# Detecting and mitigating issues in image-based COVID-19 diagnosis

João Marcos C. Silva [* 1]   Pedro Martelleto B. Rezende [* 1]   Moacir A. Ponti [1]

## Abstract

As urgency over the coronavirus disease 2019 (COVID-19) increased, many datasets with chest radiography (CXR) and chest computed tomography (CT) images emerged aiming at the detection and prognosis of COVID-19. Over the last two years, thousands of studies have been published, reporting promising results. However, a deeper analysis of the datasets and the methods employed reveals issues that may hamper conclusions and practical applicability. We investigate three major datasets commonly used in these studies, detect problems related to the existence of duplicates, address the specificity of classes within those datasets, and propose a way to perform external validation via cross-dataset evaluation. Our guidelines and findings contribute towards a trustworthy application of Machine Learning in the context of image-based diagnosis, as well as offer a more accurate assessment of models applied to the prognostication of diseases using image datasets and pave the way towards models that can be relied upon in the real world.

## 1. Introduction

Machine Learning (ML) has become a foremost approach for solving real-world problems from data. Supervised learning, in particular, has risen to be a popular method for building models that learn to classify inputs from a set of examples, amongst other tasks. When developing these models, the overarching goal is to learn a function that can generalize to examples outside of those provided during training. Indeed, fitting a model to a set of training samples is often trivial with Deep Neural Networks. However, generalizing to instances outside of training is a much more challenging problem (Mello & Ponti, 2018).

The usual method for empirically determining how well a model generalizes is dividing the original dataset into two disjoint groups: the training and testing data. The idea is to train on samples distinct from those used to evaluate the model, allowing researchers to assess whether the model is overfitting or is capable of generalizing from its original examples. However, depending on the task at hand, the manner by which this procedure is carried out can determine its practical applicability. Particularly, external validation (training on one dataset and testing in another) is recommended for sensitive applications since it may reveal weaknesses of models performing well in internal validation (i.e. those that were trained and tested on disjoint splits of data taken from the same dataset) (Cabitza et al., 2021).

When it comes to diagnostics or screening for diseases, assessing classification models in both internal and external validation settings is crucial. A great number of studies is dedicated to training deep networks with the use of images (Roberts et al., 2021) and audio (Coppock et al., 2021; Casanova et al., 2021) to detect COVID-19. In theory, Deep Convolutional Neural Networks (CNNs) could take advantage of an extensive array of example chest radiographs with positive/negative labels for COVID-19, and learn how to diagnose entirely new patients (Wang et al., 2020). Even though reverse transcription-polymerase chain reaction (RT-PCR) is the gold standard test for COVID-19, supervised learning can be applied to CXR and CT imaging to aid physicians in diagnosing COVID-19 when RT-PCR is not promptly available, as well as to look for agreement between different tests (image-based, RT-PCR, and antibodies) at each clinical stage (Hernández-Huerta et al., 2021) (Ai et al., 2020). Also, an emergent problem is the prediction of severity of the disease (Cohen et al., 2020), for which chest scans may offer valuable information to physicians.

Several works have reported impressively high accuracies when applying ML to image-based COVID-19 diagnosis, e.g., studies describing three popular datasets (Chowdhury et al., 2020) (Sait, 2021) (Wang et al., 2020). However, a review done in 2021 found that, so far, not a single paper on image-based COVID-19 diagnosis has had sufficient evidence to confirm the feasibility for clinical use (Roberts et al., 2021). This calls into question the practical use of artificial intelligence for such task (El Naqa et al., 2021).

---

[*]Equal contribution (The authors are ordered alphabetically) [1]Instituto de Ciências Matemáticas e de Computação, Universidade de São Paulo, Brazil. Correspondence to: Moacir A. Ponti <moacirponti@gmail.com>.

Although many different methods, models and settings were tested by several authors (Garg et al., 2020), a recurrent problem in previous work is their inability to conclude from testing metrics that the model is robust to enough settings to be applied in the real world. As we will be seeing in later sections, networks trained in COVID-19 CXR and CT datasets achieve high testing accuracy only in the dataset it trained on, and fail to generalize to similar testing samples from new datasets. Furthermore, despite some datasets' claims of having performed procedures to remove duplicates, we found that even a pixel-based similarity is able to find a significant amount of duplicates in all three datasets studied.

**This paper is an effort towards addressing important issues raised, for example, by Roberts et al. (2021); El Naqa et al. (2021). We strive to confirm some of those problems in practice using current and popular datasets, and report methods to mitigate such problems**. More specifically, we show contributions on:

- detecting duplicates in crowd-sourced datasets and evaluating the effects of removing them on classification metrics;

- performing external validation of models and evaluating the feasibility of transfer learning approaches, which would resemble real-world scenarios.

Therefore, we do not propose yet another method to attempt to diminish the problem. Instead, we show how to analyze datasets and handle models trained on them, offering a case study of how ML testing metrics can be misleading when important procedures are neglected.

We hope that our findings are a step towards better understanding of common pitfalls when developing critical deep learning models for healthcare applications. We also highlight the importance of attending to the underlying assumptions regarding the construction of the image datasets in the context of diagnosing COVID-19. Our code is publicly available at github.com/JoaoMarcosCSilva/issues-covid-image-diagnosis.

## 2. Related Work

Despite the enormous success of recent deep neural networks on a variety of problems, the field is plagued by serious issues on the reproducibility of the results and clarity of the methodology. Those issues may hamper the reliable deployment of many real-world deep learning models. For instance, (Wagstaff, 2012) criticizes the Machine Learning community as a whole due to its reliance on a few benchmark datasets and metrics that often do not reflect future performance in the real world. Insufficient statistical hypothesis tests, extensive hyperparameter tuning, and the lack of

satisfactory theoretical justifications often lead to the development of new techniques that have no significant impact beyond the settings they were tested in. For example, see (Schmidt et al., 2021) for neural network optimizers and (Narang et al., 2021) for modifications in the transformer architecture. When it comes to image-based COVID-19 diagnosis, (Roberts et al., 2021) exemplifies this problem by reviewing 62 studies published in 2020 (several of which reported very high testing accuracies) and concluding that none of the models are clinically usable due to risk of bias, overlapping train and test data, duplicated entries due to aggregated datasets and other severe methodological issues.

The over-reliance of the Machine Learning community on benchmark datasets may lead to wasted time, energy and effort in the search for small variations on existing techniques that result in small improvements on standard metrics, but do not actually advance the generalization capabilities of a neural network on unseen examples. Furthermore, sometimes even a slightly deviation from the distribution used in training may severely harm generalization. For example, (Recht et al., 2019) finds that CNNs trained on ImageNet (Deng et al., 2009), one of the most commonly used general-purpose image classification datasets, have a much smaller performance when evaluated on a new test set that was collected using the same procedure and sources as the original. Meanwhile, (Radford et al., 2021) found that their multi-modal architecture, which uses natural language as supervision instead of ImageNet, has an out-of-distribution performance much larger than that of models which were directly trained on the ImageNet benchmark.

However, these challenges are not exclusive to pure Machine Learning research. Scientists who work on applied ML are also often reliant on non-representative metrics and data, leading to subpar performance, which is exacerbated by the lack of contact with experts in the specific domain. For instance, (Kapoor & Narayanan, 2021) examines application papers in the social sciences and finds that several studies published in highly prestigious venues do not perform better than a simple logistical regression after correcting methodological errors such as data leakage or bad cross-validation splits.

In settings where robustness and generalization is critical, such as medical and healthcare studies, black box models (i.e. models for which we understand the inputs and outputs, but not the inner workings) that fail to generalize are particularly problematic and forbid practical utility. To mitigate this issue, several approaches have been proposed. Yang et al. (2021) provides a literature review that describes how explainability can aid in developing ML models that work in the real world, improving our understanding of the inner workings of CNN classifiers. Ter-Sarkisov (2020) takes a two step approach that first detects instances of ground

glass opacity and consolidation in CT scans, then predicts the patient's condition using only the ranked bounding box detections (thus making the model able to generalize even with a very small amount of training data). (Teixeira et al., 2021) uses a U-Net CNN architecture to perform semantic segmentation that isolates the lungs on CXR scans before trying to classify whether the patient has COVID-19, but still finds that the underlying data source may impose a strong bias which could severely hamper results.

Our work investigates current issues with some of the most popular datasets used for COVID-19 diagnosis from CXR images in an attempt to more closely understand the problems faced by researchers when dealing with those datasets, evaluate how current methods perform, and investigate how to mitigate some of those problems.

## 3. Datasets

We have used three openly available datasets (Table 1) considering the total availability of images and the number of paper citing them. Figure 1 shows samples from them.

*Table 1.* Comparison of datasets Curated (Sait, 2021), Radiography (Chowdhury et al., 2020) (Rahman et al., 2021) and COVIDx (Wang et al., 2020). Citations as of 27 May 2022.

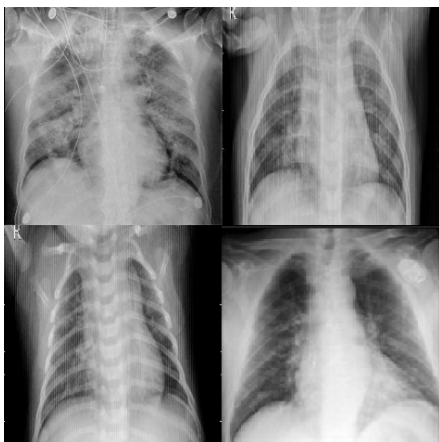|  | Curated | Radiog. | COVIDx |
|---|---|---|---|
| Scholar Citations | 59 | 713 | 1,700 |
| Number of classes | 4 | 4 | 3 |
| Fixed train/test split | No | No | Yes |
| Covid-19 Images | 1,281 | 3,616 | 16,690 |
| Total Images | 9208 | 21165 | 30530 |



*Figure 1.* Samples of chest X-Rays taken from the datasets used in this paper. Each image is labeled COVID-19 positive or not. Some datasets present extra classes, such as viral and bacterial pneumonia.

- ***Curated** Dataset for COVID-19 Posterior-Anterior Chest Radiography Images*: described in (Sait, 2021) and was created by a curated combination of 15 other publicly available sources. The authors claim to have removed duplicates and image imperfections using image similarities and the learned representations of an Inception V3 model. For simplicity, we will be referring to this dataset as *Curated Dataset*. The class frequencies of the provided data are: 3270 Normal Images; 1281 COVID-19 Images; 1656 Viral Pneumonia Images 3001 Bacterial Pneumonia Images.

- *COVID-19 **Radiography** Database*: described in (Chowdhury et al., 2020) and (Rahman et al., 2021), it was also made by combining various different sources. For simplicity, we will be referring to this dataset as *COVID-19 Radiography*. We used the second updated version of the dataset, which has the following class frequencies: 10,192 Normal Images; 3,616 COVID-19 Images; 1,345 Viral Pneumonia Images; 6,012 Lung Opacity Images.

- ***COVIDx** Dataset*: described in (Wang et al., 2020) is based on 5 different sources. The class distribution is as follows: 8,185 Normal Images, of which 100 are in the test set; 16,690 COVID-19 Images, of which 200 are in the test set; 5,655 Pneumonia Images, of which 100 are in the test set.

## 4. Methods

We propose the use of a pipeline of dataset analysis and pre-processing, as well as an external validation setup in order to investigate and mitigate issues with these datasets and point out possibilities for external validation, which was shown to be needed in the context of ML for health applications.

### 4.1. Duplicates removal

One of the main issues raised by (Roberts et al., 2021) is the existence of duplicates due to the crowd-sourcing nature of the datasets. Duplicates may cause the same instance to be at the same time in training and test sets, generating unrealistic evaluation metrics.

We employ two methods to find and remove duplicates. The first uses **pixel-by-pixel cosine similarities** between the images of the dataset and removes all images with similarity greater than some threshold. The second method uses the **similarities of image embeddings**, in particular the learned embeddings of a ResNet-18 classifier trained with the available training data. We chose the output of the second-to-last convolutional block as the embedding layer.

Let $Sim(\alpha(D_i), \alpha(D_j))$ denote the pixel-wise cosine similarity between the images $D_i$ and $D_j$ of the dataset, and

$Sim(\beta(D_i), \beta(D_j))$ the cosine similarity of the neural embeddings from the two images, where $\alpha$ is a function that outputs the image's representation as a vector of pixel values and $\beta$ outputs its representation after passing through a trained ResNet. We compute $MaxSim(D_i)$, the maximum similarity observed for some image $D_i$:

$$MaxSim(D_i) = \arg\max_{i \neq j} Sim(\Phi(D_i), \Phi(D_j)), \quad (1)$$

where $\Phi(.)$ could be either $\alpha(.)$ or $\beta(.)$. $MaxSim$ presents a way to measure candidate duplicates for all images of the dataset, according to two distinct representations. Since this similarity is 1 for identical images, it is possible to define a similarity threshold that determines which images are duplicates. Empirically, we found that the distribution of maximum similarities between images follows a distinct pattern with two peaks: one to the left and one close to the maximum similarity $(1, 00)$, as depicted in Figure 2.
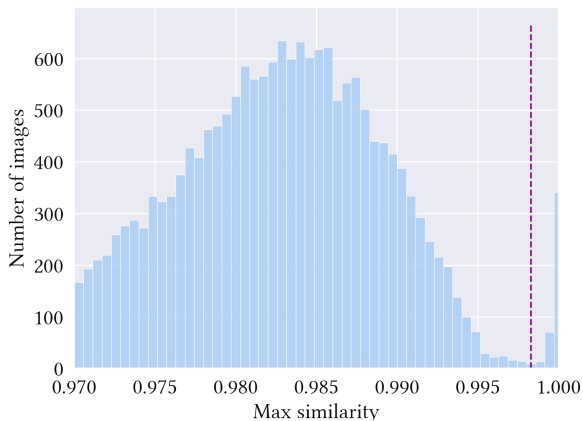


*Figure 2.* Distribution of maximum similarities between the images of the COVID-19 Radiography Database using pixel-wise similarity. The purple dashed line is the threshold chosen by our algorithm to select the images to be removed. Only pairs with maximum similarities close to 1 are shown in the graph.

To find this threshold, we employ the following algorithm: first, we divide the similarities into 50 bins between 0.95 and 1 and count the number of times the maximum similarities occur. Then, we iterate from right to left $(1, 00$ to minimum similarity) to find a region of negative slope that is close to the minimum similarity. Algorithmically, this is equivalent to finding a bin smaller than the previous bin and within some $\epsilon$ distance of the minimum. We employ this algorithm for both pixel and embedding-based similarities. Figure 2 shows the resulting threshold for an example setting.

### 4.2. External validation

A major challenge in the applicability of neural networks to the real world setting is the lack of generalization capabilities of trained models in settings even slightly different than those they were trained on Cabitza et al. (2021). For example, models trained for the diagnosis of a disease might not perform well if the target population has a different demographic composition than the one in the datasets used or if a different CRT machine is used (Oala et al., 2021). Even when trying to control for all of these factors, generalization is still a challenge.

To measure the extent that this issue affects the diagnosis of COVID-19 based on chest images, we perform experiments where we train a neural metwork for classification in one dataset, and evaluate its performance on the other two. Formally, let $S$ be a source dataset and $T$ be a target dataset. We first train a model $f$ using $S$ as input, and then perform inference on $T$. Later, we fine-tune $f$ with different percentages of data from $T$. While the first step allows us to understand the degree of generalization on external data, the second step investigates how much data from the target training data is needed to allow for a reasonable accuracy.

While it is true that the variability across the studied datasets does not reflect the true input distribution that would be found in the real world, a model that does not perform well across multiple datasets is unlikely to be of real applicability in a clinical setting (Oala et al., 2021).

Some challenges that prevent the direct use of external validation are: i) the different distributions of classes across each dataset, ii) the fact that not all sets have the same output categories, and iii) the existence of duplicates in the datasets since both could have the same instance (e.g. for multi/crowd-sourced datasets) although with different processing procedures. Thus, we adopt the duplicate detection technique from section 4.1. Also, since some classes are not present in all datasets, we simplify the classification to a binary problem, training models to determine exclusively whether a radiograph is from a patient with COVID-19 or if it is Normal, discarding other classes. The loss and accuracy are weighted based on the class distributions during both training and testing, to reduce the impact of class imbalance on the results.

Two experiments are then carried out: 1) a cross-dataset experiment using a source dataset for training and a different target dataset for testing, without any use of the target for adaptation, 2) a fine-tuning procedure, using the source datasets as pre-training and different proportions of the target dataset for tuning.

The second experiment is important since some models may still be useful despite not performing well directly on external validation (experiment 1). For instance, the

representations learned by the network could be used to obtain an initial model and improve results in a target dataset by allowing the model to adapt with data from this new dataset. Therefore, in the second experiment we measure the effects of pre-training on both accuracy and training cost.

### 4.3. Training details

We used a ResNet-18 V2 (He et al., 2016) trained from scratch with the AdamW (Loshchilov & Hutter, 2017) optimizer, starting with a learning rate of $10^{-1}$. All fine-tuning experiments were carried out with an initial learning rate of $10^{-2}$. Additionaly, we set the weight decay to $10^{-3}$, and employ a cosine decay schedule (Loshchilov & Hutter, 2016) with $NumEpochs \cdot \lfloor \frac{NumImages}{BatchSize} \rfloor$ steps. The training cross-entropy loss and the reported accuracy were weighted depending on the target class of a sample, so that all classes have equal importance on training and evaluation, regardless of different distributions across the datasets.

The number of epochs is set to 30 for all experiments, while the batch size is 128. As a pre-processing step, we scaled the images to a resolution of $256 \times 256$. We ran the experiments on a TPU v3 with JAX, Haiku, and Optax. We set the train/validation/test split proportions to 70%/10%/20% for all experiments, unless stated otherwise. All experiments were repeated 5 times using cross-validation.

## 5. Results and Discussion

### 5.1. Duplicates removal

Initially, we compare the pixel-wise and embedding approaches to determine whether they agree on which images are considered duplicates (Table 2). For this comparison, we explore the samples considered duplicates by both methods (intersection) and those found by only one of the two methods. By comparing the max similarity histograms, we find that the pattern presented in Figure 2 repeats for both pixel-wise and embedding similarity for all datasets, as shown for the Curated Dataset in Figure 3. We zoom in to show only the similarities larger than 0.97 since we are interested in finding pairs of duplicates and to emphasize the peak near similarity 1.00.

The Curated Dataset contains a high number of duplicates (28.5%), while the Radiography Database presents a smaller amount (1.2%), as does the COVIDx Dataset (0.9%). Surprisingly, both the pixel-wise and embedding-based methods were successful in finding the duplicates and agreed on almost all duplicates found. Due to this, we adopt the pixel-wise similarity going forward, as it is the simplest method of the two. However, if the datasets used contained rotated or cropped duplicates, we expect that the embeddings would have been much more reliable in detecting them. Tables

*Table 2.* Number of duplicate images detected in absolute number and percentage with respect to the whole dataset when using the pixel-wise and embedding approaches. "Intersection" is the number of duplicates discovered by both pixel-wise and the neural embeddings from the first cross-validation fold.

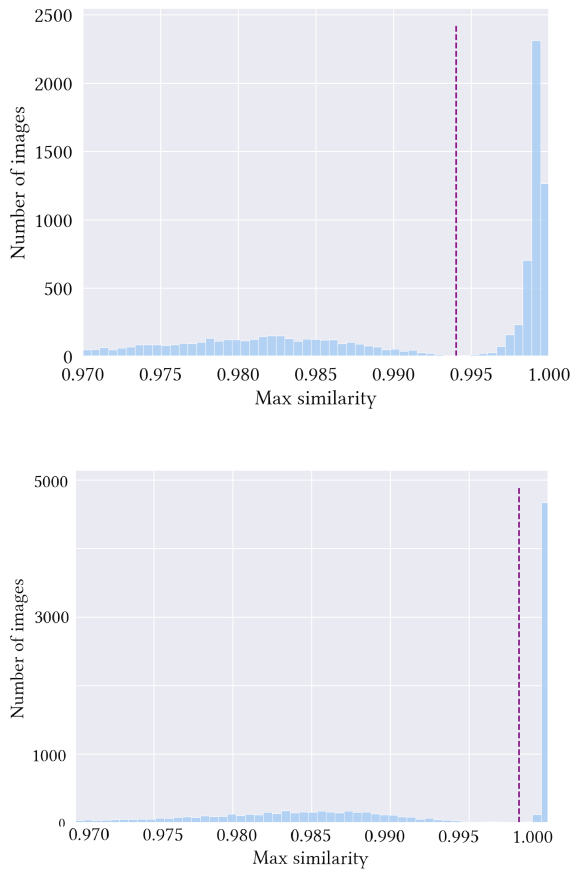|  | Curated | Radiog. | COVIDx |
|---|---|---|---|
| Pixel-wise | 2623 (28.5%) | 251 (1.2%) | 294 (1.0%) |
| Embeddings | 2634 (28.6%) | 244 (1.2%) | 276 (0.9%) |
| Intersection | 2621 (28.5%) | 244 (1.2%) | 273 (0.9%) |



*Figure 3.* Pixel-wise maximum similarity (top) and embedding-based maximum similarity (bottom) in the Curated Dataset. Despite a small difference in scale, both methods present similar patterns in their distributions. Only pairs with maximum similarities close to 1 are shown in the graphs.

3-5 show the class accuracies of a ResNet-18 V2 model trained on the dataset before and after removing the duplicates. As expected, we observe drops in accuracy on almost all instances after the duplicates were removed, an effect that is especially significant in the classes that had the most duplicates present. Tables 6-8 show a breakdown of the number of duplicates in each class, which we found to be concentrated in only one of the classes for each dataset.

*Table 3.* Testing accuracies before and after removing the duplicates from the Curated Dataset for COVID-19.

| Setting | Normal | Viral pneumonia | COVID-19 | Bacterial pneumonia |
|---|---|---|---|---|
| With duplicates | $0.960 \pm 0.010$ | $0.676 \pm 0.034$ | $0.962 \pm 0.008$ | $0.776 \pm 0.057$ |
| No duplicates | $0.928 \pm 0.013$ | $0.650 \pm 0.050$ | $0.974 \pm 0.015$ | $0.756 \pm 0.030$ |

*Table 4.* Testing accuracies before and after removing the duplicates from the COVID-19 Radiography Database.

| Setting | Normal | Viral pneumonia | COVID-19 | Lung opacity |
|---|---|---|---|---|
| With duplicates | $0.914 \pm 0.005$ | $0.956 \pm 0.018$ | $0.946 \pm 0.015$ | $0.866 \pm 0.011$ |
| No duplicates | $0.914 \pm 0.009$ | $0.952 \pm 0.011$ | $0.928 \pm 0.011$ | $0.854 \pm 0.011$ |

*Table 5.* Testing accuracies before and after removing the duplicates from the COVIDx dataset.

| Setting | Normal | Pneumonia | COVID-19 |
|---|---|---|---|
| With duplicates | $0.954 \pm 0.028$ | $0.920 \pm 0.031$ | $0.982 \pm 0.012$ |
| No duplicates | $0.872 \pm 0.027$ | $0.908 \pm 0.012$ | $0.974 \pm 0.017$ |

*Table 6.* Duplicates per class in the Curated Dataset

| Class | Number of Samples |
|---|---|
| Normal | 1721 |
| Viral Pneumonia | 298 |
| COVID-19 | 132 |
| Bacterial Pneumonia | 489 |

*Table 7.* Duplicates per class in the Radiography Database

| Class | Number of Samples |
|---|---|
| Normal | 1 |
| Viral Pneumonia | 7 |
| COVID-19 | 246 |
| Lung Opacity | 0 |

*Table 8.* Duplicates per class in the COVIDx Dataset

| Class | Number of Samples |
|---|---|
| Normal | 0 |
| Pneumonia | 4 |
| COVID-19 | 410 |

## 5.2. External validation

*Experiment 1:* After training for 30 epochs on each dataset, we observe that there is not much difficulty in achieving good performance in their corresponding test set (see the diagonals of Table 9). However, the accuracies are greatly decreased when a model is evaluated on a dataset different from the one it was trained on. We also note that the training loss diverged in one of the folds of the COVID-19 Radiography Database training set, so that fold was re-run with a different network initialization.

This drop in performance is especially significant on the Curated Dataset for COVID-19 Posterior-Anterior Chest Radiography Images, which has the least amount of training samples of all three. When evaluated on the COVIDx dataset, the accuracy is close to 46%, which is worse than a random guess. In addition to lower performance, the variance of the results also considerably increased in all external validation tests.

In every experiment, we observed a drop of at least 20% when evaluating on a target set, with the exception of evaluating a model trained on the COVID-19 Radiography Database on the Curated Dataset, in which the accuracy only dropped by around 5%. Overall, the COVID-19 Radiography Database was the dataset that best generalized, while the Curated Dataset was the worst, which can be at least partially explained by the vastly different amounts of training samples that each of them includes.

*Experiment 2:* Despite their considerable lack of generalization to different settings, the models might be useful for improving both the performance and data-efficiency of other neural networks by using their pre-trained weights as the initialization of a fine-tuning procedure. To verify whether this is the case, and to measure the influence of dataset size on this effect, we fine-tune models trained on both the Curated Dataset and COVIDx on different percentages of the COVID-19 Radiography Database. We compared the fine-tuning of models with a random initialization in order to investigate to what extent the source model contributes to the test accuracy (see Figure 4).
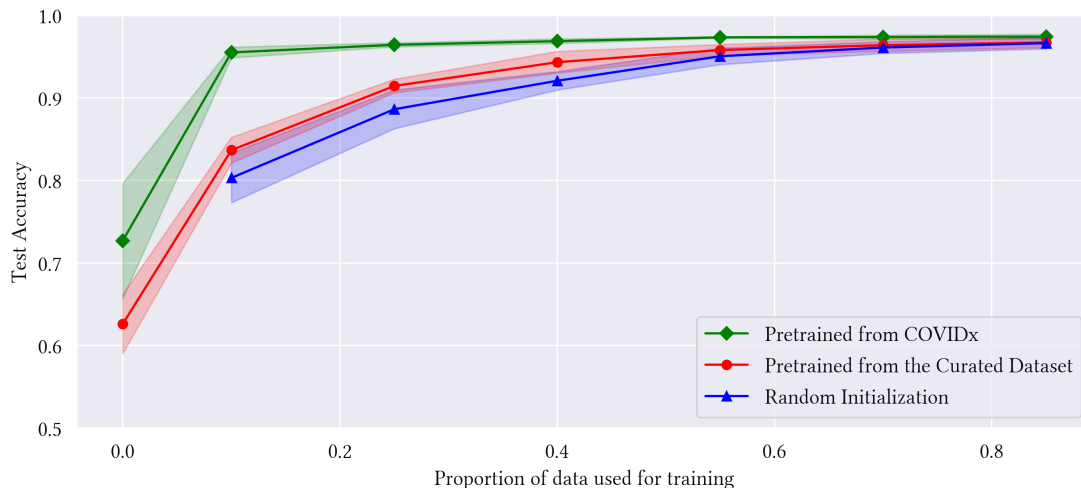
*Figure 4.* Model accuracies on the COVID-19 Radiography Database for different amounts of training samples available on the target dataset (i.e. proportional to the whole target training dataset) comparing random intialization with pre-training on each of the other two datasets. The shaded areas indicate a range of one standard deviation, estimated with 5-fold cross-validation.

In this experiment, fine-tuning on a small percentage of the data from the target distribution can greatly improve performance, especially when compared to the model without fine-tuning. Pre-training on COVIDx had a much greater effect compared to the Curated Dataset, which is expected from the fact that COVIDx has more than three times as many training samples. Both datasets, however, improved performance when compared to a randomly initialized network, in particular in the low-data regime.

*Table 9.* Class-balanced accuracy of models on each target dataset, depending on which dataset was used for training (the source dataset). The diagonals indicate each model's performance on its own test set.

| | Targets | | |
| Sources | Curated | Radiog. | COVIDx |
|---|---|---|---|
| Curated | 0.977±0.01 | 0.626±0.04 | 0.464±0.10 |
| Radiog. | 0.910±0.02 | 0.963±0.01 | 0.723±0.02 |
| COVIDx | 0.442±0.05 | 0.727±0.07 | 0.985±0.00 |

## 6. Conclusions

Our study reveals important issues when dealing with medical imaging for the automated diagnostic of COVID-19. First, mainly due to the fact that large datasets obtain images from multiple sources, a high rate of duplicates exist even in popular curated datasets. We show that the accuracy across different categories often reduces after the removal of duplicates. Therefore, a duplicate removal step is recommended to be adopted as default in future work.

Furthermore, we show the importance of assessing the models in an external validation setting. Although we easily obtained accuracies higher than 95% within a single dataset, when testing on external data we observed accuracies as low as 46% even when constraining the problem to a binary classification one. Designing methods that are less prone to memorization and actually capture the patterns associated with the disease are of foremost importance in this context.

Nevertheless, our study also shows promising results by fine-tuning the initial model with fractions of a dataset in the target distribution. Studies focusing on adapting models to target datasets could be an important aspect for future research, facilitating clinical viability in a variety of scenarios. Future work could also investigate models that determine the severity of the disease, building a more transparent framework to be used in practice by physicians.

## Acknowledgements

## References

Ai, T., Yang, Z., Hou, H., Zhan, C., Chen, C., Lv, W., Tao, Q., Sun, Z., and Xia, L. Correlation of chest CT and RT-PCR testing for coronavirus disease 2019 (COVID-19) in china: a report of 1014 cases. *Radiology*, 296(2): E32–E40, 2020.

Cabitza, F., Campagner, A., Soares, F., de Guadiana-Romualdo, L. G., Challa, F., Sulejmani, A., Seghezzi, M., and Carobene, A. The importance of being external. methodological insights for the external validation of machine learning models in medicine. *Computer Methods and Programs in Biomedicine*, 208:106288, 2021.

Casanova, E., Candido, A., Fernandes, R., Finger, M., Gris, L. R. S., Ponti, M. A., da Silva, D., et al. Transfer learning and data augmentation techniques to the COVID-19 identification tasks in ComParE 2021. In *22nd Annual Conference of the International Speech Communication Association, INTERSPEECH 2021*, pp. 4301–4305, 2021.

Chowdhury, M. E. H., Rahman, T., Khandakar, A., Mazhar, R., Kadir, M. A., Mahbub, Z. B., Islam, K. R., Khan, M. S., Iqbal, A., Emadi, N. A., Reaz, M. B. I., and Islam, M. T. Can ai help in screening viral and COVID-19 pneumonia? *IEEE Access*, 8:132665–132676, 2020. doi: 10.1109/ACCESS.2020.3010287.

Cohen, J. P., Dao, L., Roth, K., Morrison, P., Bengio, Y., Abbasi, A. F., Shen, B., Mahsa, H. K., Ghassemi, M., Li, H., et al. Predicting covid-19 pneumonia severity on chest x-ray with deep learning. *Cureus*, 12(7), 2020.

Coppock, H., Jones, L., Kiskin, I., and Schuller, B. Bias and privacy in ai's cough-based COVID-19 recognition–authors' reply. *The Lancet Digital Health*, 3(12):e761, 2021.

Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., and Fei-Fei, L. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009.

El Naqa, I. M., Li, H., Fuhrman, J. D., Hu, Q., Gorre, N., Chen, W., and Giger, M. L. Lessons learned in transitioning to ai in the medical imaging of COVID-19. *Journal of Medical Imaging*, 8(S1):010902, 2021.

Garg, T., Garg, M., Mahela, O. P., and Garg, A. R. Convolutional neural networks with transfer learning for recognition of COVID-19: A comparative study of different approaches. *AI*, 1(4):586–606, 2020.

He, K., Zhang, X., Ren, S., and Sun, J. Identity mappings in deep residual networks. In *European conference on computer vision*, pp. 630–645. Springer, 2016.

Hernández-Huerta, M. T., Mayoral, L. P.-C., Navarro, L. M. S., Mayoral-Andrade, G., Mayoral, E. P.-C., Zenteno, E., and Pérez-Campos, E. Should RT-PCR be considered a gold standard in the diagnosis of COVID-19? *Journal of Medical Virology*, 2021.

Kapoor, S. and Narayanan, A. (Ir)Reproducible Machine Learning: A Case Study. https://reproducible.cs.princeton.edu/, 2021. URL https://reproducible.cs.princeton.edu/.

Loshchilov, I. and Hutter, F. Sgdr: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983*, 2016.

Loshchilov, I. and Hutter, F. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.

Mello, R. F. and Ponti, M. A. *Machine learning: a practical approach on the statistical learning theory*. Springer, 2018.

Narang, S., Chung, H. W., Tay, Y., Fedus, W., Fevry, T., Matena, M., Malkan, K., Fiedel, N., Shazeer, N., Lan, Z., Zhou, Y., Li, W., Ding, N., Marcus, J., Roberts, A., and Raffel, C. Do transformer modifications transfer across implementations and applications?, 2021.

Oala, L., Murchison, A. G., Balachandran, P., Choudhary, S., Fehr, J., Leite, A. W., Goldschmidt, P. G., Johner, C., Schörverth, E. D., Nakasi, R., et al. Machine learning for health: Algorithm auditing & quality control. *Journal of medical systems*, 45(12):1–8, 2021.

Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al. Learning transferable visual models from natural language supervision. *arXiv preprint arXiv:2103.00020*, 2021.

Rahman, T., Khandakar, A., Qiblawey, Y., Tahir, A., Kiranyaz, S., Abul Kashem, S. B., Islam, M. T., Al Maadeed, S., Zughaier, S. M., Khan, M. S., and Chowdhury, M. E. Exploring the effect of image enhancement techniques on COVID-19 detection using chest x-ray images. *Computers in Biology and Medicine*, 132:104319, 2021. ISSN 0010-4825. doi: https://doi.org/10.1016/j.compbiomed.2021.104319.

Recht, B., Roelofs, R., Schmidt, L., and Shankar, V. Do imagenet classifiers generalize to imagenet? In *International Conference on Machine Learning*, pp. 5389–5400. PMLR, 2019.

Roberts, M., Driggs, D., Thorpe, M., Gilbey, J., Yeung, M., Ursprung, S., Aviles-Rivero, A. I., Etmann, C., McCague, C., Beer, L., et al. Common pitfalls and recommendations for using machine learning to detect and prognosticate for COVID-19 using chest radiographs and CT scans. *Nature Machine Intelligence*, 3(3):199–217, 2021.

Sait, U. Curated dataset for COVID-19 posterior-anterior chest radiography images (x-rays)., 2021. URL https://data.mendeley.com/datasets/9xkhgts2s6/3.

Schmidt, R. M., Schneider, F., and Hennig, P. Descending through a crowded valley - benchmarking deep learning optimizers. In Meila, M. and Zhang, T. (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 9367–9376. PMLR, 18–24 Jul 2021.

Teixeira, L. O., Pereira, R. M., Bertolini, D., Oliveira, L. S., Nanni, L., Cavalcanti, G. D., and Costa, Y. M. Impact of lung segmentation on the diagnosis and explanation of COVID-19 in chest x-ray images. *Sensors*, 21(21):7116, 2021.

Ter-Sarkisov, A. Covid-ct-mask-net: Prediction of covid-19 from ct scans using regional features. *medRxiv*, 2020.

Wagstaff, K. Machine learning that matters, 2012.

Wang, L., Lin, Z. Q., and Wong, A. COVID-net: a tailored deep convolutional neural network design for detection of COVID-19 cases from chest x-ray images. *Scientific Reports*, 10(1):19549, Nov 2020. ISSN 2045-2322. doi: 10.1038/s41598-020-76550-z. URL https://doi.org/10.1038/s41598-020-76550-z.

Yang, G., Ye, Q., and Xia, J. Unbox the black-box for the medical explainable ai via multi-modal and multi-centre data fusion: A mini-review, two showcases and beyond. *arXiv preprint arXiv:2102.01998*, 2021.