

Quantification of pulmonary involvement in COVID-19 pneumonia by means of a cascade of two U-nets: training and assessment on multiple datasets using different annotation criteria

Francesca Lizzi ^{1,2}, Abramo Agosti ⁶, Francesca Brero ^{4,5}, Raffaella Fiamma Cabini ^{4,6}, Maria Evelina Fantacci ^{2,3}, Silvia Figini ^{4,11}, Alessandro Lascialfari ^{4,5}, Francesco Laruina ^{1,2}, Piernicola Oliva ^{8,9}, Stefano Piffer ^{7,10}, Ian Postuma ⁴, Lisa Rinaldi ^{4,5}, Cinzia Talamonti ^{7,10}, Alessandra Retico ²

¹ Scuola Normale Superiore, Pisa;

² National Institute of Nuclear Physics (INFN), Pisa division, Pisa, IT

³ Department of Physics, University of Pisa, Pisa, IT;

⁴ INFN, Pavia division, Pavia, IT;

⁵ Department of Physics, University of Pavia, Pavia, IT;

⁶ Department of Mathematics, University of Pavia, Pavia, IT;

⁷ Department of Biomedical Experimental Clinical Science "M. Serio", University of Florence, Florence, IT

⁸ Department of Chemistry and Pharmacy, University of Sassari, Sassari, IT.

⁹ INFN, Cagliari division, Cagliari, IT

¹⁰ INFN, Florence division, Florence, IT

¹¹ Department of Social and Political Science, University of Pavia, Pavia, IT;

May 7, 2021

Abstract

Purpose The automatic assignment of a severity score to the CT scans of patients affected by COVID-19 pneumonia could reduce the workload in radiology departments. This study aims at exploiting Artificial intelligence (AI) for the identification, segmentation and quantification of COVID-19 pulmonary lesions. The limited data availability and the annotation quality are relevant factors in training AI-methods. We investigated the effects of using multiple datasets, heterogeneously populated and annotated according to different criteria.

Methods We developed an automated analysis pipeline, the *LungQuant* system, based on a cascade of two U-nets. The first one (U-net₁) is devoted to the identification of the lung parenchyma, the second one (U-net₂) acts on a bounding box enclosing the segmented lungs to identify the areas affected by COVID-19 lesions. Different public datasets were used to train

the U-nets and to evaluate their segmentation performances, which have been quantified in terms of the Dice index. The accuracy in predicting the CT-Severity Score (CT-SS) of the *LungQuant* system has been also evaluated.

Results Both Dice and accuracy showed a dependency on the quality of annotations of the available data samples. On an independent and publicly available benchmark dataset (COVID-19-CT-Seg), the Dice values measured between the masks predicted by *LungQuant* system and the reference ones were 0.95 ± 0.01 and 0.66 ± 0.13 for the segmentation of lungs and COVID-19 lesions, respectively. The accuracy of 90% in the identification of the CT-SS on this benchmark dataset was achieved.

Conclusion We analysed the impact of using data samples with different annotation criteria in training an AI-based quantification system for pulmonary involvement in COVID-19 pneumonia. In terms of the Dice index, the U-net segmentation quality strongly depends on the quality of the lesion annotations. Nevertheless, the CT-SS can be accurately predicted on independent validation sets, demonstrating the satisfactory generalization ability of the *LungQuant*.

Keywords COVID-19, Chest Computed Tomography, Ground-glass opacities, Segmentation, Machine Learning, U-net

1 Introduction

The task of segmenting the abnormalities of the lung parenchyma related to COVID-19 infection is a typical segmentation problem that can be addressed with methods based on DL. CT findings of patients with COVID-19 infection may include bilateral distribution of ground-glass opacifications (GGO), consolidations, crazy-paving patterns, reversed halo sign and vascular enlargement [2]. Due to the extremely heterogeneous appearance of COVID-19 lesions in density, textural pattern, global shape and location in the lung, an analytical approach is definitely hard to code, whereas it is preferable to learn directly from examples. The potential of DL-based segmentation approaches is particularly suited in this case, provided that a sufficient number of annotated examples are available for training the models.

Few fully automated software tools devoted to this task have been recently proposed [8, 4, 9]. Lessmann *et al.* [8] developed a U-net model for lesion segmentation trained on semi-automatically annotated COVID-19 cases. The output of this system was then combined with the lung lobe segmentation algorithm reported in Xie *et al.* [12]. The approach proposed in Fang *et al.* [4] implements the automated lung segmentation method provided in the work of Hofmanninger *et al.* [6], together with a lesion segmentation strategy based on multiscale feature extraction [5]. The specific problem related to the development of fully automated DL-based segmentation strategies with limited annotated data samples has been explicitly tackled by Ma *et al.* [9]. The authors studied how to train and evaluate a DL-based system for lung and COVID-19 lesion segmentation on a poorly populated samples of CT scans. They also made the data collected

Table 1: A summary of the datasets used in this study. The CT Severity Score (CT-SS) information is not available for all datasets, but it can be computed for data which has both lung masks and ground-glass opacification (GGO) masks.

Dataset name	Lung mask	GGO mask	CT-SS	N. of cases
Plethora [7]	Yes	No	No	402
Lung CT Segmentation Challenge [13]	Yes	No	No	60
COVID-19 Challenge [1]	No	Yes	No	199
MosMed [10]	No	No	Inferable	1110
MosMed (annotated subsample)	No	Yes	Inferable	50
MosMed (in-house annotated subsample)	Yes	No	Inferable	91
COVID-19-CT-Seg [9]	Yes	Yes	Inferable	10

for their experiment publicly available and tested their algorithm on a public dataset, allowing for a fair comparison with their system.

In this work we present the DL-based fully automated system to segment both lungs and lesions associated with COVID-19 pneumonia, the *LungQuant* system, which provides the percentage of lung volume compromised by the infection. We extended the study proposed by Ma *et al.* [9] in two main directions: 1) we investigated the impact of the training annotation style on the prediction across different datasets; 2) we translated the segmentation problem into a CT-SS assessment problem and we evaluated the reliability of the automated CT-SS assignment across different data samples. This paper is structured as follows: we list all the publicly accessible data samples we used to develop and validate the *LungQuant* system; then, we describe the image analysis pipeline we set up and the training and cross-validation strategies we adopted; finally, we show and discuss the quantification performance obtained either against a voxel-wise ground truth or in terms of the CT severity scores, according to the information available for each data sample.

2 Material and Methods

2.1 Datasets

Five public available datasets have been used to train and evaluate our segmentation pipeline. Most of them include image annotations, but each annotation has been associated to patients using different criteria, which are described in the Supplementary Materials. In Table 1, a summary of available labels for each dataset is reported.

2.2 *LungQuant*: the Deep-Learning based quantification analysis pipeline

In this section, we describe the fully-automated pipeline we set up for the quantification of lung involvement in patients affected by COVID-19 pneumonia. The analysis pipeline, which is referred to in what follows as the *LungQuant* system, provides in output the percentage P of lung volume affected by COVID-19 lesions and the corresponding CT severity score (CT-SS=1 for $P < 5\%$, CT-SS=2 for $5\% \leq P < 25\%$, CT-SS=3 for $25\% \leq P < 50\%$, CT-SS=4 for $50\% \leq P < 75\%$, CT-SS=5 for $P \geq 75\%$).

A summary of our image analysis pipeline is reported in Fig. 1. The central analysis module is a U-net for image segmentation [11] (see sec. 2.2.1), which is implemented in a cascade of two different U-nets: the first network, U-net₁, is trained to segment the lung and the second one, U-net₂, is trained to segment the COVID lesions in the CT scans. In the following sections, the whole process is described step by step.

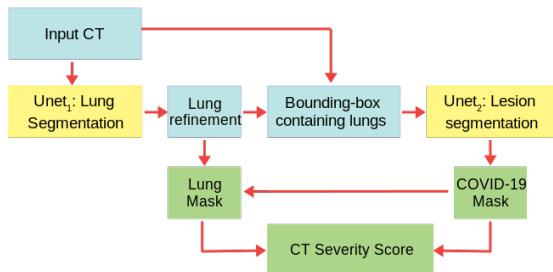


Figure 1: A summary of the whole analysis pipeline: the input CT scans are used to train U-net₁, which is devoted to lung segmentation; its output is refined by a morphology-based method. A bounding box containing the segmented lungs is made and applied to all CT scans for training U-net₂, which is devoted to COVID-19 lesion segmentation. Finally, the output of U-net₂ is the definitive COVID-19 lesion mask, whereas the definitive lung mask is obtained as the union between the outputs of U-net₁ and U-net₂. The ratio between the COVID-19 lesion mask and the lung mask provides the CT-SS for each patient.

2.2.1 U-net

For both lung and COVID-19 lesion segmentation, we implemented a fully automated method inspired by the U-net, the fully-convolutional neural networks for image segmentation developed by Ronneberger *et al.* [11]. We implemented a U-net using Keras [3], a Python deep-learning API that uses Tensorflow as backend. In Figure 2 a simplified scheme of our U-net is reported.

Each block of layers in the compression path (left) is made by 3 convolutional layers, ReLu activation functions and instance normalization layers. The input of each block is added to the block output in order to implement a residual

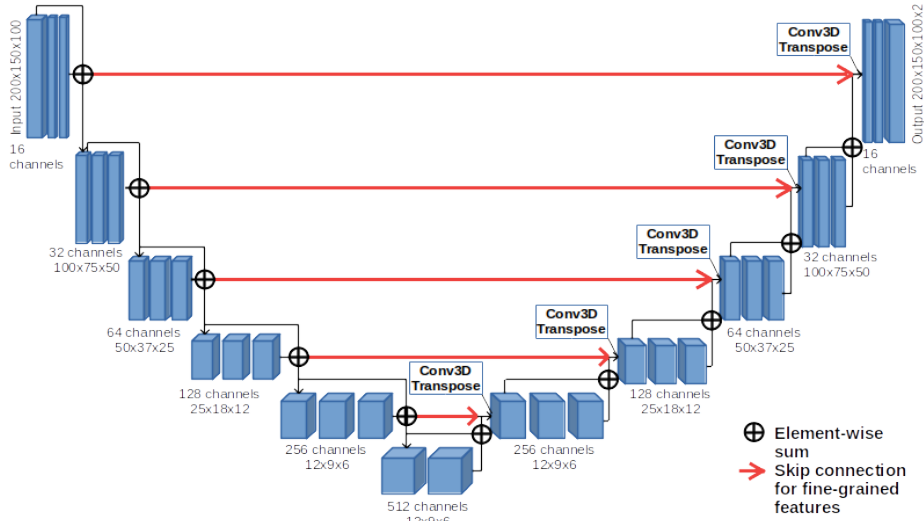


Figure 2: U-net scheme: the neural network is made of 6 levels of depth. In the compression path (left), the input is processed through convolutions, activation layers (ReLU) and instance normalization layers, while in the decompression one (right), in addition to those already mentioned, 3D Transpose Convolution (de-convolution) layers are also introduced.

connection. In the decompression path (right), one convolutional layer has been replaced by a de-convolutional layer to upsample the images to the input size. In the last layer of the U-nets, a softmax is applied to the final feature map and then the loss is computed.

2.2.2 The U-net cascade for lesion quantification and severity score assignment

We started by training U-net₁, which is devoted to lung segmentation, using the three datasets containing original CT scans and lung masks (see Table. 1). The input CT scans, whose number of slices is highly variable, are resampled to matrices of 200x150x100 voxels. The output of U-net₁ was then refined using a connected-component labeling strategy, which helps to remove small regions of the segmented mask not connected with the main objects identified as the lungs (see Supplementary Materials for further details). We then built for each CT a bounding box enclosing the morphologically refined segmented lungs, adding a conservative padding of 2.5 cm. The bounding boxes were used to crop the training images for U-net₂, which has the same architecture as U-net₁. The cropped images were resized to a matrix of 200x150x100 voxels. We applied a windowing on the grey-level values of the CT scans to optimize the image contrast for the two segmentation problems we focused on in this analysis. In particular, we selected the [-1000, 1000] HU window range for the U-net₁ and

the [-1000, 300] HU range for U-net₂. The first window highlights the contrast between the lung parenchyma and the surrounding tissues, whereas the second one enhances the heterogeneous structure of the lung abnormalities related to the COVID-19 infection. We implemented a data augmentation strategy, relying on the most commonly used data augmentation techniques for DL (see Supplementary Materials for further details) to overcome the problem of having a limited amount of labelled data.

The *LungQuant* system returns the infection mask as the output of U-net₂ and the lung mask as the union between the output of U-net₁ and U-net₂. This choice has been made *a priori* by design, as U-net₁ has been trained to segment the lungs relying on the available annotated data, which are almost totally of patients not affected by COVID-19 pneumonia. Thus, U-net₁ is expected to be unable to accurately segment the areas affected by GGO or consolidations; as also these areas are part of the lungs, they should be instead included in the mask. Training U-net₂ to recognize the COVID-19 lesions on a conservative bounding box containing only the lungs has two main advantages: it allows to restrict the action volume of the U-net to the region where the lung parenchyma (either normal or affected by COVID-19 lesions) is supposed to be, thus avoiding false-positive findings outside the chest; it facilitates the U-net training phase, as the dimensions of the lungs of different patients are normalized, thus the U-net learning process can be focused on the textural patterns characterizing the COVID-19 lesions.

Finally, once lung and lesion masks have been identified, the *LungQuant* system computes the percentage of lung volume affected by COVID-19 lesions as the ratio between the total number of voxels of the infection mask and the total number of voxels of the lung mask, and converts it into the corresponding CT severity score.

2.3 Training details and evaluation strategy for the U-nets

A detailed description of the metrics used and on the data augmentation strategies implemented is provided in the Supplementary Materials, whereas the data-splitting criterion adopted is described below.

2.3.1 Cross-validation strategy

To train, validate and test the performances of each of the two U-nets, we partitioned the available datasets into the training, validation and test sets, and we evaluated the network performance separately and globally. U-net₁ has been trained and evaluated on CT scans coming from three different datasets: Plethora, MosMed and LCTSC. U-net₂ has been trained and evaluated on samples made of CT scans coming from the COVID-19-Challenge dataset and from the MosMed dataset. The amount of CT scan used for train, validation and test sets for each U-net is reported in Table 2. U-net₂ has been trained twice, i.e. on both 60% and 90% of the CT scans of COVID-19-Challenge and Mosmed datasets to investigate the effect of maximizing training set size on the ability

of the system to properly segment the lesions. In the former case, $U\text{-net}_2^{60\%}$ training has been evaluated on a validation set made of 20% of cases and tested on the remaining 20%. As regard the latter, $U\text{-net}_2^{90\%}$, the remaining 10% of CT scans has been used as validation set. The trained segmentation networks ($U\text{-net}_1$ and both $U\text{-net}_2^{60\%}$ and $U\text{-net}_2^{90\%}$) have been validated on a completely independent validation set consisting of the 10 CT scans of the COVID-19-CT-Seg dataset, which is the only public available dataset containing both lung and infection mask annotations.

Table 2: Number of CT scans assigned to the train, validation (val) and test sets used during the training and performance assessment of the $U\text{-net}_1$ and the $U\text{-net}_2$ networks.

$U\text{-net}_1$	train	val	test
Plethora	319	40	40
MosMed (91 CT-0)	55	18	18
LCTSC	36	12	12
Coronacases	/	/	10
$U\text{-net}_2^{60\%}$	train (60%)	val (20%)	test
COVID-19 Challenge	119	40	40
MosMed (50 CT-1)	30	10	10
Coronacases	/	/	10
$U\text{-net}_2^{90\%}$	train (90%)	val (10%)	test
COVID-19 Challenge	179	20	/
MosMed (50 CT-1)	45	5	/
Coronacases	/	/	10

The *LungQuant* system has been set up by integrating all analysis modules, as reported in Fig. 1. In this work we built and analyzed two *LungQuant* systems, obtained by integrating alternately $U\text{-net}_2^{60\%}$ or $U\text{-net}_2^{90\%}$ into the analysis pipeline. The systems have been evaluated in terms of the ability to predict the percentage of affected lung parenchyma and CT-SS on the fully annotated COVID-19-CT-Seg dataset, which is completely independent of the system training phase.

3 Results

We report in this section, first, the performance achieved by each of the segmentation networks we trained, $U\text{-net}_1$ and $U\text{-net}_2$, then, the quantification performance of the integrated *LungQuant* system, evaluated on completely independent test sets. We trained both U-nets for 300 epochs on a NVIDIA V100 GPU using ADAM as optimizer, and we kept the models trained at the epoch where the best evaluation metric on the validation set was obtained.

3.1 U-net₁: Lung segmentation performance

U-net₁ for lung segmentation was trained using three different datasets, as specified in Table 2: the Plethora, a subsample of 91 CT-0 cases of the MosMed dataset and the 60 CT scans of the LCTSC datasets. For the MosMed dataset, as reported in Table 1, the lung mask annotations were provided by an in-house developed segmentation software (see Supplementary Materials). Out of the 254 CT scans of the CT-0 MosMed sample, the 91 CT scan we considered here are those on which the in-house segmentation algorithm provided an accurate segmentation, as judged by an experienced medical imaging data analyst. Then, we tested U-net₁ on each of the three independent test sets, and we reported in Table 3 the performance achieved in terms of Dice values computed between the segmented and the reference masks. The average Dice value obtained on all test samples is also reported. We evaluated the lung segmentation performances in three cases: 1) on CT scans and masks resized to the 200x150x100 voxel array size; 2) on CT scans and masks in the original size before undergoing the morphological refinement; 3) on CT scans and masks in the original size and after the morphological refinement. Even if segmentation refinement has a small effect on Dice score, as shown in Table 3, it is a fundamental step to allow the definition of precise bounding boxes enclosing the lungs, and thus to facilitate the U-net₂ learning process.

Table 3: Performances achieved by U-net₁ in lung segmentation on different test sets, evaluated in terms of the Dice metric at three successive stages of the segmentation procedure.

Test set	Masks of U-net size (Dice coefficient)	Masks before refinement (Dice coefficient)	Masks after refinement (Dice coefficient)
Plethora	0.96 ± 0.02	0.95 ± 0.02	0.95 ± 0.04
MosMed	0.97 ± 0.02	0.97 ± 0.02	0.97 ± 0.02
LCTSC	0.96 ± 0.03	0.95 ± 0.03	0.96 ± 0.01
Coronacases	0.96 ± 0.01	0.95 ± 0.01	0.95 ± 0.01

3.2 U-net₂: COVID-19 lesion segmentation performance

U-net₂ for COVID-19 lesion segmentation has been trained and evaluated separately on the COVID-19-Challenge dataset and on the annotated subset of the MosMed dataset, following the train/validation/test partitioning reported in Table 2. The segmentation performances achieved on the test sets are reported in terms of the Dice metric in Table 4. As reported in the table, the performances of U-net₂ were evaluated also according to a cross-sample validation scheme.

As expected, the U-net₂ performances are higher when both the training set and independent test sets belong to the same data cohort. By contrast, when a U-net₂ is trained on COVID-19-Challenge data and tested on Mosmed (and the other way around) performances significantly decrease. This effect is due to the fact that the two datasets have been collected and annotated with different

Table 4: Performances achieved by U-net₂ in COVID-19 lesion segmentation, evaluated in terms of the Dice metric. The composition of the train and test sets is reported in Table 2.

U-net	Trained on	Test set	U-net size (Dice coefficient)	Original CT size (Dice Coefficient)
U-net ₂ ^{60%}	COVID-19 Challenge	COVID-19 challenge	0.51 ± 0.24	0.51 ± 0.25
	COVID-19 Challenge	MosMed	0.39 ± 0.19	0.40 ± 0.19
	MosMed	MosMed	0.54 ± 0.22	0.55 ± 0.22
	MosMed	COVID-19 challenge	0.25 ± 0.23	0.25 ± 0.23
	COVID-19 challenge + MosMed	COVID-19 challenge + MosMed	0.49 ± 0.21	0.50 ± 0.21
U-net ₂ ^{90%}	COVID-19 Challenge + MosMed	COVID-19 Challenge + MosMed	0.64 ± 0.23	0.65 ± 0.23

criteria. We obtained a better result with the U-net₂ trained on the COVID-19 Challenge dataset and tested on the MosMed test set, since the network has been trained on a larger data sample and hence it has a higher generalization capability. When using data from both the COVID-19-Challenge and MosMed datasets in the training and test sets, the Dice index on the test set stands on 0.50 ± 0.21 , which is similar to the performance obtained in training the U-net₂ on COVID-19-Challenge data only. The best segmentation performances have been obtained by the U-net₂ trained using the 90% of the available data, U-net₂^{90%}, which reaches a Dice value of 0.65 ± 0.23 on the test set. This result suggests the need to train U-net models on the largest possible data samples in order to achieve higher segmentation performance.

3.3 Evaluation of the quantification performance of the *LungQuant* system

3.3.1 Evaluation of lung and COVID-19 lesion segmentations

Once the two U-nets have been trained and the whole analysis pipeline has been integrated in the *LungQuant* system, we tested it on an independent set (COVID-19-CT-Seg dataset) of CT scans in order to quantify the performances of the whole process.

Table 5: Performances of the *LungQuant* system on the independent COVID-19-CT-Seg test dataset. The Dice metric computed between the reference lung and lesion masks and those respectively predicted by the *LunQuant* system are reported.

<i>LungQuant</i> system	Lung segmentation (Dice coefficient)	Infection segmentation (Dice coefficient)
<i>LungQuant</i> with U-net ₂ ^{60%}	0.96 ± 0.01	0.62 ± 0.09
<i>LungQuant</i> with U-net ₂ ^{90%}	0.95 ± 0.01	0.66 ± 0.13

Figure 3 allows a visual comparison between the lung and lesion masks pro-

vided by the *LungQuant* systems integrating U-net₂^{90%} and the reference ones. Three axial slices of the first CT scan of the COVID-19-CT-Seg test dataset (*coronacases001.nii*) are shown, together with two overlays of the lung and lesion masks, respectively. A very good overlap between the predicted and reference lung masks is observable, whereas a partial overlap occurs between the predicted and reference lesion masks.

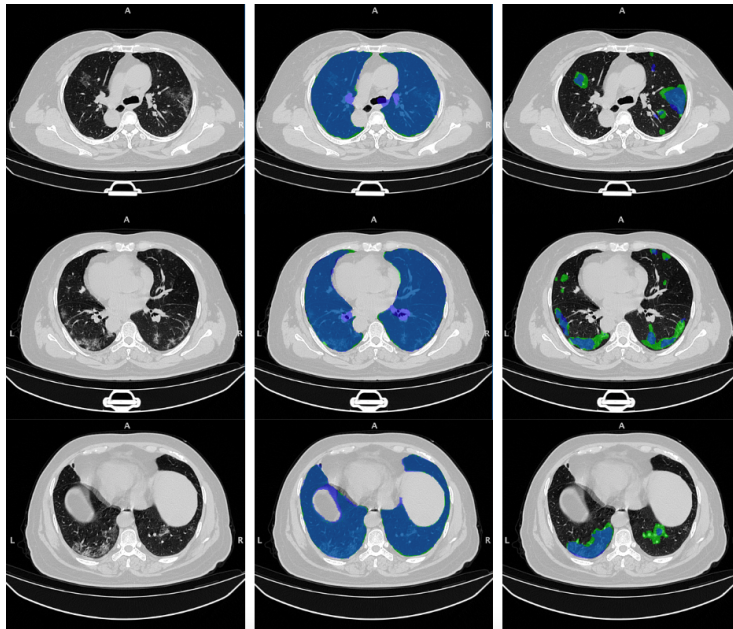


Figure 3: On the rows: three axial slices of the first CT scan on the COVID-19-CT-Seg test dataset (*coronacases001.nii*) are shown. On the columns: original images (left); overlays between the predicted (blue) and the reference (green) lung (center) and COVID-19 lesion (right) masks. The predicted masks were obtained by the *LungQuant* system integrating U-net₂^{90%}.

3.3.2 Percentage of affected lung volume and CT-SS estimation

The lung and lesion masks provided by the *LungQuant* system can be further processed to derive the physical volumes of each mask and the ratios between the lesion and lung volumes. We show in Fig. 4 the relationship between the percentage of lung involvement as predicted by the *LungQuant* system vs. the corresponding values computed on the reference masks, for both the *LungQuant* systems where the U-net₂^{60%} and the U-net₂^{90%} were alternatively integrated. As test samples, we considered the COVID-19-CT-Seg fully independent test dataset, and we complemented it with the partially annotated sample of 50 CT scans of the MosMed collection belonging to the CT-1 class, for which the lesion

masks were provided. It has to be noticed that this MosMed subsample is not fully independent of the training process since part of this data was used to train the U-net₂ networks. As shown in Fig. 4, the dataset available for this test consists of CT scans with a low percentage P of affected lung, which in most cases is below then 10%. Despite the limited range of P values to carry out this test, an agreement between the *LungQuant* system output and the reference values is observed for both the systems where either U-net₂^{60%} or U-net₂^{90%} were integrated. In terms of the mean absolute error (MAE) among the estimated and the reference P values, we obtained as an average on the 60 test cases: MAE=2.4% (4.6% on COVID-19-CT-Seg and 1.9% on MosMed) for the *LungQuant* system with U-net₂^{60%} and MAE=2.1% (4.2% on COVID-19-CT-Seg and 1.7% on MosMed) for the system with U-net₂^{90%}.

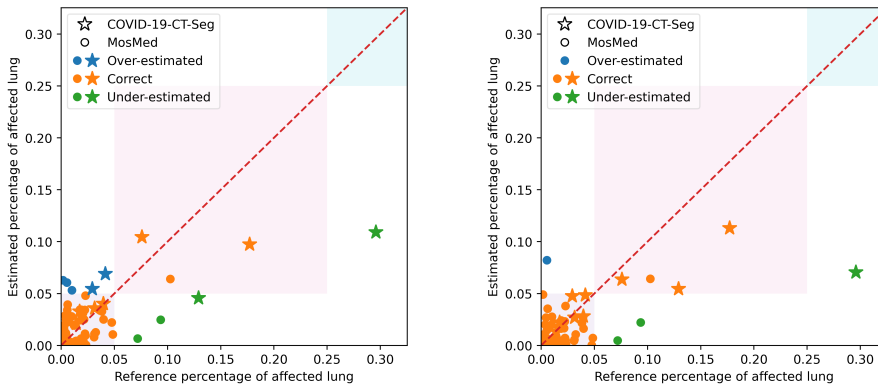


Figure 4: Estimated percentages P of affected lung volume versus the ground truth percentages, as obtained by the *LungQuant* system integrating U-net₂^{60%} (left) and U-net₂^{90%} (right). The colored areas in the plot backgrounds guide the eye to recognize the CT-SS values assigned to each value of P (from left to right: CT-SS=1, CT-SS=2, CT-SS=3).

The percentage of lung volume affected by COVID-19 lesions can also be directly converted into the CT-SS values. The accuracy in assigning the correct CT-SS class is reported in Table 6, together with the number of misclassified cases, for the 10 cases of the COVID-19-CT-Seg dataset and for the subset of 50 annotated MosMed CT scans. As reported in the table, an accuracy of 85% is achieved for the *LungQuant* system with U-net₂^{60%} and of 93% for the *LungQuant* system with U-net₂^{90%}. In all cases, the system misclassifies the examples (9 cases out of 60 for U-net₂^{60%} and 4 out of 60 cases for the U-net₂^{90%}) of 1 class at most.

Table 6: Classification performances of the whole system in predicting CT Severity Score on MosMed and Coronacases datasets.

U-net	Dataset	Accuracy	1-class misclassification	2-class misclassification
U-net ₂ ^{60%}	MosMed	45/50	5/50	0
	Coronacases	6/10	4/10	0
U-net ₂ ^{90%}	MosMed	47/50	3/50	0
	Coronacases	9/10	1/10	0

4 Discussion and Conclusion

We developed a fully automated quantification pipeline, the *LungQuant* system, for the identification and segmentation lungs and pulmonary lesions related to COVID-19 pneumonia in CT scans. The system returns the COVID-19 related lesions, the lung mask and the ratio between their volumes, which is converted into a CT Severity Score.

The performance obtained against a voxel-wise segmentation ground truth was evaluated in terms of the Dice index, which provides a measure of the overlap between the predicted and the reference masks. The *LungQuant* system achieved a Dice index of 0.95 ± 0.01 in the lung segmentation task and of 0.66 ± 0.13 in segmenting the COVID-19 related lesions on the fully annotated publicly available benchmark COVID-19-CT-Seg dataset of 10 CT scans.

Regarding the correct assignment of the CT-SS, the *LungQuant* system showed an accuracy of 93%, considering the subjects for which the ground truth information was either directly available or derivable within our analysis pipeline. The *LungQuant* system misclassified only the 7% of cases of one CT-SS class at most. Despite this result is encouraging, it was obtained on a rather small dataset, constituted by COVID-19-CT-Seg and MosMed CT scans, which involves most subjects with low disease severity, thus, a broader validation of larger data sample with more heterogeneous composition in terms of disease severity is required.

Nonetheless, the *LungQuant* image analysis system we developed can be a useful support tool to assist clinicians in their workflows during the COVID-19 pandemic.

acknowledgements

This work has been carried out within the Artificial Intelligence in Medicine (AIM) project funded by INFN (CSN5, 2019-2021), <https://www.pi.infn.it/aim>. We are grateful to the staff of the Data Center of the INFN Division of Pisa. We thank the CINECA Italian computing center for making available part of the computing resources used in this paper; in particular, Dr. Tommaso Boccali (INFN, Pisa) as PI of PRACE Project Access #2018194658 and a 2021 ISCRA-C grant. Moreover, we thank the EOS cluster of Department of Mathematics

”F. Casorati” (Pavia) for computing resources.

Conflict of interest

The authors declare that they have no conflict of interest.

Ethical approval and informed consent

All procedures performed in studies involving human participants were in accordance with the ethical standards of the institutional and/or national research committee and with the 1964 Helsinki Declaration and its later amendments or comparable ethical standards.

Informed consent was obtained from all individual participants included in the study.

References

- [1] An, P., Xu, S., Harmon, S.A., Turkbey, E.B., Sanford, T.H., Amalou, A., Kassin, M., Varble, N., Blain, M., Anderson, V., Patella, F., Carrafiello, G., Turkbey, B.T., Wood, B.J.: CT Images in COVID-19 (2020). DOI <https://doi.org/10.7937/tcia.2020.gqry-nc81>
- [2] Carotti, M., Salaffi, F., Sarzi-Puttini, P., Agostini, A., Borgheresi, A., Minorati, D., Galli, M., Marotto, D., Giovagnoni, A.: Chest CT features of coronavirus disease 2019 (COVID-19) pneumonia: key points for radiologists. *Radiologia Medica* **125**(7), 636–646 (2020). DOI [10.1007/s11547-020-01237-4](https://doi.org/10.1007/s11547-020-01237-4). URL <https://doi.org/10.1007/s11547-020-01237-4>
- [3] Chollet, F.: Keras. <https://keras.io> (2015)
- [4] Fang, X., Kruger, U., Homayounieh, F., Chao, H., Zhang, J., Digmurthy, S.R., Arru, C.D., Kalra, M.K., Yan, P.: Association of AI quantified COVID-19 chest CT and patient outcome. *International Journal of Computer Assisted Radiology and Surgery* (2021). DOI [10.1007/s11548-020-02299-5](https://doi.org/10.1007/s11548-020-02299-5). URL <http://www.ncbi.nlm.nih.gov/pubmed/33484428>
- [5] Fang, X., Yan, P.: Multi-organ segmentation over partially labeled datasets with multi-scale feature abstraction. *IEEE Transactions on Medical Imaging* **39**(11), 3619–3629 (2020). DOI [10.1109/TMI.2020.3001036](https://doi.org/10.1109/TMI.2020.3001036)
- [6] Hofmanninger, J., Prayer, F., Pan, J., Röhrich, S., Prosch, H., Langs, G.: Automatic lung segmentation in routine imaging is a data diversity problem, not a methodology problem. *arXiv* **2** (2020)

- [7] Kiser, K.J., Ahmed, S., Stieb, S., Mohamed, A.S., Elhalawani, H., Park, P.Y., Doyle, N.S., Wang, B.J., Barman, A., Li, Z., Zheng, W.J., Fuller, C.D., Giancardo, L.: PleThora: Pleural effusion and thoracic cavity segmentations in diseased lungs for benchmarking chest CT processing pipelines. *Medical Physics* **47**(11), 5941–5952 (2020). DOI 10.1002/mp.14424
- [8] Lessmann, N., Sánchez, C.I., Beenen, L., Boulogne, L.H., Brink, M., Calli, E., Charbonnier, J.P., Dofferhoff, T., van Everdingen, W.M., Gerke, P.K., Geurts, B., Gietema, H.A., Groeneveld, M., van Harten, L., Hendrix, N., Hendrix, W., Huisman, H.J., Išgum, I., Jacobs, C., Kluge, R., Kok, M., Krdzalic, J., Lassen-Schmidt, B., van Leeuwen, K., Meakin, J., Overkamp, M., van Rees Vellinga, T., van Rikxoort, E.M., Samperna, R., Schaefer-Prokop, C., Schalekamp, S., Scholten, E.T., Sital, C., Stöger, J.L., Teuwen, J., Venkadesh, K.V., de Vente, C., Vermaat, M., Xie, W., de Wilde, B., Prokop, M., van Ginneken, B.: Automated assessment of COVID-19 reporting and data system and chest CT severity scores in patients suspected of having COVID-19 using artificial intelligence. *Radiology* **298**(1), E18–E28 (2021). DOI 10.1148/RADIOL.2020202439
- [9] Ma, J., Wang, Y., An, X., Ge, C., Yu, Z., Chen, J., Zhu, Q., Dong, G., He, J., He, Z., Cao, T., Zhu, Y., Nie, Z., Yang, X.: Toward data-efficient learning: A benchmark for COVID-19 CT lung and infection segmentation. *Medical Physics* (2020). DOI 10.1002/mp.14676
- [10] Morozov, S.P., Andreychenko, A.E., Pavlov, N.A., Vladzimirskyy, A.V., Ledikhova, N.V., Gomboleviskiy, V.A., Blokhin, I.A., Gelezhe, P.B., Gonchar, A.V., Chernina, V.: MosMedData: Chest CT Scans with COVID-19 Related Findings Dataset. medRxiv p. 2020.05.20.20100362 (2020). DOI 10.1101/2020.05.20.20100362. URL <http://medrxiv.org/content/early/2020/05/22/2020.05.20.20100362.abstract>
- [11] Ronneberger, O., Fischer, P., Brox, T.: U-net: Convolutional networks for biomedical image segmentation. *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)* **9351**, 234–241 (2015). DOI 10.1007/978-3-319-24574-4_28
- [12] Xie, W., Jacobs, C., Charbonnier, J.P., van Ginneken, B.: Relational Modeling for Robust and Efficient Pulmonary Lobe Segmentation in CT Scans. *IEEE Transactions on Medical Imaging* pp. 1–1 (2020). DOI 10.1109/tmi.2020.2995108
- [13] Yang, J., Sharp, G., Veeraraghavan, H., van Elmpft, W., Dekker, A., Lustberg, T., Gooding, M.: Data from Lung CT Segmentation Challenge. The Cancer Imaging Archive. (2017). DOI <http://doi.org/10.7937/K9/TCIA.2017.3r3fvz08>

SUPPLEMENTARY MATERIALS

Quantification of pulmonary involvement in COVID-19 pneumonia by means of a cascade of two U-nets: training and assessment on multiple datasets using different annotation criteria

Francesca Lizzi ^{1,2} et al.

¹ Scuola Normale Superiore, Pisa;

² National Institute of Nuclear Physics (INFN), Pisa division, Pisa, IT

May 7, 2021

1 Additional descriptions of Materials and Methods

1.1 Characteristics of the public datasets used in the study

1.1.1 The Plethora dataset

The PleThora dataset [4] is a chest CT scan collection with thoracic volume and pleural effusion segmentations, delineated on 402 CT studies of the Non-Small Cell Lung Cancer (NSCLC) radiomics dataset, available through the The Cancer Imaging Archive (TCIA) repository [3]. This dataset has been made publicly available to facilitate improvement of the automatic segmentation of lung cavities, which is typically the initial step in the development of automated or semi-automated algorithms for chest CT analysis. In fact, automatic lung identification struggles to perform consistently in subjects with lung diseases. The PleThora lung annotations have been produced with a U-net based algorithm trained on chest CT of subjects without cancer, manually corrected by a medical student and revised by a radiation oncologist or a radiologist.

1.1.2 The 2017 Lung CT Segmentation Challenge dataset

The Lung CT Segmentation Challenge (LCTSC) dataset consists of CT scans of 60 patients, acquired from 3 different institutions and made publicly available in the context of the 2017 Lung CT Segmentation Challenge [9]. Since the aim of

this challenge was to foster the development of auto-segmentation methods for organs at risk in radiotherapy, the lung annotations followed the RTOG 1106 contouring atlas.

1.1.3 The 2020 COVID-19 Lung CT Lesion Segmentation Challenge dataset

The 2020 COVID-19 Lung CT Lesion Segmentation Challenge dataset (COVID-19 Challenge) is a public dataset consisting of unenhanced chest CT scans of 199 patients with positive RT-PCR for SARS-CoV-2 [1]. Each CT is accompanied with the ground truth annotations for COVID-19 lesions. Data has been provided in NIFTI format by The Multi-national NIH Consortium for CT AI in COVID-19 via the NCI TCIA public website [3]. Annotations have been made using a COVID-19 lesion segmentation model provided by NVIDIA, which takes a full CT chest volume and produces pixel-wise segmentations of COVID-19 lesions. These segmentations have been adjusted manually by a board of certified radiologists, in order to give 3D consistency to lesion masks. The dataset annotation was made possible through the joint work of Children’s National Hospital, NVIDIA and National Institutes of Health for the COVID-19-20 Lung CT Lesion Segmentation Grand Challenge.

The dataset and the annotations have been made available in the context of a MICCAI-endorsed international challenge (<https://covid-segmentation.grand-challenge.org/>) which had the aim to compare AI-based approaches to automated segmentation of COVID-19 lung lesions.

1.1.4 The MosMed dataset

MosMed [6] is a COVID-19 chest CT dataset collected by the Research and Practical Clinical Center for Diagnostics and Telemedicine Technologies of the Moscow Health Care Department. It includes CT studies taken from 1110 patients. Each study is represented by one series of images reconstructed into soft tissue mediastinal window. MosMed provides 5 labeled categories, based on the percentage of lung parenchyma affected by COVID-19 lesions. The 5 categories of lung involvement and their correspondence to the CT-SS scale are described in Table 1. The first category (CT-0) contains cases with normal lung tissue and no CT-signs of viral pneumonia, whereas the other categories contain GGO (CT-1 and CT-2) and both GGO and regions of consolidation in the higher classes (CT-3 and CT-4).

A small subset of class CT-1 cases (50 patients) had been annotated by expert radiologists with the support of MedSeg software (2020 Artificial Intelligence AS). The annotations consist of binary masks in which white voxels represent both ground-glass opacifications and consolidations. Both CT scans and annotations were provided in NIFTI format. During the DICOM-to-NIFTI conversion process, only one slice out of ten was preserved and, as a result, MosMed CT scans have a reduced total number of slices with respect to the other datasets.

Table 1: MosMed severity categories defined on the basis of the percentage P of lung volume affected by COVID-19 lesions. The correspondence to the CT-SS scale is reported.

MosMed CT category	N. of cases	Percentage P of involved lung parenchyma	Corresponding CT-SS
0	254	$P = 0$	0
1	684	$0 < P \leq 25$	1, 2
2	125	$25 < P \leq 50$	3
3	45	$50 < P \leq 75$	4
4	2	$75 < P \leq 100$	5

1.1.5 The COVID-19-CT-Seg dataset

The COVID-19-CT-Seg dataset is a collection of CT scans taken from the Coronacases Initiative and Radiopaedia [5]. It contains 20 CT scans tested positive for COVID-19 infection. This public dataset contains both lung and infection annotations. The ground truth has been made in three steps: first, junior radiologists (1-5 years of experience) delineated lungs and infections annotations, then two radiologists (5-10 years of experience) refined the labels and finally the annotations have been verified and optimized by a senior radiologist (more than 10 years of experience in chest radiology). The annotations have been produced with the ITK-SNAP software. Ten CT images of this dataset were provided in 8-bit depth, therefore, we decided to not use them.

1.2 Additional training details and evaluation strategy for the U-nets

1.2.1 Evaluation metrics

The segmentation performances for both U-nets have been evaluated with the Dice coefficient, computed between the true mask (M_{true}) and the predicted mask ($M_{predict}$), as follows;

$$\text{Dice}_{metric} = \frac{2 \cdot |M_{true} \cap M_{predict}|}{|M_{true}| + |M_{pred}|} \quad (1)$$

The loss function used to train the U-net for lung segmentation is the DICE loss, defined as follows

$$\text{Dice}_{loss} = 1 - \frac{2 \cdot |M_{true} \cap M_{pred}|}{|M_{true}| + |M_{pred}|} \quad (2)$$

and computed only on the foreground (white voxels). We used this strategy in order to avoid giving excessive weight to the background (black voxels), since the number of black and white voxels is quite unbalanced in favor of the former.

For U-net₂, we used a loss function (L) consisting of the sum of the Dice loss and a weighted cross-entropy (CE), defined as follows:

$$L = Dice_{loss} + CE_{weighted} \quad (3)$$

$$CE_{weighted} = w(x) \sum_{x \in \Omega} \log(M_{true}(x) \cdot M_{pred}(x)) \quad (4)$$

where $w(x)$ is the weight map which takes into account the frequency of white voxels, x is the current sample and Ω is the training set.

Since the background class is larger than the foreground class on the order 10^3 , we computed the weight map $w(x)$ for each ground-truth segmentation to increase the relevance of the underrepresented class, following the approach described in [7]. The weight map was defined as $w(x) = w_0/f_j$ where f_j is the average number of voxels of the j^{th} class over the entire training data set ($j = 0, 1$) and w_0 is the the average between the frequencies f_j .

1.2.2 Data augmentation

Data augmentation is a strategy to increase the size of the training set by synthetically generating additional training images through geometric transformations. This technique is particularly important to improve the generalization capability of the model, especially in the case of a limited number of training samples. In our work, we applied data augmentation during the data pre-processing phase (after defining the bounding boxes enclosing the segmented lungs) in order to generate a fixed number of augmented images for each original data. We chose an augmentation factor equal to 2 which means that the number of artificially generated images is twice the number of the original training set. For each image in the training set, two of the following geometric transformations were randomly chosen:

- **Zooming.** The CT image and the ground truth masks were zoomed in the axial plane, using a third-order spline interpolation and the k-nearest neighbor method, respectively. The zooming factor was randomly chosen among the following values: 1.05, 1.1, 1.15, 1.2.
- **Rotation.** The CT image and the ground truth mask were rotated in the axial plane, using a third-order spline interpolation and the k-nearest neighbor method, respectively. The rotation angle was randomly sampled among the following values: -15° , -10° , -5° , 5° , 10° , 15° .
- **Gaussian noise.** An array of noise terms randomly drawn from a normal distribution was added to the original CT image. For each image, the mean of the Gaussian distribution was randomly sampled in the $[-400, 200]$ HU range and the standard deviation randomly chosen among 3 values: 25, 50, 75 HU.
- **Elastic deformation.** An elastic distortion was applied to the original 3D CT and mask arrays following the approach of Simard *et al.* [8]. This

transformation has two parameters: the elasticity coefficient which we fixed to 12 and the scaling factor, fixed to 1000.

- Motion blurring. Slice by slice, we convolved the CT image with a linear kernel (i.e. ones along the central row and zero elsewhere for a matrix of size $k \times k$) through the function `filter2D`, defined in the OpenCV Python library [2], keeping the output image size the same as the input image. The filter is applied with a kernel size of 4, 3, and 3, in the anterior-posterior, latero-lateral and cranio-caudal direction, respectively.

An example of the application of these augmentation techniques to one CT scan of the dataset is provided in Fig. 1.

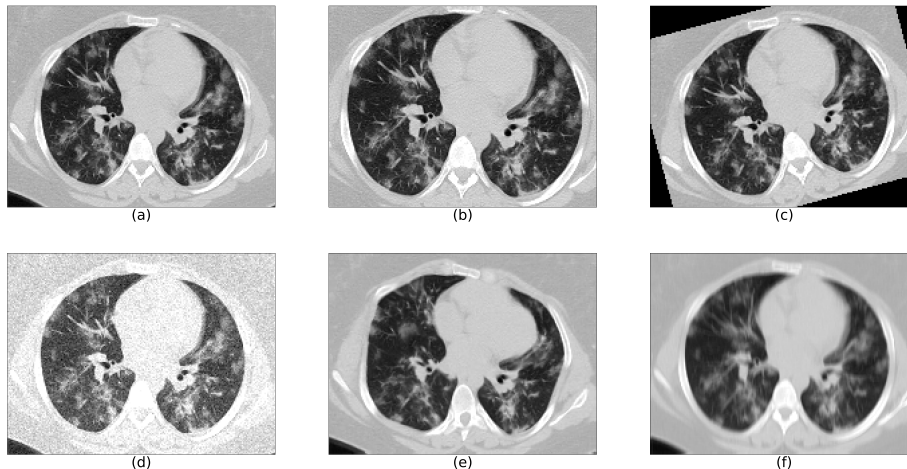


Figure 1: Data augmentation to increase the diversity of dataset: a) Image without data augmentation; b) Zooming; c) Rotation; d) Gaussian noise; e) Elastic deformations; f) Motion blurring.

1.3 Morphological refinement of U-net₁ lung segmentation

In order to remove false-positive regions (*i.e.* voxels misclassified as lung parts), at first, we identified the connected components in the lung masks generated by U-net₁, then, we excluded those components whose number of voxels was below an empirically-fixed threshold. This threshold was set to the 40% of the foreground mask, and it was reduced to 30% whether the resulting number of voxels was found to be lower than the 65% of the initial mask provided by U-net₁. Figure 2 shows some examples of how this procedure works on real CT scans.

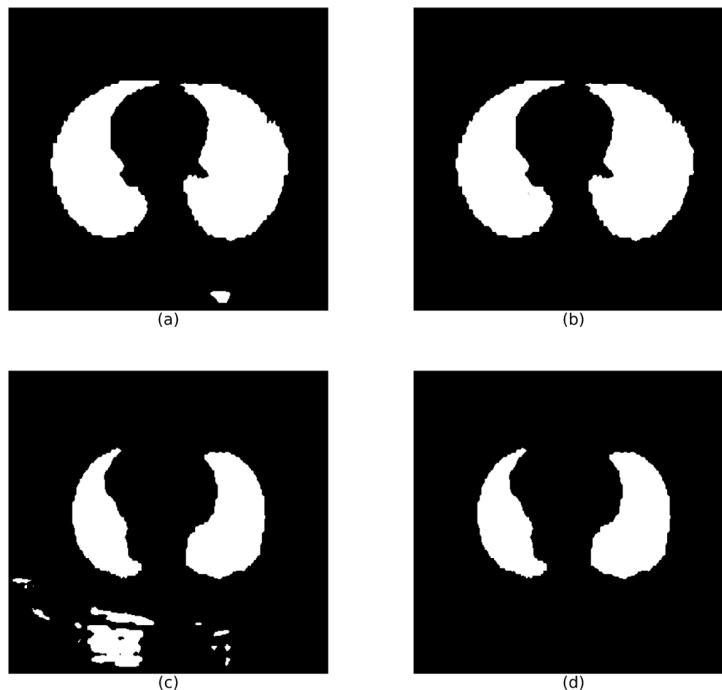


Figure 2: Morphological refinement of the $U\text{-net}_1$ output: a) and c) lung masks as generated by $U\text{-net}_1$; b) and d) refined masks after the connected component selection.

1.4 Generation of a set of reference lung segmentation for model training

As reported in Table 1 (main paper), the available datasets with lung mask annotations, which were necessary to train the U-net for lung segmentation, are mainly of subjects affected by lung cancer (Plethora and LCTSC datasets). To complement this sample with subjects without lesions, and, at the same time, to expose to U-net to the acquisition characteristics of the MosMed CT scans, we generated the lung mask annotations for a subset of subjects of the CT-0 MosMed category, i.e. that of subjects without COVID-19 lesions.

An in-house lung segmentation algorithm was developed for this purpose and implemented in *matlab* (The MathWorks, Inc.). It is based on the following steps: 1) CT windowing in the $[-1000,1000]$ HU range; 2) rough segmentation of the lungs on a central coronal slice (Otsu binary thresholding and removal of components connected with the image border) to define the minimum and maximum axial coordinates of the lung region; 3) 2D rough segmentation of the lungs on each axial slice (same procedure as the previous step) to generate a 3D seed mask for the following step; 4) segmentation of the lung parenchyma by an active contour model (*activecontour* matlab function); 5) filling holes (e.g.

vessels and airway walls) with 3D morphological operators (*imclose* matlab function).

This algorithm, which accurately segments the lung parenchyma in absence of lesions, has very limited performance on CT scans of subjects with COVID-19 lesions.

References

- [1] An, P., Xu, S., Harmon, S.A., Turkbey, E.B., Sanford, T.H., Amalou, A., Kassin, M., Varble, N., Blain, M., Anderson, V., Patella, F., Carrafiello, G., Turkbey, B.T., Wood, B.J.: CT Images in COVID-19 (2020). DOI <https://doi.org/10.7937/tcia.2020.gqry-nc81>
- [2] Bradski, G.: The OpenCV Library. Dr. Dobb's Journal of Software Tools (2000)
- [3] Clark, K., Vendt, B., Smith, K., Freymann, J., Kirby, J., Koppel, P., Moore, S., Phillips, S., Maffitt, D., Pringle, M., Tarbox, L., Prior, F.: The cancer imaging archive (TCIA): Maintaining and operating a public information repository. *Journal of Digital Imaging* **26**(6), 1045–1057 (2013). DOI 10.1007/s10278-013-9622-7
- [4] Kiser, K.J., Ahmed, S., Stieb, S., Mohamed, A.S., Elhalawani, H., Park, P.Y., Doyle, N.S., Wang, B.J., Barman, A., Li, Z., Zheng, W.J., Fuller, C.D., Giancardo, L.: PleThora: Pleural effusion and thoracic cavity segmentations in diseased lungs for benchmarking chest CT processing pipelines. *Medical Physics* **47**(11), 5941–5952 (2020). DOI 10.1002/mp.14424
- [5] Ma, J., Wang, Y., An, X., Ge, C., Yu, Z., Chen, J., Zhu, Q., Dong, G., He, J., He, Z., Cao, T., Zhu, Y., Nie, Z., Yang, X.: Toward data-efficient learning: A benchmark for COVID-19 CT lung and infection segmentation. *Medical Physics* (2020). DOI 10.1002/mp.14676
- [6] Morozov, S.P., Andreychenko, A.E., Pavlov, N.A., Vladzimirskyy, A.V., Ledikhova, N.V., Gombolevskiy, V.A., Blokhin, I.A., Gelezhe, P.B., Gonchar, A.V., Chernina, V.: MosMedData: Chest CT Scans with COVID-19 Related Findings Dataset. medRxiv p. 2020.05.20.20100362 (2020). DOI 10.1101/2020.05.20.20100362. URL <http://medrxiv.org/content/early/2020/05/22/2020.05.20.20100362.abstract>
- [7] Phan, T.H., Yamamoto, K.: Resolving class imbalance in object detection with weighted cross entropy losses. arXiv preprint arXiv:2006.01413 (2020)
- [8] Simard, P.Y., Steinkraus, D., Platt, J.C.: Best practices for convolutional neural networks applied to visual document analysis. In: *Icdar*, vol. 2003-January, pp. 958–963. IEEE Computer Society (2003). DOI 10.1109/ICDAR.2003.1227801

- [9] Yang, J., Sharp, G., Veeraraghavan, H., van Elmpt, W., Dekker, A., Lustberg, T., Gooding, M.: Data from Lung CT Segmentation Challenge. The Cancer Imaging Archive. (2017). DOI <http://doi.org/10.7937/K9/TCIA.2017.3r3fvz08>