of detecting out-of-distribution data (thus making the model more robust). This is particularly relevant for public web-based inference services. In this paper, we present all results *omitting the*

<sup>1</sup>[https://github.com/jannisborn/covid19\\_pocus\\_ultrasound/tree/master/data](https://github.com/jannisborn/covid19_pocus_ultrasound/tree/master/data)

*uninformative class*, as it is not relevant for the analysis of differential diagnosis performance and would bias the results (please refer to appendix A.4.1 for results including uninformative data). Furthermore, we use data augmentation techniques (horizontal and vertical flips, rotations up to  $10^\circ$  and translations of up to 10%) to diversify the dataset and prevent overfitting.

**Frame-based models.** Our backbone neural architecture is a VGG-16 [43] that is compared to NasNET Mobile, a light-weight alternative [55] that uses less than  $\frac{1}{3}$  of the parameters of VGG and was optimized for applications on portable devices. Both models are pre-trained on Imagenet and fine-tuned on the frames sampled from the videos. Specifically, we use two variants of VGG-16 that we name VGG and VGG-CAM. VGG-CAM has a single dense layer following the convolutions, thus enabling the usage of plain CAMs, class activation maps [53], whereas VGG has an additional dense layer with ReLU activation and batch normalization.

Considering the recent work of Roy et al. [39] on lung US segmentation and severity prediction for COVID-19, we investigated whether a segmentation-targeted network can also add value to the prediction in differential diagnosis. We implemented two approaches building upon the pre-trained model of Roy et al. [39], an ensemble of three separate U-Net-based models (U-Net, U-Net++, and DeepLabv3+, with a total of  $\sim 19.5$ M parameters). First, VGG-Segment is identical to VGG, however instead of training on the raw US data, we train on the segmented images from the ensemble (see example in Appendix A.2). Although it might seem unconventional, we hypothesized that the colouring entails additional information that might simplify classification. Secondly, in Segment-Enc the bottleneck layer of each of the three models is used as a feature encoding of the images, resulting in 560 filter maps that are fed through two dense layers of size 512 and 256 respectively. The encoding weights are fixed during training. Both settings are compared to the other models that directly utilize the raw images. For more details on the architectures and the training procedure, please refer to appendix A.1.

**Video-based model.** In comparison to a naïve, frame-based video classifier (obtained by averaging scores of all frames), we also investigate Models Genesis, a generic model for 3D medical image analysis pretrained on lung CT scans [54]. For Models Genesis, the videos are split into chunks of 5 frames each, sampled at a frame rate of 5Hz. 5-fold cross validation is performed using the same split as for frame-based classifiers. Individual images were excluded, leaving aside 86 videos (from which 10 were excluded due to too many frames with artifacts such as moving pointers) which were split into 292 video chunks.

### 3.2 Frame-based experiments

Table 2 shows a detailed comparison of the three best models in terms of recall, precision, specificity and F1-scores, as well as MCC. Overall, both VGG and VGG-CAM achieve promising performance with an accuracy of  $90 \pm 2\%$  and  $90 \pm 5\%$  respectively on a 5-fold CV of 1,365 frames. Concerning per-class prediction accuracies, it is evident that bacterial pneumonia infections are distinguished best, with recall, precision, and specificity above 0.93 for VGG and VGG-CAM, indicating the models' ability to recognize strong irregularities in lung images. Although VGG slightly outperforms VGG-CAM, we explored the latter more in detail, due to its higher sensitivity for COVID-19 and its better performance when taking into account the class activation maps. Figure 2a visualizes the results of the VGG-CAM model for each binary detection task as a ROC curve, showing ROC-AUC scores of 0.94 and above for COVID-19 and the other two classes, while depicting the point where the accuracy is maximal for each class. The false positive rate at the maximal-accuracy point is larger for COVID-19 than for pneumonia and healthy patients. In a clinical setting, where false positives are less problematic than false negatives, this property is highly desirable. Since the data is imbalanced, we also plot the precision-recall curve in Figure 2b, which confirms that pneumonia is the class that is predicted most easily. In addition, the confusion matrices in Figure 2c and Figure 2d further detail the predictions of VGG-CAM; we observe that the high sensitivity for COVID-19 (0.93, 642 out of 693 frames) comes at a cost of 22% false positives from the healthy class. For further results including the ROC- and precision-recall curves of all three models see Appendix A.4.

**Ablation study with segmentation models.** Lung US recordings are noisy and operator-dependent, posing difficulties for the classification of raw data. Hence, we compare VGG and VGG-CAM to VGG-Segment where all frames are segmented (i.e. classified on a pixel level into pathological

	Class	Recall	Precision	F1-score	Specificity	MCC
<b>VGG</b> Acc.: 0.90, Bal.: <b>0.90</b> # Par.: 14 747 971	COVID-19	0.89 ± 0.06	<b>0.91</b> ± 0.05	0.90 ± 0.03	<b>0.92</b> ± 0.04	0.80 ± 0.03
	Pneumonia	0.94 ± 0.06	0.93 ± 0.05	0.94 ± 0.05	0.97 ± 0.02	0.91 ± 0.07
	Healthy	<b>0.85</b> ± 0.11	0.83 ± 0.09	<b>0.83</b> ± 0.07	0.95 ± 0.03	0.79 ± 0.08
<b>VGG-CAM</b> Acc.: 0.90, Bal.: 0.88 # Par.: 14 716 227	COVID-19	0.93 ± 0.05	0.87 ± 0.07	0.90 ± 0.05	0.87 ± 0.06	0.79 ± 0.09
	Pneumonia	0.94 ± 0.05	0.95 ± 0.05	0.94 ± 0.04	0.98 ± 0.02	0.92 ± 0.06
	Healthy	0.78 ± 0.10	0.86 ± 0.08	0.81 ± 0.05	0.96 ± 0.02	0.77 ± 0.06
<b>NASNetMobile</b> Acc.: 0.76, Bal.: 0.71 # Par.: 4 814 487	COVID-19	0.87 ± 0.10	0.74 ± 0.12	0.80 ± 0.10	0.72 ± 0.07	0.59 ± 0.15
	Pneumonia	0.79 ± 0.15	0.88 ± 0.08	0.83 ± 0.10	0.96 ± 0.03	0.77 ± 0.15
	Healthy	0.47 ± 0.03	0.61 ± 0.13	0.53 ± 0.05	0.92 ± 0.04	0.43 ± 0.07
<b>VGG-Segment</b> Acc.: <b>0.91</b> , Bal.: 0.89 # Par.: 34 018 074	COVID-19	<b>0.96</b> ± 0.05	0.89 ± 0.06	<b>0.92</b> ± 0.04	0.88 ± 0.07	<b>0.84</b> ± 0.07
	Pneumonia	<b>0.96</b> ± 0.03	<b>0.97</b> ± 0.03	<b>0.96</b> ± 0.02	<b>0.99</b> ± 0.01	<b>0.95</b> ± 0.03
	Healthy	0.77 ± 0.14	<b>0.91</b> ± 0.08	0.82 ± 0.08	<b>0.97</b> ± 0.03	0.79 ± 0.09
<b>Segment-Enc</b> Acc.: 0.90, Bal.: 0.89 # Par.: 19 993 307	COVID-19	0.92 ± 0.09	0.91 ± 0.06	0.91 ± 0.03	0.90 ± 0.06	0.82 ± 0.04
	Pneumonia	0.95 ± 0.04	0.89 ± 0.12	0.92 ± 0.07	0.95 ± 0.06	0.89 ± 0.08
	Healthy	0.79 ± 0.17	0.89 ± 0.10	0.82 ± 0.11	0.98 ± 0.02	<b>0.80</b> ± 0.12

Table 2: **Performance comparison.** Comparison of the tested classification models on 5-fold cross validation for each class. Acc. abbreviates accuracy, Bal. balanced accuracy, MCC Matthews Correlation Coefficient and Par. the number of parameters. For each class and each column the best model is highlighted in bold. While VGG and VGG-CAM achieve an accuracy of 90%, the performance can, for the price of doubled computational costs, be further pushed with pre-trained lung US segmentation models from Roy et al. [39]. NASNetMobile instead is inferior but significantly smaller.

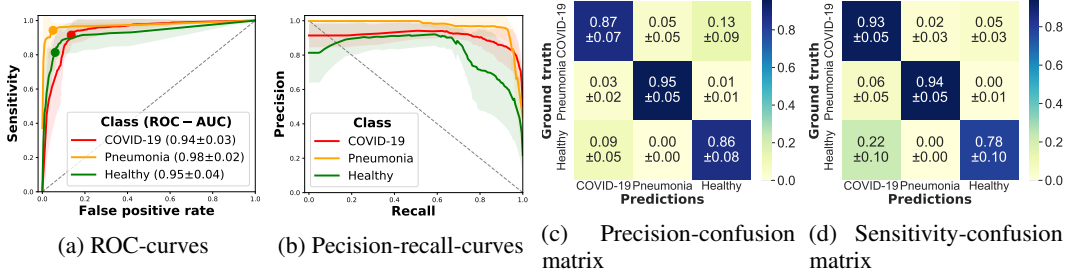


Figure 2: **Performance of the VGG-CAM model.** Per-class ROC-AUC, sensitivity, and precision are shown on the diagonals of the normalized confusion matrices, highlighting the model’s ability to distinguish COVID-19 from pneumonia and healthy lung images.

patterns) with the model from Roy et al. [39]; see Appendix A.1) for an example input. The relevant rows in Table 2 exhibit mixed results: while training on segmented images *improves* most relevant performance metrics slightly (higher accuracy, COVID-19 sensitivity, and MCC scores), balanced accuracy is *inferior* compared to VGG. Since this small increase in predictive performance comes at the cost of a large increase in model size (due to the ensemble of three independent models; selecting only one of the models resulted in inferior performance), we considered Segment-Enc, i.e. a dense model classifying the 560-dimensional encoding produced by the pre-trained segmentation models. Segment-Enc achieved comparable performance for most metrics, apart from lower scores for pneumonia detection. Since the difference in performance is only marginal, and the architectures of VGG-Segment and Segment-Enc prohibit the computation of class activation maps, we prefer to focus on the analysis of VGG-CAM in the following.

**Ablation study on other architectures.** Initially, we had tested further models proposed for medical image analysis, such as COVID-Net (previously used for the classification of X-Ray images [50]), and an architecture following [28] based on a Res-Net [22], but we observed that the experiments on our data resulted in *significantly worse* results. Last, we tested several smaller networks such as MobileNet[23] as an additional ablation study, with NASNetMobile [55] performing best. As most ultrasound devices are portable and real-time inference on the devices is technically feasible, resource-efficient networks are highly relevant and could supersede web-based inference. Due to low precision and recall on healthy data, our fine-tuned NASNetMobile is less performant than VGG-CAM,

	Class	Recall	Precision	F1-score	Specificity	MCC
<b>VGG-CAM</b> Acc.: <b>0.94</b> , Bal.: <b>0.93</b> # Par.: 14, 716 227	COVID-19	<b>0.98</b> $\pm$ 0.04	<b>0.91</b> $\pm$ 0.08	<b>0.94</b> $\pm$ 0.04	0.91 $\pm$ 0.08	<b>0.89</b> $\pm$ 0.06
	Pneumonia	<b>1.00</b> $\pm$ 0.00	<b>1.00</b> $\pm$ 0.00	<b>1.00</b> $\pm$ 0.00	<b>1.00</b> $\pm$ 0.00	<b>1.00</b> $\pm$ 0.00
	Healthy	0.80 $\pm$ 0.16	<b>0.96</b> $\pm$ 0.08	<b>0.86</b> $\pm$ 0.08	<b>0.98</b> $\pm$ 0.03	<b>0.84</b> $\pm$ 0.09
<b>Models Genesis</b> Acc.: 0.87, Bal.: 0.87 # Par.: 7,559,043	COVID-19	0.80 $\pm$ 0.19	0.90 $\pm$ 0.09	0.83 $\pm$ 0.11	<b>0.92</b> $\pm$ 0.08	0.74 $\pm$ 0.11
	Pneumonia	<b>1.00</b> $\pm$ 0.00	<b>1.00</b> $\pm$ 0.00	<b>1.00</b> $\pm$ 0.00	<b>1.00</b> $\pm$ 0.00	<b>1.00</b> $\pm$ 0.00
	Healthy	<b>0.82</b> $\pm$ 0.15	0.75 $\pm$ 0.21	0.75 $\pm$ 0.07	0.89 $\pm$ 0.10	0.69 $\pm$ 0.09

Table 3: **Video classification results.** The frame-based model VGG-CAM outperforms the 3D CNN Model Genesis, showing high accuracy (94%), recall, precision for COVID-19 and pneumonia detection.

but also requires less than a third of the parameters, thus providing a first step towards real-time on-device inference.

### 3.3 Video-based experiments

To investigate the need for a model with the ability to detect spatiotemporal patterns in lung US, we explored Model Genesis, a pretrained 3D-CNN designed for 3D medical image analysis [54]. Table 3 contrasts the frame-based performance of VGG-CAM model to Model Genesis. The video classifier is outperformed by VGG-CAM, with a video accuracy of 94% compared to 87%. Note that all videos of pneumonia-infections are classified correctly, while especially Model Genesis struggles with the prediction of healthy patients. Considering that only 292 video-chunks were available for training Model Genesis, while 1356 images are used to train VGG-CAM, even extended through data augmentation techniques, it is likely that video-based classification may improve with increasing data availability.

### 3.4 Evaluation on independent test data

Very recently, the ICLUS initiative released 60 COVID-19 lung US recordings from Italian patients<sup>2</sup>[39]. The data is not annotated, but was initially assumed to contain only COVID-19 videos, based on its general description. We evaluated the performance of the VGG-CAM model on all 40 convex probes from ICLUS, alongside 24 recordings from healthy controls (6 viewpoints each) and 2 videos from public sources (healthy), jointly comprising an independent test dataset of 66 videos.

We predicted all frames as an average of the five VGG-CAM models trained in cross-validation. The model achieves a frame-prediction accuracy of 83.3%, divided into 89.5% for healthy-patient data and 74% for COVID-19 videos. Furthermore, averaging the class probabilities over all frames, VGG-CAM achieves a video classification accuracy of 92.2% and 77.5%, respectively. Notably, the four healthy patients are *all* classified correctly if summarized across viewpoints. Combining both datasets, the sensitivity of detecting COVID-19 corresponds to the accuracy (0.775) with a precision score of 0.94 (no video was classified as bacterial pneumonia).

**Evaluation by domain experts.** We further investigated the comparably low sensitivity on the COVID-19 data (ICLUS) with the help of two medical experts. When asked for their unbiased diagnosis of the incorrectly-predicted videos, they independently reported for 6 out of 9 videos that no disease-specific patterns can be observed (“A-lines, normal lung”). While these findings support the performance of our model, the *true* label of the data remains unclear. In addition, the dataset may contain further healthy-patient data which was incorrectly predicted as COVID-19. At this point, we can safely conclude that test data performance is highly promising, in particular considering the high accuracy for healthy patients, but requires further validation with independent and labeled data.

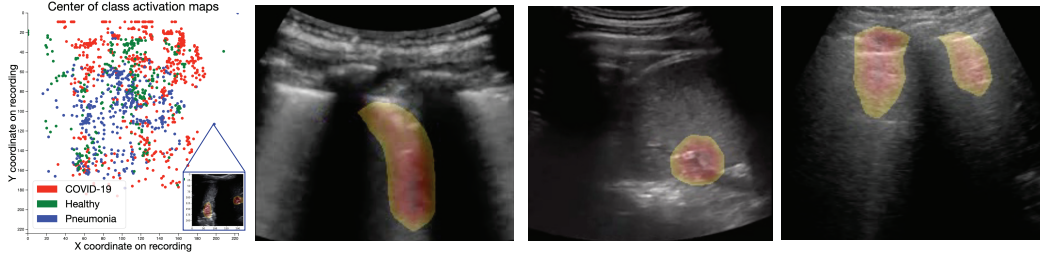


Figure 3: **Class activation maps.** **Left:** Interactive scatterplot of the origins of the CAMs across the entire dataset, colored by class. While the data seems rather unstructured, pneumonia-CAMs have lower  $y$ -coordinates than COVID-19 and healthy samples. **Rest:** Exemplary CAMs for COVID-19 (highlighting a B-line), bacterial pneumonia (highlighting pleural consolidations) and healthy lungs (highlighting A-lines).

## 4 Model explainability

### 4.1 Class activation maps

Class activation mapping (CAM) are a popular technique for model explainability that exploits global average pooling and allows to compute class-specific heatmaps that indicate the discriminative regions of the image that caused the particular class activity of interest [53]. For healthcare applications, CAMs, or their generalization Grad-CAMs [40], can provide valuable decision support by unravelling whether a model’s prediction was based on visible pathological patterns. Moreover, CAMs can guide doctors and point to informative patterns, especially relevant in time-sensitive (triage) or knowledge-sensitive (third-world countries) situations.

**Results.** Figure 3 shows representative CAMs in the three rightmost panels. They highlight the most frequent US pattern for the three classes, COVID-19 (vertical B-lines), bacterial pneumonia (consolidations), and healthy (horizontal A-line). For a more quantitative estimate, we computed the points of maximal activation of the CAMs for each class (abbreviated as **C**, **P**, and **H**) and all samples of the dataset (see Figure 3 left). While, in general, the heatmaps are fairly distributed across the probe, pneumonia related features were rather found in the center and bottom part, especially compared to COVID-19 and healthy patterns<sup>3</sup>. Please refer to Appendix A.5 for a density plot. To assess to what extent the differences between the individual distributions are significant, we employed *maximum mean discrepancy* (MMD), a metric between statistical distributions [21] that enables the comparison of distributions via kernels, i.e. generic similarity functions. Given two coordinates  $x, y \in \mathbb{R}^2$  and a smoothing parameter  $\sigma \in \mathbb{R}$ , we use a Gaussian kernel  $k(x, y) := \exp(-\|x-y\|^2/\sigma^2)$  to assess the dissimilarity between  $x$  and  $y$ . Following Gretton et al. [21], we set  $\sigma$  to the median distance in the aggregated samples (i.e. all samples, without considering labels). We then calculate MMD values for the distance between the three classes, i.e.  $\text{MMD}(\mathbf{C}, \mathbf{P}) \approx 0.0051$ ,  $\text{MMD}(\mathbf{C}, \mathbf{H}) \approx 0.0061$ , and  $\text{MMD}(\mathbf{P}, \mathbf{H}) \approx 0.0065$ . Repeating this calculation for 5000 bootstrap samples per class (see Figure 9 for the resulting histograms), we find that the observe achieved significance levels of the intra-class MMD values of well below an  $\alpha = 0.05$  significance level.

**Expert validation of CAMs for Human-in-the-loop settings.** A potential application of our framework is a human-in-the-loop (HITL) setting with CAMs as a core component of the decision support tool that highlights pulmonary biomarkers and guides the decision makers. Since the performance of qualitative methods like CAMs can only be validated with the help of doctors, we conducted a blind-folded study with two medical experts experienced in the diagnostic process with ultrasound recordings. The experts were shown 50 videos (14 COVID-19, 21 pneumonia, 14 regular) comprising all non-proprietary video data which was correctly classified by the model. The class activation map for the respective class was computed two times, first with an average of all five models that were trained, and second only with the model that did not see any frame of the video during training (called train- and test-CAMs in the following). Both experts were asked to compare

<sup>2</sup>40 convex + 20 linear probes are available from <https://iclus-web.bluetensor.ai/>

<sup>3</sup>The interactive HTML and a few exemplary CAM videos are available at: <https://bit.ly/3eASPC8>

both activation maps for all 50 videos, and to score them on a scale of  $-3$  (“the heatmap is only distracting”) to  $3$  (“the heatmap is very helpful for diagnosis”).

First, the CAMs were overall perceived useful and the train and test CAMs were assigned a *higher* average score of 0.45 and 0.81 respectively. Second, disagreeing in only 8% of the cases, both experts independently decided for the test-CAM with a probability of 56%. Hence, the test-CAMs are not inferior to the train-CAMs, however non-significant in a Wilcoxon signed-rank test. However, train- and test-CAM both scored best for videos of bacterial pneumonia, lacking performance for videos of healthy and COVID-19 patients. Specifically, test-CAM received an average score of 0.81, divided into  $-0.25$  for COVID-19, 2.05 for pneumonia, and 0 for healthy patients. Third, the experts were asked to name the pathological patterns visible in general, as well as the patterns that were highlighted by the heatmap. Figure 4 shows the average ratio of pattern that were correctly highlighted by the CAM model, where the patterns listed by the more senior expert are taken as the ground truth for each video.

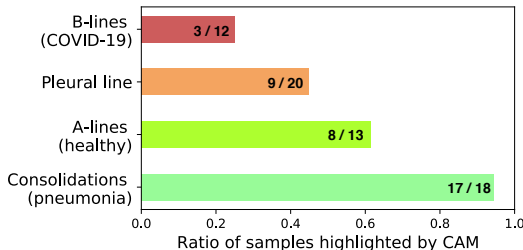


Figure 4: Patterns highlighted by CAMs compared to visible patterns in the video.

Interestingly, the high performance of our model in classifying videos of bacterial pneumonia is probably explained by the model’s ability to detect consolidated areas, where 17 out of 18 are correctly classified. Moreover, A-lines are highlighted in  $\sim 60\%$  of the normal lung recordings. Problematically, in 13 videos mostly fat, muscles or skin is highlighted, which has to be studied and improved in future work.

## 4.2 Confidence estimates

The ability to quantify states of high uncertainty is of crucial importance for medical image analysis and computer vision applications in healthcare. We assessed this via independent measures of epistemic (model) uncertainty (by drawing Monte Carlo samples from the approximate predictive posterior [20]) and aleatoric (data) uncertainty (by means of test time data augmentation [4]). The sample standard deviation of 10 forward passes is interpreted as inverse, empirical confidence score  $\in [0, 1]$  (for details see appendix). The epistemic confidence estimate was found to be highly correlated with the correctness of the predictions ( $\rho = 0.41, p < 4e-124$ , mean confidence of 0.75 and 0.26 for correct and wrong predictions), while the aleatoric confidence was found correlated to a lesser extent ( $\rho = 0.29, p < 6e-35$ , mean confidence of 0.88 and 0.73, respectively). Across the entire dataset, both scores are highly correlated ( $\rho = 0.52$ ), suggesting to exploit them jointly to detect and remove predictions of low confidence in a possible application.

## 5 Discussion

Ultrasound as an established diagnosis tool that is both safe and highly available constitutes a method with potentially huge impact that has nevertheless been neglected by the machine learning community. This work presents methods and analyses that pave the way towards computer vision-assisted differential diagnosis of COVID-19 from US, providing an extensive analysis of (interpretable) methods that are relevant not only in the context of COVID-19, but in general for the diagnosis of viral and bacterial pneumonia.

We provide strong evidence that automatic detection of COVID-19 is a promising future endeavour and competitive compared to CT and CXR based models, with a sensitivity of 98% and a specificity of 91% for COVID-19, achieved on our dataset of 106 lung US videos. In comparison, sensitivity up to 98% and specificity up to 92% was reported for CT [12, 33]. We verified our results with independent test data, studied model uncertainty and concluded a significant ability of our model



to recognize low-confidence situations. We combined our approach with the only available related work, lung US segmentation models from Roy et al. [39], and found mild performance improvement in most metrics. It however remains unclear whether this gain can be attributed to the segmentation itself or is a side-effect of the increased parametrization. Certainly, there are many approaches yet to be explored in order to improve on the results presented here, including further work on video classification, but also exploiting the higher availability of CT or X-ray scans with transfer learning or adapting generative models to complement the scarce data about COVID-19 as proposed in [32]. Furthermore, we investigated the value of interpretable methods in a quantitative manner with the implementation and validation of class activation mapping in a study involving medical experts. While the analysis provides excellent evidence for the successful detection of pathological patterns like consolidations, A-lines and effusion, it reveals problems in the model’s “focal point” (e.g. missing B-lines and sometimes highlighting muscles instead of the lung) which should be further addressed using ultrasound segmentation techniques [49].

Our published database is constantly updated and verified by medical experts researchers are invited to contribute to our initiative. We envision the proposed tool as a decision support system to accelerate diagnosis or provide a “second opinion” to increase reliability. We started to collaborate with radiologists and an intensive care unit and are currently designing a controlled, clinical study to investigate the predictive power of US for automatic detection of COVID-19, especially in comparison to CT and CXR. As a preliminary demonstration, we have built a web service (link not anonymized) where users can screen ultrasound images, querying our averaged prediction model. We aim to extend the functionality of the website in the future to offer interpretable video inference, aiming for an accessible and validated tool that enables medical doctors to draw inference from their US images with unprecedented ease, convenience and speed.

## References

- [1] W. Abdalla, M. Elgendy, A. Abdelaziz, and M. Ammar. Lung ultrasound versus chest radiography for the diagnosis of pneumothorax in critically ill patients: a prospective, single-blind study. *Saudi journal of anaesthesia*, 10(3):265, 2016.
- [2] T. Ai, Z. Yang, H. Hou, C. Zhan, C. Chen, W. Lv, Q. Tao, Z. Sun, and L. Xia. Correlation of chest ct and rt-pcr testing in coronavirus disease 2019 (covid-19) in china: a report of 1014 cases. *Radiology*, page 200642, 2020.
- [3] Y. Amatya, J. Rupp, F. M. Russell, J. Saunders, B. Bales, and D. R. House. Diagnostic use of lung ultrasound compared to chest radiograph for suspected pneumonia in a resource-limited setting. *International journal of emergency medicine*, 11(1):8, 2018.
- [4] M. S. Ayhan and P. Berens. Test-time data augmentation for estimation of heteroscedastic aleatoric uncertainty in deep neural networks. *MIDL 2018 Conference*, 2018.
- [5] C. Bao, X. Liu, H. Zhang, Y. Li, and J. Liu. Covid-19 computed tomography findings: A systematic review and meta-analysis. *Journal of the American College of Radiology*, 2020.
- [6] J. Born, D. Beymer, D. Rajan, A. Coy, V. V. Mukherjee, M. Manica, P. Prasanna, D. Ballah, P. L. Shah, E. Karteris, J. L. Robertus, M. Gabrani, and M. Rosen-Zvi. On the role of artificial intelligence in medical imaging of covid-19. *medRxiv*, 2020. doi: 10.1101/2020.09.02.20187096.
- [7] J. Born, G. Brändle, M. Cossio, M. Disdier, J. Goulet, J. Roulin, and N. Wiedemann. Pocovid-net: Automatic detection of covid-19 from a new lung ultrasound imaging dataset (pocus). *arXiv preprint arXiv:2004.12084*, 2020.
- [8] J.-E. Bourcier, J. Paquet, M. Seinger, E. Gallard, J.-P. Redonnet, F. Cheddadi, D. Garnier, J.-M. Bourgeois, and T. Geeraerts. Performance comparison of lung ultrasound and chest x-ray for the diagnosis of pneumonia in the ed. *The American journal of emergency medicine*, 32(2):115–118, 2014.
- [9] J.-E. Bourcier, S. Braga, and D. Garnier. Lung ultrasound will soon replace chest radiography in the diagnosis of acute community-acquired pneumonia. *Current infectious disease reports*, 18(12):43, 2016.
- [10] E. Brogi, E. Bignami, A. Sidoti, M. Shawar, L. Gargani, L. Vetrugno, G. Volpicelli, and F. Forfori. Could the use of bedside lung ultrasound reduce the number of chest x-rays in the intensive care unit? *Cardiovascular ultrasound*, 15(1):23, 2017.

- [11] D. Buonsenso, D. Pata, and A. Chiaretti. Covid-19 outbreak: less stethoscope, more ultrasound. *The Lancet Respiratory Medicine*, 2020.
- [12] C. Butt, J. Gill, D. Chun, and B. A. Babu. Deep learning system to screen coronavirus disease 2019 pneumonia. *Applied Intelligence*, page 1, 2020.
- [13] M. Castillo. The industry of ct scanning, 2012.
- [14] M. A. Chavez, N. Shams, L. E. Ellington, N. Naithani, R. H. Gilman, M. C. Steinhoff, M. Santosham, R. E. Black, C. Price, M. Gross, et al. Lung ultrasound for the diagnosis of pneumonia in adults: a systematic review and meta-analysis. *Respiratory research*, 15(1):50, 2014.
- [15] A.-S. Claes, P. Clapuyt, R. Menten, N. Michoux, and D. Dumitriu. Performance of chest ultrasound in pediatric pneumonia. *European journal of radiology*, 88:82–87, 2017.
- [16] D. Dong, Z. Tang, S. Wang, H. Hui, L. Gong, Y. Lu, Z. Xue, H. Liao, F. Chen, F. Yang, et al. The role of imaging in the detection and management of covid-19: a review. *IEEE Reviews in Biomedical Engineering*, 2020.
- [17] L. E. Ellington, R. H. Gilman, M. A. Chavez, F. Pervaiz, J. Marin-Concha, P. Compen-Chang, S. Riedel, S. J. Rodriguez, C. Gaydos, J. Hardick, et al. Lung ultrasound as a diagnostic tool for radiographically-confirmed pneumonia in low resource settings. *Respiratory medicine*, 128:57–64, 2017.
- [18] Y. Fang, H. Zhang, J. Xie, M. Lin, L. Ying, P. Pang, and W. Ji. Sensitivity of chest ct for covid-19: comparison to rt-pcr. *Radiology*, page 200432, 2020.
- [19] M. Fiala. Ultrasound in covid-19: a timeline of ultrasound findings in relation to ct. *Clinical Radiology*, 2020.
- [20] Y. Gal and Z. Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *international conference on machine learning*, pages 1050–1059, 2016.
- [21] A. Gretton, K. M. Borgwardt, M. J. Rasch, B. Schölkopf, and A. Smola. A kernel two-sample test. *Journal of Machine Learning Research*, 13(Mar):723–773, 2012.
- [22] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [23] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017.
- [24] R. Kalkreuth and P. Kaufmann. Covid-19: A survey on public medical imaging data resources. *arXiv preprint arXiv:2004.04569*, 2020.
- [25] J. P. Kanne, B. P. Little, J. H. Chung, B. M. Elicker, and L. H. Ketai. Essentials for radiologists on covid-19: an update—radiology scientific expert panel, 2020.
- [26] L. M. Kucirka, S. A. Lauer, O. Laeyendecker, D. Boon, and J. Lessler. Variation in false-negative rate of reverse transcriptase polymerase chain reaction–based sars-cov-2 tests by time since exposure. *Annals of Internal Medicine*, 2020.
- [27] G. Lepri, M. Orlandi, C. Lazzeri, C. Bruni, M. Hughes, M. Bonizzoli, Y. Wang, A. Peris, and M. Matucci-Cerinic. The emerging role of lung ultrasound in covid-19 pneumonia. *Eur J Rheumatol*, pages 10–5152, 2020.
- [28] L. Li, L. Qin, Z. Xu, Y. Yin, X. Wang, B. Kong, J. Bai, Y. Lu, Z. Fang, Q. Song, et al. Artificial intelligence distinguishes covid-19 from community acquired pneumonia on chest ct. *Radiology*, page 200905, 2020.
- [29] Q. Li, X. Guan, P. Wu, X. Wang, L. Zhou, Y. Tong, R. Ren, K. S. Leung, E. H. Lau, J. Y. Wong, et al. Early transmission dynamics in wuhan, china, of novel coronavirus–infected pneumonia. *New England Journal of Medicine*, 2020.
- [30] D. Lichtenstein, I. Goldstein, E. Mourgeon, P. Cluzel, P. Grenier, J.-J. Rouby, et al. Comparative diagnostic performances of auscultation, chest radiography, and lung ultrasonography in acute respiratory distress syndrome. *Anesthesiology-Philadelphia then Hagerstown*, 100(1):9–15, 2004.
- [31] D. A. Lichtenstein. *Lung ultrasound in the critically ill: the BLUE protocol*. Springer, 2015.

- [32] M. Loey, F. Smarandache, and N. E. M. Khalifa. Within the lack of covid-19 benchmark dataset: A novel gan with deep transfer learning for corona-virus detection in chest x-ray images. 2020.
- [33] X. Mei, H.-C. Lee, K.-y. Diao, M. Huang, B. Lin, C. Liu, Z. Xie, Y. Ma, P. M. Robson, M. Chung, et al. Artificial intelligence-enabled rapid diagnosis of patients with covid-19. *Nature Medicine*, pages 1–5, 2020.
- [34] M. Mossa-Basha, C. C. Meltzer, D. C. Kim, M. J. Tuite, K. P. Kolli, and B. S. Tan. Radiology department preparedness for covid-19: radiology scientific expert panel. *Radiology*, page 200988, 2020.
- [35] M.-Y. Ng, E. Y. Lee, J. Yang, F. Yang, X. Li, H. Wang, M. M.-s. Lui, C. S.-Y. Lo, B. Leung, P.-L. Khong, et al. Imaging profile of the covid-19 infection: radiologic findings and literature review. *Radiology: Cardiothoracic Imaging*, 2(1):e200034, 2020.
- [36] A. Pagano, F. G. Numis, G. Visone, C. Pirozzi, M. Masarone, M. Olibet, R. Nasti, F. Schiraldi, and F. Paladino. Lung ultrasound for diagnosis of pneumonia in emergency department. *Internal and emergency medicine*, 10(7):851–854, 2015.
- [37] Q.-Y. Peng, X.-T. Wang, L.-N. Zhang, C. C. C. U. S. Group, et al. Findings of lung ultrasonography of novel corona virus pneumonia during the 2019–2020 epidemic. *Intensive Care Medicine*, page 1, 2020.
- [38] F. Reali, G. F. S. Papa, P. Carlucci, P. Fracasso, F. Di Marco, M. Mandelli, S. Soldi, E. Riva, and S. Centanni. Can lung ultrasound replace chest radiography for the diagnosis of pneumonia in hospitalized children? *Respiration*, 88(2):112–115, 2014.
- [39] S. Roy, W. Menapace, S. Oei, B. Luijten, E. Fini, C. Saltori, I. Huijben, N. Chennakeshava, F. Mento, A. Sentelli, et al. Deep learning for classification and localization of covid-19 markers in point-of-care lung ultrasound. *IEEE Transactions on Medical Imaging*, 2020.
- [40] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision*, pages 618–626, 2017.
- [41] F. Shan+, Y. Gao+, J. Wang, W. Shi, N. Shi, M. Han, Z. Xue, D. Shen, and Y. Shi. Lung infection quantification of covid-19 in ct images with deep learning. *arXiv preprint arXiv:2003.04655*, 2020.
- [42] F. Shi, J. Wang, J. Shi, Z. Wu, Q. Wang, Z. Tang, K. He, Y. Shi, and D. Shen. Review of artificial intelligence techniques in imaging data acquisition, segmentation and diagnosis for covid-19. *arXiv preprint arXiv:2004.02731*, 2020.
- [43] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In *International Conference on Learning Representations*, 2015.
- [44] M. Smith, S. Hayward, S. Innes, and A. Miller. Point-of-care lung ultrasound in patients with covid-19—a narrative review. *Anaesthesia*, 2020.
- [45] K. A. Stewart, S. M. Navarro, S. Kambala, G. Tan, R. Poondla, S. Lederman, K. Barbour, and C. Lavy. Trends in ultrasound use in low and middle income countries: A systematic review. *International Journal*, 9(1):103–120, 2020.
- [46] L. R. Sultan and C. M. Sehgal. A review of early experience in lung ultrasound (lus) in the diagnosis and management of covid-19. *Ultrasound in Medicine & Biology*, 2020.
- [47] L. Tutino, G. Cianchi, F. Barbani, S. Batacchi, R. Cammelli, and A. Peris. Time needed to achieve completeness and accuracy in bedside lung ultrasound reporting in intensive care unit. *Scandinavian journal of trauma, resuscitation and emergency medicine*, 18(1):44, 2010.
- [48] A. Ulhaq, A. Khan, D. Gomes, and M. Pau. Computer vision for covid-19 control: A survey. *arXiv preprint arXiv:2004.09420*, 2020.
- [49] R. J. van Sloun and L. Demi. Localizing b-lines in lung ultrasonography by weakly-supervised deep learning, in-vivo results. *IEEE Journal of Biomedical and Health Informatics*, 2019.
- [50] L. Wang and A. Wong. Covid-net: A tailored deep convolutional neural network design for detection of covid-19 cases from chest radiography images. *arXiv preprint arXiv:2003.09871*, 2020.
- [51] M. Weinstock, A. Echenique, J. Russell, et al. Chest x-ray findings in 636 ambulatory patients with covid-19 presenting to an urgent care center: a normal chest x-ray is no guarantee. *J Urgent Care Med*, 14(7):13–18, 2020.

- [52] W. Yang and F. Yan. Patients with rt-pcr-confirmed covid-19 and normal chest ct. *Radiology*, 295(2): E3–E3, 2020.
- [53] B. Zhou, A. Khosla, A. Lapedriza, A. Oliva, and A. Torralba. Learning deep features for discriminative localization. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2921–2929, 2016.
- [54] Z. Zhou, V. Sodha, M. M. Rahman Siddiquee, R. Feng, N. Tajbakhsh, M. B. Gotway, and J. Liang. Models genesis: Generic autodidactic models for 3d medical image analysis. In *Medical Image Computing and Computer Assisted Intervention – MICCAI 2019*, pages 384–393, Cham, 2019. Springer International Publishing. ISBN 978-3-030-32251-9. URL [https://link.springer.com/chapter/10.1007/978-3-030-32251-9\\_42](https://link.springer.com/chapter/10.1007/978-3-030-32251-9_42).
- [55] B. Zoph, V. Vasudevan, J. Shlens, and Q. V. Le. Learning transferable architectures for scalable image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8697–8710, 2018.

## A Appendix

### A.1 Model architectures and hyperparameter

As a base, we use the convolutional part of the established VGG-16 [43], pre-trained on Imagenet. The model we call VGG is followed by one hidden layer of 64 neurons with ReLU activation, dropout of 0.5, batch normalization and the output layer with softmax activation. The CAMs for this model were computed with Grad-CAM [40]. To compare Grad-CAMs with regular CAMs [53], we also tested VGG-CAM, a CAM-compatible VGG with a single dense layer following the global average pooling after the last convolutional layer. For both models, during training only the weights of the last three layers were fine-tuned, while the other ones were frozen to the values from pre-training. This results in a total of  $\sim 2.4\text{M}$  trainable and  $\sim 12.4\text{M}$  non-trainable parameters. The model is trained with a cross entropy loss function on the softmax outputs, and optimized with Adam with an initial learning rate of  $1e-4$ . All models were implemented in TensorFlow and trained for 40 epochs with a batch size of 8 and early stopping was enabled.

### A.2 Pretrained segmentation models

Figure 5 gives an example for the segmented ultrasound image with the model from Roy et al. [39]. In our work the segmented image serves as input to the VGG-Segment model.

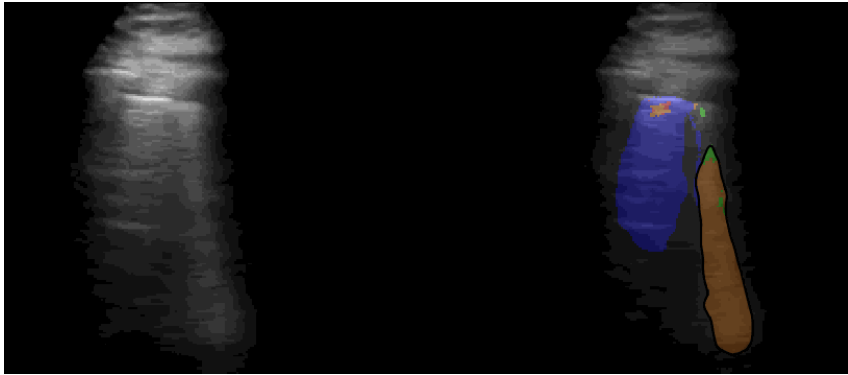


Figure 5: **Example snapshot from lung segmentation of COVID-19 patient.** Left side shows the raw US recording and the right side shows the segmentation method from Roy et al. [39] highlighting the B-line. The images shown on the right were used as input for the VGG-Segment model.

### A.3 Uncertainty estimation

For both aleatoric and epistemic uncertainty, the confidence estimate  $c_i$  of sample  $i$  is computed by scaling the sample’s standard deviation to  $\in [0, 1]$  and interpreting it as an inverse precision:

$$c_i = -\left(\frac{\sigma_{i,j} - \sigma_{min}}{\sigma_{max} - \sigma_{min}}\right) + 1, \tag{1}$$

where  $\sigma_{i,j}$  is the sample standard deviation of the ten class probabilities of the winning class  $j$ ,  $\sigma_{min}$  is the minimal standard deviation (0, i.e. all probabilities for the winning class are identical) and  $\sigma_{max}$  is the maximal standard deviation, i.e. 0.5. Practically, for epistemic uncertainty, dropout was set to 0.5 across the VGG model and for aleatoric uncertainty the same transformations as during training are employed.

### A.4 Results

Re-formulating the classification as a binary task, the ROC-curve and precision-recall curves can be computed for each class. Figure 6 and Figure 7 depict the performance per class, comparing all proposed models. While pneumonia is distinguished well by all models, NASNet has difficulties with the correct classification of normal lung images. Figure 7a and Figure 7b show that COVID-19 is predicted better than healthy lung images, but not as distinct as pneumonia infections.

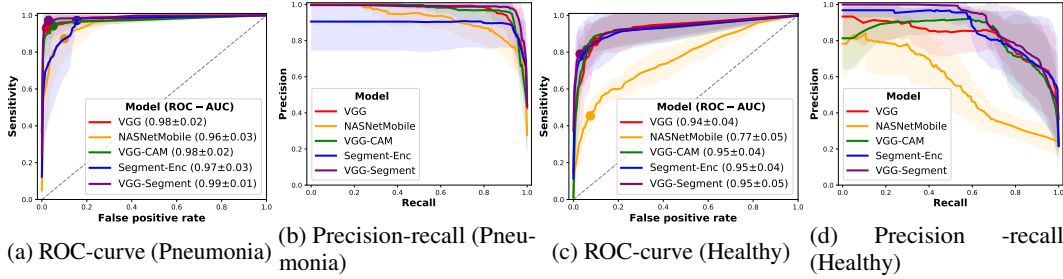


Figure 6: **Binary classification results.** All models achieve good precision and recall in pneumonia detection, but lower scores and higher variances are observed for data of healthy patients

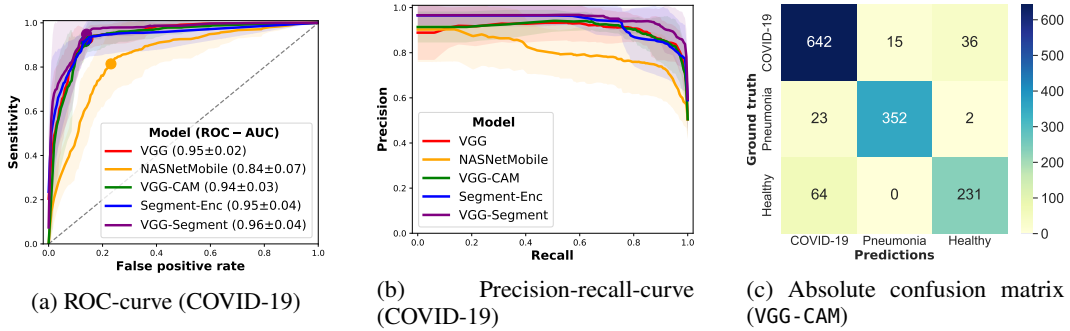


Figure 7: **COVID-19 detection and absolute confusion matrix.**

Furthermore, in addition to the normalized confusion matrices we provide the absolute values here in Figure 7c (referring to VGG-CAM). Note that most of our data shows COVID-19 infected lungs, despite the novelty of the disease. Problematically, healthy and COVID-19 patients are confused in 100 images, whereas bacterial pneumonia is predicted rather reliably.

#### A.4.1 Uninformative class

Although the main task is defined as differentiating COVID-19, bacterial pneumonia and healthy, we trained the model actually with a fourth "uninformative" class in order to identify out-of-distribution samples. This concerns both entirely different pictures (no ultrasound), as well as ultrasound images not showing the lung. Thus, we added 200 images from Tiny ImageNet (one per class taken from the test set) together with 200 neck ultrasound scans taken from the Kaggle ultrasound nerve segmentation challenge. Note that the latter is data recorded with linear ultrasound probes, leading to very different ultrasound images.

Table 4 lists the results including these uninformative samples, where better accuracy is achieved due to the ease of distinguishing the uninformative samples from other data. In all cases, precision and recall are higher than 0.98 with low standard deviation.

	Class	Recall	Precision	F1-score	Specificity	MCC
<b>VGG</b> Acc.: 0.92, bal.: 0.91 Par.: 14 747 971	COVID-19	0.89 ± 0.06	0.91 ± 0.05	0.9 ± 0.03	0.95 ± 0.03	0.84 ± 0.04
	Pneumonia	0.94 ± 0.05	0.92 ± 0.06	0.93 ± 0.05	0.98 ± 0.02	0.91 ± 0.06
	Healthy	0.85 ± 0.11	0.83 ± 0.09	0.83 ± 0.07	0.96 ± 0.02	0.81 ± 0.07
	Uninformative	0.99 ± 0.01	1.0 ± 0.01	0.99 ± 0.01	1.0 ± 0.0	0.99 ± 0.01
<b>VGG-CAM</b> Acc.: 0.9, bal.: 0.88 Par.: 14 716 227	COVID-19	0.93 ± 0.05	0.87 ± 0.07	0.9 ± 0.05	0.92 ± 0.04	0.83 ± 0.07
	Pneumonia	0.93 ± 0.05	0.95 ± 0.06	0.94 ± 0.05	0.99 ± 0.01	0.92 ± 0.06
	Healthy	0.78 ± 0.1	0.86 ± 0.08	0.81 ± 0.05	0.97 ± 0.02	0.78 ± 0.05
	Uninformative	0.99 ± 0.01	0.99 ± 0.01	0.99 ± 0.01	1.0 ± 0.0	0.99 ± 0.01
<b>NASNetMobile</b> Acc.: 0.81, bal.: 0.78 Par.: 4 814 487	COVID-19	0.87 ± 0.1	0.74 ± 0.11	0.8 ± 0.1	0.82 ± 0.05	0.67 ± 0.13
	Pneumonia	0.79 ± 0.15	0.88 ± 0.08	0.83 ± 0.1	0.97 ± 0.02	0.79 ± 0.13
	Healthy	0.47 ± 0.03	0.61 ± 0.13	0.53 ± 0.05	0.94 ± 0.03	0.46 ± 0.07
	Uninformative	0.99 ± 0.01	1.0 ± 0.0	0.99 ± 0.01	1.0 ± 0.0	0.99 ± 0.01
<b>VGG-Segment</b> Acc.: 0.93, bal.: 0.91 Par.: 34 018 074	COVID-19	0.96 ± 0.05	0.89 ± 0.06	0.92 ± 0.04	0.92 ± 0.04	0.87 ± 0.06
	Pneumonia	0.96 ± 0.03	0.95 ± 0.03	0.95 ± 0.02	0.98 ± 0.01	0.94 ± 0.03
	Healthy	0.77 ± 0.14	0.91 ± 0.08	0.82 ± 0.08	0.98 ± 0.02	0.8 ± 0.08
	Uninformative	0.97 ± 0.03	1.0 ± 0.0	0.99 ± 0.01	1.0 ± 0.0	0.98 ± 0.02
<b>Segment-Enc</b> Acc.: 0.92, bal.: 0.91 Par.: 19 993 307	COVID-19	0.92 ± 0.09	0.91 ± 0.06	0.91 ± 0.03	0.94 ± 0.04	0.86 ± 0.04
	Pneumonia	0.95 ± 0.04	0.89 ± 0.12	0.92 ± 0.07	0.96 ± 0.04	0.9 ± 0.08
	Healthy	0.79 ± 0.17	0.89 ± 0.1	0.82 ± 0.11	0.98 ± 0.01	0.81 ± 0.12
	Uninformative	1.0 ± 0.0	1.0 ± 0.0	1.0 ± 0.0	1.0 ± 0.0	1.0 ± 0.0

Table 4: **Performance comparison.** Acc. abbreviates accuracy, Bal. balanced accuracy and Par. the number of parameters. The raw results are listed, including the uninformative class. Clearly, this fourth class is very distinctive and is learnt successfully, with almost all scores above 0.89

## A.5 Class activation maps

In addition to the scatter plot in Figure 3 we present the corresponding density plot in Figure 8, showing the area of the ultrasound image where the class activation is maximal for each class. It can be observed that the activation on healthy and COVID-19 videos is located further in the upper part of the image, where usually only muscles and skin are observed. Further work is thus necessary to analyze and improve the qualitative results of the model.

However, with respect to pathological patterns visible, the model does in many cases focus on the patterns that are interesting to medical experts. Table 5 breaks down the results presented in Figure 4 more in detail, and in particular separately for both medical experts. Note that with respect to the pleural line, we only consider the opinion of expert 2 since expert 1 did not mention it. With the exception of consolidations, the difference in responses is quite large, which is however unsurprising for such a qualitative task. Besides the patterns that were already named in Figure 4, the heatmaps also correctly highlighted air bronchograms (2 cases according to expert 1) and a pleural effusion in 1 out of 7 cases.

	Consolidations	A-lines	B-lines	Bronchogram	Effusion	Pleural line
<b>Specific for</b>	Bacterial pne.	Healthy	COVID-19, viral pne.	Bacterial pne.	Pne.	Pne. if irregular
<b>Total visible (expert 1)</b>	18	13	12	2	7	20 (expert 2)
<b>CAM highlighted (expert 1)</b>	17	6	0	2	1	0
<b>CAM highlighted (expert 2)</b>	17	10	6	0	0	9

Table 5: Pathological patterns visible and highlighted by class activation maps of our model. Pneu abbreviated pneumonia. The model focuses on consolidations, A-lines and the pleural line.

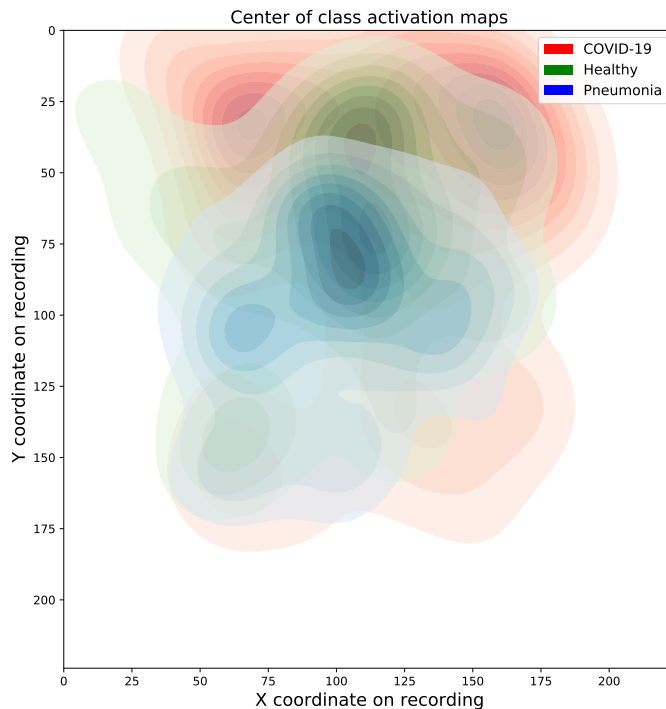


Figure 8: **Density plot of centers of class activation maps.** Pneumonia-CAMs are rather centralized compared to other CAMs. Problematically, COVID-CAMs seem to exhibit a tendency for upper regions of the probe that do not necessarily belong to the lung.

### A.6 Maximum mean discrepancy analysis

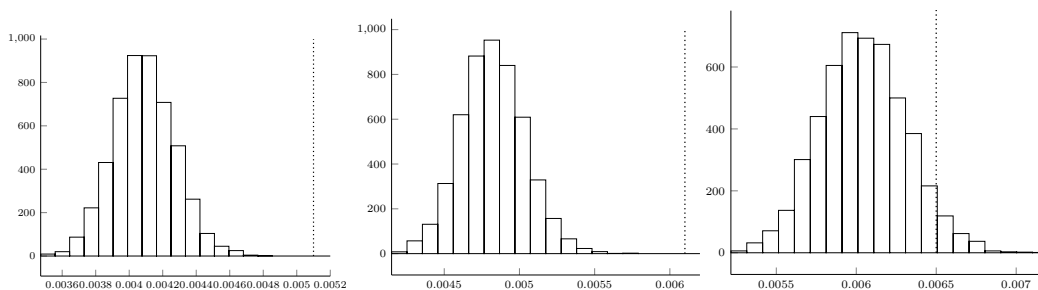


Figure 9: Histograms depicting the empirical null distribution, obtained via bootstrapping 1000 samples, of the MMD values (from left to right)  $MMD(\mathbf{C}, \mathbf{P}) \approx 0.0051$ ,  $MMD(\mathbf{C}, \mathbf{H}) \approx 0.0061$ , and  $MMD(\mathbf{P}, \mathbf{H}) \approx 0.0065$ , respectively. The corresponding true MMD values, i.e. the ones we obtain by looking at the labels, is indicated as a dashed line in each histogram. We observe that these values are highly infrequent under the null distribution, indicating that the differences between the three classes are significant. Notably, the statistical distance between patients suffering from bacterial pneumonia and healthy patients (rightmost histogram) achieves a slightly lower empirical significance of  $\approx 0.04$ . We speculate that this might be related to *other* pre-existing conditions in healthy patients that are not pertinent to this study, though.