# Deep Learning-based Patient Re-identification Is able to Exploit the Biometric Nature of Medical Chest X-ray Data

**Kai Packhäuser[1,*], Sebastian Gündel[1], Nicolas Münster[1], Christopher Syben[1], Vincent Christlein[1], and Andreas Maier[1]**

[1]Pattern Recognition Lab, Department of Computer Science, Friedrich-Alexander-Universität Erlangen-Nürnberg, 91058 Erlangen, Germany
[*]kai.packhaeuser@fau.de

## ABSTRACT

**With the rise and ever-increasing potential of deep learning techniques in recent years, publicly available medical datasets became a key factor to enable reproducible development of diagnostic algorithms in the medical domain. Medical data contains sensitive patient-related information and is therefore usually anonymized by removing patient identifiers, e. g., patient names before publication. To the best of our knowledge, we are the first to show that a well-trained deep learning system is able to recover the patient identity from chest X-ray data. We demonstrate this using the publicly available large-scale ChestX-ray14 dataset, a collection of 112,120 frontal-view chest X-ray images from 30,805 unique patients. Our verification system is able to identify whether two frontal chest X-ray images are from the same person with an AUC of 0.9940 and a classification accuracy of 95.55 %. We further highlight that the proposed system is able to reveal the same person even ten and more years after the initial scan. When pursuing a retrieval approach, we observe an mAP@R of 0.9748 and a precision@1 of 0.9963. Furthermore, we achieve an AUC of up to 0.9870 and a precision@1 of up to 0.9444 when evaluating our trained networks on external datasets such as CheXpert and the COVID-19 Image Data Collection. Based on this high identification rate, a potential attacker may leak patient-related information and additionally cross-reference images to obtain more information. Thus, there is a great risk of sensitive content falling into unauthorized hands or being disseminated against the will of the concerned patients. Especially during the COVID-19 pandemic, numerous chest X-ray datasets have been published to advance research. Therefore, such data may be vulnerable to potential attacks by deep learning-based re-identification algorithms.**

Chest radiography (X-ray) is a modality that is routinely used for diagnostic procedures around the world[1]. It became the most common medical imaging examination for pulmonary diseases and allows a clear investigation of the thorax[2]. Chest X-ray imaging is therefore well-suited for diagnosing several pathologies including pulmonary nodules, masses, pleural effusions, pneumonia, COPD, and cardiac abnormalities[3]. It is also used for COVID-19[4] screening, as abnormalities typical of those infected with the coronavirus can be detected in radiographs[5]. While chest radiography plays a crucial role in clinical care, discovering certain diseases and abnormalities in chest radiographs can be a challenging task for radiologists, which potentially results in undesirable misdiagnoses[6]. Therefore, computer-aided detection (CAD) systems based on deep learning (DL)[7] techniques have been developed in recent years to facilitate radiology workflows. These systems, characterized by their enormous benefits, can be utilized for a wide range of applications, e. g., for the automatic recognition of abnormalities in chest radiographs[3,8] and the detection of tumors in mammography[9]. Some techniques even show the potential to exceed human performance[10]. However, the CAD systems are only treated as an additional source to support the radiologists and to increase certainty in their reading decisions.

On the one hand, the large variety of medical applications allows DL to grow and tackle real-life problems that were previously not solvable or improving solutions offered by traditional machine learning methods[7]. On the other hand, DL is a data-driven approach and well-known for its need for big data to train the neural networks[11,12]. For these reasons, a vast amount of medical datasets have been published in recent years that enable researchers to develop diagnostic algorithms in the medical field in a reproducible way[13]. These include several large-scale chest radiography datasets, e. g., the *CheXpert*[14], the *PLCO*[15] and the *ChestX-ray14*[16] datasets. But especially during the COVID-19 pandemic[17,18], the number of publicly available chest radiography datasets increased rapidly. A few selected examples are the *COVID-19 Image Data Collection*[19], the *Figure 1 COVID-19 Chest X-ray Dataset Initiative*[20], the *ActualMed COVID-19 Chest X-ray Dataset Initiative*[21], and the *COVID-19 Radiography Database*[22].
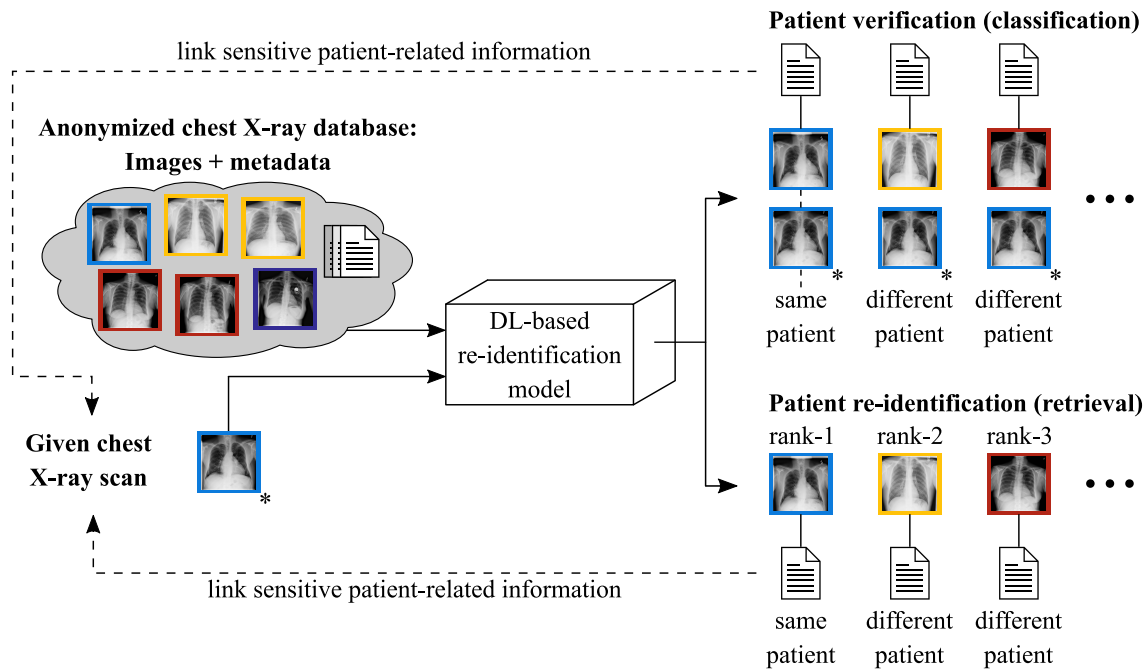
**Figure 1.** General problem scenario: Comparing a given chest radiograph to publicly available dataset images by means of DL techniques would either result in discrete labels indicating whether or not the dataset images belong to the same patient as the given radiograph (verification scenario) or yield a ranked list of the most similar radiographs related to the given scan (retrieval scenario). Images belonging to the same patient are highlighted with the same color. The given radiograph is marked with an asterisk. The shown cases would enable a potential attacker to link sensitive patient-related information contained in the dataset to the image of interest.

Chest radiography datasets typically consist of two parts: First, the image data itself, which provides clinical information about the anatomical structure of the thorax. Second, the associated metadata, which contains sensitive patient-related information that is either stored in a separate file or embedded directly in the images[23]. Proper data anonymization constitutes an important step when preparing medical data for public usage to ensure that a patient's identity cannot be revealed in publicly available datasets[23]. In practice, any personally identifiable information is attempted to be removed from the data before it is shared. These objectives and requirements are specified, e. g., by the Health Insurance Portability and Accountability Act (HIPAA)[24] in the United States or the General Data Protection Regulation (GDPR)[25] in Europe.

In 2017, Google entered into a project with the National Institutes of Health (NIH) to publish a dataset containing 100,000 chest radiographs. However, the release was canceled two days before publication after Google was informed by the NIH that the radiographs still contained personal information which indicates that the data was incorrectly anonymized[26,27]. This major incident highlights that many potential pitfalls can arise when clinical and technological institutions collect and share large medical datasets to revolutionize health-care.

In the past, various data de-identification techniques have been proposed, including commonly-used methods like pseudonymization[28] and $k$-anonymity[29]. Pseudonymization describes a technique that replaces a true identifier, e. g., the name or the patient identification number by a pseudonym that is unique to the patient but has no relation to the person[28]. However, pseudonymization is a rather weak anonymization technique as the patient's identity may still be revealed, e. g., by cross-referencing with other publicly available datasets. In contrast, $k$-anonymity modifies the data before sharing in such a way that every sample in the published dataset can be associated with at least $k$ different subjects. In this way, the probability of performing identity disclosure is limited to at most $^1/k$[29,30]. Nevertheless, when background knowledge is available, $k$-anonymity is susceptible to many attacks.

To date, little attention has been paid to the possibility of re-identifying patients in large medical datasets by means of DL techniques. However in theory, medical data disclosure, as illustrated in Figure 1, could be facilitated for potential attackers by using suitable DL approaches. Consider a publicly available dataset that is supposedly anonymized but contains further sensitive patient-related information, e. g., diagnosis, treatment history, and clinical institution. If a radiograph of known identity is accessible to a potential attacker and a properly working verification or re-identification model exists, then the model could be used to compare the given radiograph to each image in the dataset which would essentially result in a set of images

**Table 1.** Overview of the obtained verification results for our experiments using varying training set sizes $N_s$ at different learning rates $\eta$. Moreover, different data handling techniques were used (FTS: Fixed training set; RNP: Randomized negative pairs). For each experiment, the training sets were balanced with respect to the amount of positive and negative image pairs. In this table, we present the area under the curve (AUC) (together with the lower and upper bounds of the 95 % confidence intervals from 10,000 bootstrap runs), the accuracy, the specificity, the recall, the precision, and the F1-score. Bold text emphasizes the overall highest AUC value.

| Data handling | $N_s$ | $\eta$ | AUC + 95 % CI | Accuracy ($\frac{TP+TN}{P+N}$) | Specificity ($\frac{TN}{N}$) | Recall ($\frac{TP}{P}$) | Precision ($\frac{TP}{TP+FP}$) | F1-score |
|---|---|---|---|---|---|---|---|---|
| FTS | 100,000 | $10^{-3}$ | 0.8610 $^{0.8632}_{0.8588}$ | 0.7782 ($\frac{77,815}{100,000}$) | 0.7710 ($\frac{38,548}{50,000}$) | 0.7853 ($\frac{39,267}{50,000}$) | 0.7742 ($\frac{39,267}{50,719}$) | 0.7797 |
| | 200,000 | $10^{-3}$ | 0.9448 $^{0.9461}_{0.9435}$ | 0.8743 ($\frac{87,428}{100,000}$) | 0.8685 ($\frac{43,426}{50,000}$) | 0.8800 ($\frac{44,002}{50,000}$) | 0.8700 ($\frac{44,002}{50,576}$) | 0.8750 |
| | 400,000 | $10^{-4}$ | 0.9587 $^{0.9599}_{0.9575}$ | 0.8755 ($\frac{87,546}{100,000}$) | 0.9290 ($\frac{46,452}{50,000}$) | 0.8219 ($\frac{41,094}{50,000}$) | 0.9205 ($\frac{41,094}{44,642}$) | 0.8684 |
| | 800,000 | $10^{-4}$ | 0.9896 $^{0.9901}_{0.9891}$ | 0.9537 ($\frac{95,367}{100,000}$) | 0.9541 ($\frac{47,705}{50,000}$) | 0.9532 ($\frac{47,662}{50,000}$) | 0.9541 ($\frac{47,662}{49,957}$) | 0.9536 |
| RNP | 800,000 | $10^{-4}$ | **0.9940** $^{0.9944}_{0.9937}$ | 0.9555 ($\frac{95,545}{100,000}$) | 0.9822 ($\frac{49,111}{50,000}$) | 0.9287 ($\frac{46,434}{50,000}$) | 0.9812 ($\frac{46,434}{47,323}$) | 0.9542 |

belonging to the same patient (patient verification) or yield a ranked list of the most similar images to the given radiograph (patient re-identification). In this way, the patient's identity may be linked to sensitive data contained in the dataset. As a result, more patient-related information may have been leaked, highlighting the enormous data security and data privacy issues involved.

In our work, we investigated whether conventional anonymization techniques are secure enough and whether it is possible to re-identify and de-anonymize individuals from their medical data using DL-based methods. Therefore, we considered the public ChestX-ray14 dataset[16], which is one of the most widely used research datasets in the field of radiographic problems. Our algorithms are trained to determine whether two arbitrary chest radiographs can be recognized to belong to the same patient or not. Moreover, we showed that our proposed methods are able to perform a successful linkage attack on publicly available chest radiography datasets. Furthermore, this work aims to draw attention to the massive problem of releasing medical data without considering that DL systems can easily be used to reveal a patient's identity. Therefore, we call for reconsidering conventional anonymization techniques and developing more secure methods that resist potential attacks by DL algorithms.

## Patient Verification

First, we trained a siamese neural network (SNN) architecture on the ChestX-ray14 dataset to determine whether two individual chest radiographs correspond to the same patient or not. Our model was designed to process the two input images in two identical network branches, which are then combined by a merging layer. The fused information is fed through further network layers resulting in a single output score indicating the identity similarity.

Table 1 summarizes the outcomes of our evaluation. We analyzed a multitude of different experimental setups with varying learning rates $\eta$ and differing balanced training set sizes $N_s$. Moreover, we investigated the effect of using epoch-wise randomized negative pairs (RNP) versus fixed training sets (FTS) for the entire learning procedure. When using RNP as the data handling technique, the negative image pairs were randomly constructed in each epoch, meaning that much more negative pairs could be utilized in a complete training run compared to FTS where the generated image pairs remain the same for the entire learning procedure. For all experiments on the ChestX-ray14 dataset, we used the same balanced validation and testing set with 50,000 and 100,000 image pairs, respectively, without patient overlap between any split. To assess the performance of the trained models, we performed a receiver operating characteristic (ROC) analysis by computing the AUC value together with the 95 % confidence intervals from 10,000 bootstrap runs. Moreover, we calculated the accuracy, specificity, recall, precision, and F1-score.

The results indicate that the amount of training data plays a crucial role in the patient verification task. We observe a significant performance increase as the training set size grows. For instance, when using a subset of 100,000 image pairs for training, we obtain an AUC value of 0.8610. In contrast, by enlarging the training set size to 800,000 image pairs (i. e. the total number of 400,000 positive image pairs combined with 400,000 negative pairs), we receive an AUC score of 0.9896. These findings have been visualized in the ROC curves shown in Figure 2 which illustrates the effect of the training set size on the verification performance when using fixed training sets. Note that Table 1 only shows the best experiments per training set size $N_s$. Additional experiments were conducted to investigate the effect of the learning rate (LR). The corresponding results are provided in a separate table in the appendix (see Supplementary Table 1).

We also observed that randomly constructing the negative image pairs in each epoch led to further improvements in the final model performance. By using this data handling technique, we achieved our overall best results. The respective outcomes are reported in Table 1. When training our network architecture with a total of 800,000 training samples with epoch-wise randomly
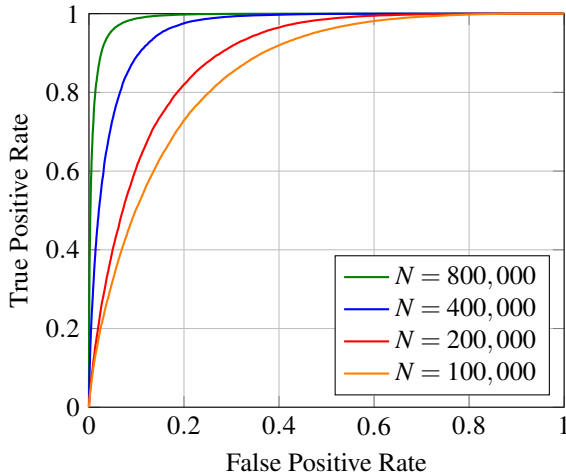
**Figure 2.** ROC curves for different training set sizes $N_s$ and a fixed LR of $\eta = 10^{-4}$. During training, the fixed data handling technique was employed.

|  |  | **Predicted** | | |
|---|---|---|---|---|
|  |  | neg | pos | **Total** |
| **Actual** | neg | $49,111$ (TN) | $889$ (FP) | $50,000$ |
|  | pos | $3,566$ (FN) | $46,434$ (TP) | $50,000$ |
|  | **Total** | $52,677$ | $47,323$ | $100,000$ |

**Figure 3.** Confusion matrix corresponding to the best experiment shown in Table 1 (last row) giving clear insights into the performance of our trained model.
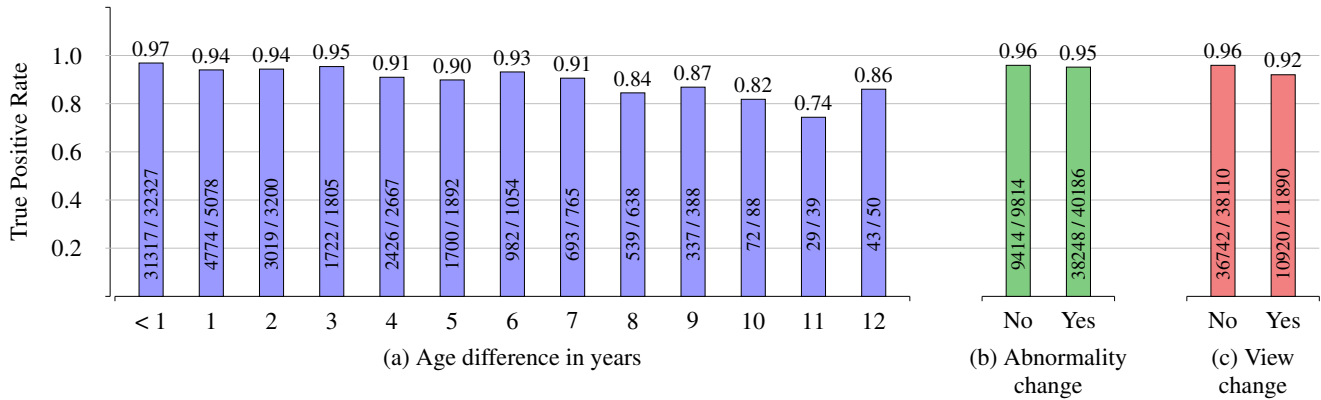


**Figure 4.** True positive rates (TPRs) for image pairs with (a) age differences, (b) with changes in the disease pattern, and (c) with changes in the projection view. The absolute numbers of true positives and overall positives are given for each bin. Note that the number of image pairs with age differences of more than 12 years is comparatively small, which is why the corresponding TPRs are neglected in this figure.

constructed negative pairs, the AUC score improved from 0.9896 to 0.9940. Besides, the other reported evaluation metrics apart from the recall also increased compared to the results achieved by the model trained with the fixed set. Figure 3 depicts the confusion matrix resulting from our best-trained model listed in Table 1 (last row), thus giving clear insights into the patient verification performance.

We also analyzed how the model with the best recall (fourth row in Table 1) behaves when comparing images of the same patient where the acquisition dates are several years apart. The results are illustrated in Figure 4a. We received a TPR of 0.97 for image pairs that had small age differences of less than one year. As the age variation between the follow-up images and the initial scan increases, we observe a slight decrease in the TPR values. Nevertheless, our model still shows competitive results even if the patient's age in two images differs by several years. Even for an age difference of twelve years, we can verify that two images belong to the same patient by 86 %. We only report the TPRs for image pairs with follow-up intervals of up to 12 years in Figure 4a as the number of pairs with larger intervals is relatively small.

Additionally, we investigated the model's verification capability in the case of new abnormality patterns appearing in follow-up scans that did not occur in previously acquired chest radiographs. Figure 4b shows that regardless of the abnormality, we nearly observe no decline in the TPR values, emphasizing the robustness of our trained SNN architecture. Note that the disease labels in the ChestX-ray14 dataset were extracted using natural language processing techniques. This could potentially
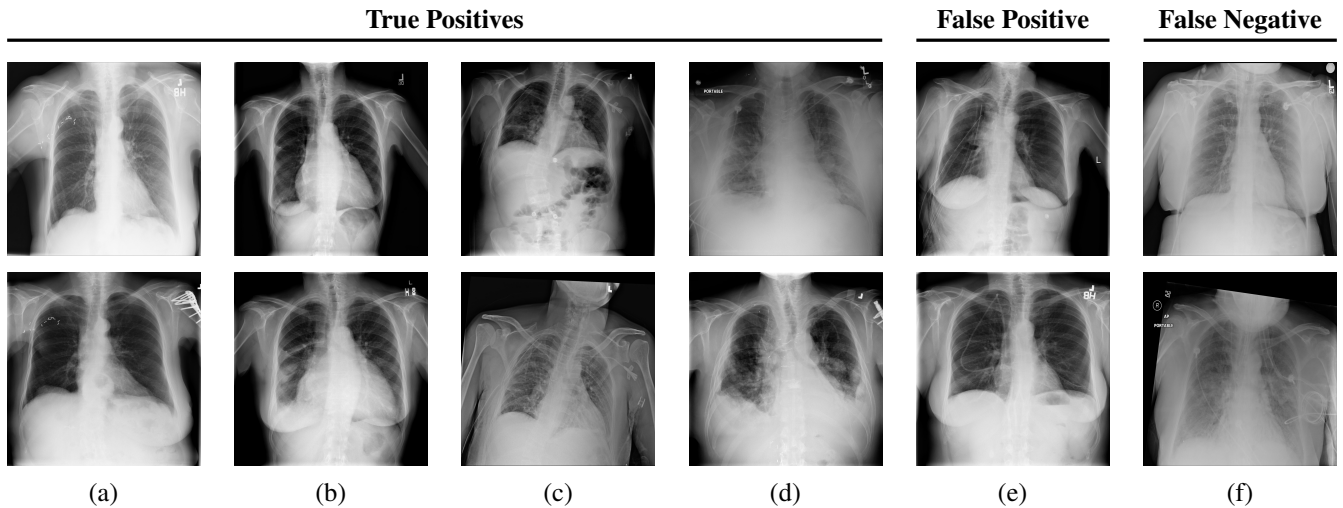
|  | **True Positives** | | | | **False Positive** | **False Negative** |
|---|---|---|---|---|---|---|
| (a) | (b) | (c) | (d) | (e) | (f) | |

**Figure 5.** Exemplary image pairs are classified by our best performing verification model. Each column represents one image pair. The first four columns (a)–(d) show true positive classifications. The last two columns (e) and (f) depict a false positive and a false negative classification, respectively.

**Table 2.** Comparison of the verification performance on two different subsets of the ChestX-ray14 dataset that either contain foreign material or not (first two rows). Furthermore, we show the verification results for the CheXpert dataset and the COVID-19 Image Data Collection (last two rows). We present the AUC (together with the lower and upper bounds of the 95 % confidence intervals from 10,000 bootstrap runs), the accuracy, the specificity, the recall, the precision, and the F1-score.

| Dataset | Subset | AUC + 95 % CI | Accuracy ($\frac{TP+TN}{P+N}$) | Specificity ($\frac{TN}{N}$) | Recall ($\frac{TP}{P}$) | Precision ($\frac{TP}{TP+FP}$) | F1-score |
|---|---|---|---|---|---|---|---|
| ChestX-ray14 | w/ foreign material | $0.9970 \; {}^{0.9993}_{0.9938}$ | $0.9796 \; \left(\frac{672}{686}\right)$ | $0.9854 \; \left(\frac{338}{343}\right)$ | $0.9738 \; \left(\frac{334}{343}\right)$ | $0.9853 \; \left(\frac{334}{339}\right)$ | 0.9795 |
| | w/o foreign material | $0.9972 \; {}^{0.9999}_{0.9909}$ | $0.9862 \; \left(\frac{430}{436}\right)$ | $0.9908 \; \left(\frac{216}{218}\right)$ | $0.9817 \; \left(\frac{214}{218}\right)$ | $0.9907 \; \left(\frac{214}{216}\right)$ | 0.9862 |
| CheXpert | - | $0.9870 \; {}^{0.9884}_{0.9855}$ | $0.9440 \; \left(\frac{15,562}{16,486}\right)$ | $0.9629 \; \left(\frac{7,937}{8,243}\right)$ | $0.9250 \; \left(\frac{7,625}{8,243}\right)$ | $0.9614 \; \left(\frac{7,625}{7,931}\right)$ | 0.9429 |
| COVID-19 | - | $0.9763 \; {}^{0.9825}_{0.9696}$ | $0.9180 \; \left(\frac{1,421}{1,548}\right)$ | $0.9780 \; \left(\frac{757}{774}\right)$ | $0.8579 \; \left(\frac{664}{774}\right)$ | $0.9750 \; \left(\frac{664}{681}\right)$ | 0.9127 |

have caused label noise which could have affected the results shown in Figure 4b. Furthermore, Figure 4c illustrates that changes in the projection view (e. g., one image taken in the anterior-posterior position and the other image acquired using the posterior-anterior view) hardly lead to any deterioration in the performance.

Moreover, we perform a qualitative evaluation where we visually inspect some exemplary image pairs evaluated using our best-performing verification model. In Figure 5, we show four true positive (TP) classifications (a)–(d), one pair that has been classified as a false positive (FP) (e), and one example for a false negative (FN) image pair (f). The shown images clearly illustrate the high technical variance present in the ChestX-ray14 dataset. The first image pair (a) shows two images belonging to the same patient with a difference of seven years. Clear differences in pixel intensities and lung shape are observed. However, both images belong to the same person, cf. the small vascular clips in the area of the upper right lung. Also, image pairs with large difference in scaling (b) or rotation (c) are verified correctly. Our model is also robust to the patients' pathology: While the upper image of (b) shows characteristics of pneumothorax, the patient suffered from cardiomegaly, effusion, and masses in the lower image, according to the provided annotations. Similarly in (c), where the upper image indicates the presence of infiltration and pneumothorax, whereas the lower scan shows signs of infiltration and nodules. Figure 5e shows an exemplary image pair that has falsely been classified as positive. Conversely, (f) depicts a positive image pair that has been incorrectly classified as negative. To visually demonstrate which parts of the images are responsible for the verification task, we applied a siamese attention mechanism[31] to our network architecture which utilizes the Grad-CAM algorithm[32]. The obtained attention maps can be seen in Figures 6 and 7. They clearly indicate that the human anatomy, especially the shape of the lungs and ribs, is the driving factor for the network decisions.

To investigate how foreign material (see Figure 5a) affects the verification performance, we evaluated our trained network on two small manually created subsets of around 200 images. The first one consisted only of images in which foreign material is visible, whereas the second one solely contained images without foreign material. When constructing the subsets, we selected the patients at random and then assigned the corresponding patient images to the respective subset after visual assessment. Furthermore, we ensured that no more than 5 images were used per patient. Table 2 summarizes the results indicating that the
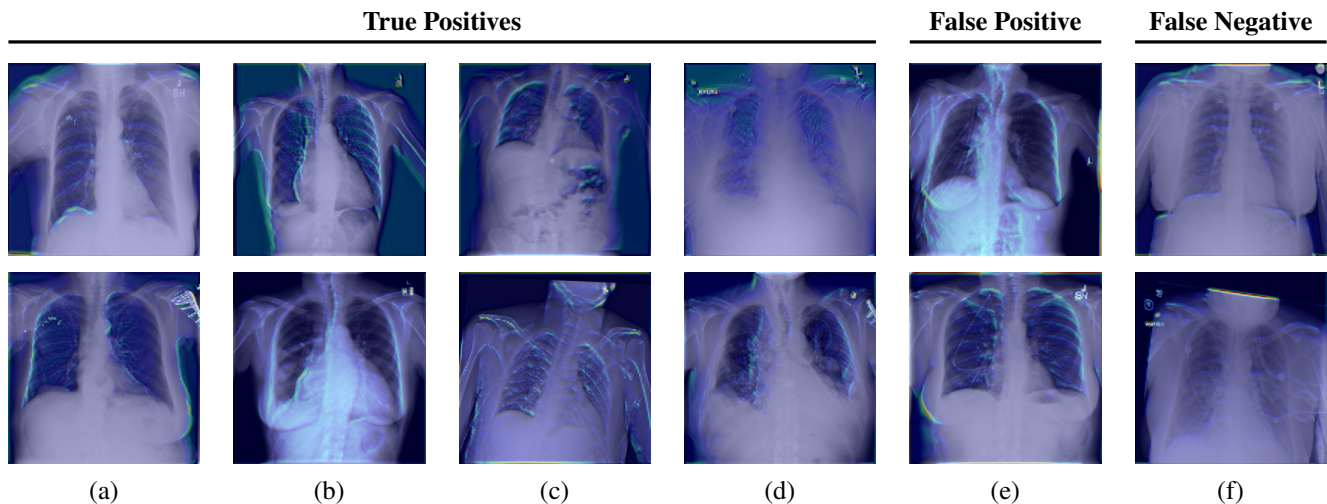
**Figure 6.** Grad-CAM visualizations for the first convolutional layer of the ResNet-50 incorporated in our SNN. Each column represents one image pair. The first four columns (a)–(d) show true positive classifications. The last two columns (e) and (f) depict a false positive and a false negative classification, respectively. The shown images illustrate that the anatomical structure of, e. g. the breast (cf. (a), (b), (e)), the lungs (cf. (a), (b), (c), (e)), and the heart (cf. (a), (b)) have a high impact on the final model prediction. Furthermore, it can be seen that our network focuses on the collarbones (cf. (a), (c), (d), (f)) and the ribs (cf. (b), (c), (e)). The upper images of (a) and (f) also highlight that our network pays attention to the contour of the diaphragm.
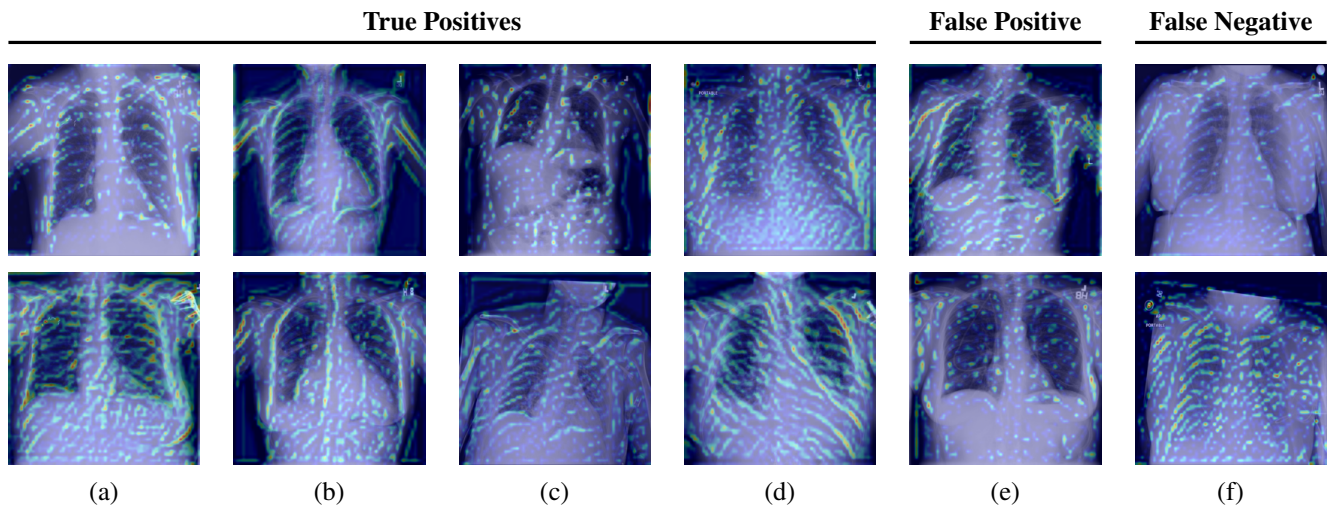


**Figure 7.** Grad-CAM visualizations for an intermediate convolutional layer of the ResNet-50 incorporated in our SNN. Each column represents one image pair. The first four columns (a)–(d) show true positive classifications. The last two columns (e) and (f) depict a false positive and a false negative classification, respectively. The obtained attention maps clearly illustrate that the selected network layer focuses on the ribs and the outline of the thorax.

patient verification works with high performance regardless of the occurrence of foreign material. We even observe a slight improvement in performance for the subset where no foreign material is visible in the images.

Finally, to analyze whether our trained model is able to generalize to other datasets which have not been used during training, we evaluated our network on external datasets such as CheXpert and the COVID-19 Image Data Collection. For this, we utilized 16,486 image pairs from the CheXpert dataset and 1,548 image pairs from the COVID-19 Image Data Collection. The results are summarized in the last two rows of Table 2. It can be seen that our network still yields high AUC values of 0.9870 (CheXpert) and 0.9763 (COVID-19) for the verification task although the model has not been fine-tuned on the respective datasets. Also the other presented evaluation metrics show competitive values without deteriorating too much.

**Table 3.** Overview of the obtained results for our image retrieval experiments. In this table, we report the *mAP@R*, the *R*-Precision, and the Precision@1. The first 4 rows show the results on the ChestX-ray14 dataset for different image resolutions used for evaluation. The fifth row shows the outcomes on the CheXpert dataset. The last row indicates the results on the COVID-19 Image Data Collection. Bold text represents the overall highest performance metrics.

| Dataset | Input dimensions | mAP@R | R-Precision | Precision@1 |
|---|---|---|---|---|
| ChestX-ray14 | 1024×1024 | **0.9748** | **0.9763** | **0.9963** |
| | 800×800 | 0.9709 | 0.9726 | 0.9958 |
| | 512×512 | 0.9572 | 0.9601 | 0.9945 |
| | 224×224 | 0.7730 | 0.7979 | 0.9756 |
| CheXpert | original | 0.8001 | 0.8148 | 0.9444 |
| COVID-19 | original | 0.8569 | 0.8707 | 0.8821 |

**Table 4.** Comparison of the re-identification performance on two different subsets (ChestX-ray14) that either contain foreign material or not. We report the *mAP@R*, the *R*-Precision, and the Precision@1.

| Subset | mAP@R | R-Precision | Precision@1 |
|---|---|---|---|
| w/ foreign material | 0.9925 | 0.9925 | > 0.9999 |
| w/o foreign material | > 0.9999 | > 0.9999 | > 0.9999 |

## Patient Re-identification

For our patient re-identification experiments, we trained another SNN architecture on the ChestX-ray14 dataset. In contrast to the verification model, we omitted all the layers from the merging layer onwards. The main objective was to learn appropriate feature representations instead of directly determining whether the inputs belong to the same patient or not. After training the network, we used the ResNet-50 backbone as a feature extractor for the actual image retrieval task. By computing the Euclidean distance between the embeddings of the query image and each other image, we obtained for each query image a ranked list of its most similar images in terms of identity. The used training, validation, and testing set consisted of 61,755, 10,815, and 25,596 images, respectively.

The results of the corresponding image retrieval experiments are summarized in Table 3. When using the original image size of 1024×1024 pixels for evaluation, we obtain a precision@1 of more than 99 % showing that the closest match nearly always is the same patient. The high mean average precision at *R* (mAP@R) of about 97 % further depicts that most of the most similar images are correctly identified. We observe a slight decrease in performance as the image size was reduced. Nevertheless, when the images were downsampled to a resolution of 512×512 pixels, we still obtain high performance values. When the image size was reduced too aggressively, e. g., to 224×224 pixels, the *mAP@R* and the *R*-Precision rates drop. Yet, we still observed a high Precision@1 of more than 97 %.

Similar to the experiments in the patient verification section, we evaluated our best-trained re-identification model on two small subsets, one of which only contained images with visible foreign material and the other consisted exclusively of images without the presence of foreign material. The obtained results are presented in Table 4. As can be seen, we achieve high performance values for both subsets. Thus, we hypothesize that our outcomes are independent of foreign material which may occur only for specific patients.

Lastly, we analyzed the re-identification performance on external datasets such as CheXpert and the COVID-19 Image Data Collection. For this, we utilized 6,454 images from the CheXpert dataset and 781 images from the COVID-19 dataset. As can be seen in the last two rows of Table 3, we also obtain high retrieval values although we haven't performed any fine-tuning on both datasets, which demonstrates the feasibility of the trained re-identification network on previously unseen datasets.

## Discussion and Conclusion

In this paper, we investigated the patient verification and re-identification capabilities of DL techniques on chest radiographs. We have shown that well-trained SNN architectures are able to compare two individual frontal chest radiographs and reliably predict whether these images belong to the same patient or not. Moreover, we have shown that DL models have the potential to accurately retrieve relevant images in a ranked list. Our models have been evaluated on the publicly available ChestX-ray14 dataset and showed competitive results with an AUC of up to 0.9940 and classification accuracy of more than 95 % in the verification scenario and an *mAP@R* of 97 % and a precision@1 of about 99 % in the image retrieval scenario. Especially the

fact that basic SNNs have the capability to re-identify patients despite potential age differences, disease changes or differing projection views demonstrated the effectiveness of DL techniques for this task. However, note that the shown results were obtained empirically, i. e. they do not necessarily reflect true measures of certainty.

As shown in Figure 5, the used dataset suffers from a high technical variance which may occur due to various windowing techniques applied to the images. In a real-life scenario, the resulting variations in image contrast and brightness could be significantly mitigated by using dynamic normalization approaches[33]. Furthermore, we believe that variations in rotation and scaling can be counteracted by appropriate alignment algorithms. Nevertheless, even without such pre-processing steps, we were able to show that patient matching for chest radiographs is possible with a high performance by using DL techniques.

Moreover, we hypothesize that special noise patterns characteristic for unique patients appear in the images which might unintentionally improve the re-identification performance. For example, the initial anonymization strategy may be biased towards the clinical institution and, therefore, also towards follow-up images. To get a better impression of the re-identification capability of our SNN architecture, we also intend to investigate other datasets which show less or ideally no correlation between potential noise patterns and the patient identity. Therefore, further research on multiple datasets should ideally be considered. For our experiments, we already evaluated our models on two completely different datasets, the CheXpert dataset and the COVID-19 Image Data Collection. While the evaluation metrics are lower (possibly due to a domain shift or the severity of diseases), we still obtain AUC scores of over 97 % (COVID-19) and 98 % (CheXpert) and precision@1 values of more than 88 % (COVID-19) and 94 % (CheXpert) without fine-tuning on these datasets. This indicates that patient verification and re-identification is also applicable for data that was acquired in various hospitals around the world where other pre-processing steps may be taken before data publication compared to the ChestX-ray14 dataset.

The COVID-19 Image Data Collection is very heterogeneous, containing, e. g., images of different sizes, both gray-scale and color images, and images with visible markers, arrows or date displays. For our experiments, only those images in the COVID-19 Image Data Collection were used that were acquired using the anterior-posterior or the posterior-anterior view, while images taken in the lateral position and CT scans were discarded. Apart from this, no further steps were taken to ensure the quality of the dataset. Although some of the factors mentioned above (e. g., brandings such as markers, arrows or dates) may facilitate the patient re-identification, we hypothesize that the COVID-19 Image Data Collection poses a realistic example of a public medical dataset and we therefore consider the conducted experiment as an authentic real-life application scenario.

Furthermore, we want to accentuate that our trained network architectures are able to handle non-rigid transformations that may appear between two images of the same person in the ChestX-ray14 dataset. Such deformations can occur due to various breath states in follow-up scans or due to different positioning during X-ray acquisition. Hence, the shape of the heart and lungs, or the contours of the ribs may appear deformed compared to an initial scan. The obtained results lead to the assumption that our trained SNN architectures can withstand such deformations and can therefore be used for reliable patient re-identification on chest radiographs.

We conclude that DL techniques render medical chest radiographs biometric for everyone and allow a re-identification similar to a fingerprint. Therefore, publicly available medical chest X-ray data is not entirely anonymous. Using a DL-based re-identification network enables an attacker to compare a given radiograph with public datasets and to associate accessible metadata with the image of interest. The strength of our proposed method is that patients can be re-identified in a fully automated way without the need of expert knowledge. Thus, sensitive patient data is exposed to a high risk of falling into the unauthorized hands of an attacker who may disseminate the gained information against the will of the concerned patient. At this point, we want to emphasize that data leakage of this kind requires that the attacker has previously gained access to an image of a known person. This could happen, for example, through a stolen CD containing raw medical data of a specific patient, or through accidental data release by a radiological facility. Furthermore, data breaches due to inadequate data security measures at, e. g., healthcare institutions or health insurance companies, represent a possibility for attackers to obtain images of known patients, which could subsequently be utilized for a linkage attack as presented in our work. However, even if the attacker owns an image of an unknown identity, a re-identification model can be used to find the same patient across various datasets. Assuming multiple datasets contain the same patient but different metadata, an attacker would be able to obtain a more complete picture of the respective patient. We hypothesize that collecting patient information by this means could significantly help an attacker infer the true identity of the patient. We therefore urge that conventional anonymization techniques be reconsidered and that more secure methods be developed to resist the potential attacks by DL-based algorithms.

At this point, we would like to draw attention to the analogy of our work to the field of automatic speaker verification (ASV). Speech signals contain a large amount of private data, e. g., age, gender, health and emotional state, ethnic origin, and more[34]. As such information is embedded in speech data itself, it can be exploited to reveal the speaker's identity by applying attack models in the form of ASV systems[34]. Therefore, in the speech community, raw speech signals are not considered anonymous. Instead, privacy challenges, such as the VoicePrivacy Challenge[35], were formed to develop solutions for the preservation of privacy in this field. With our work, we were able to demonstrate that the privacy issue with speech data is 1-to-1 transferable to the privacy issue with chest radiographs.

While the proposed algorithms may be used maliciously to produce harm in terms of patient privacy and data anonymity, we also want to draw attention to a potentially positive application area. Often, different datasets are used for training and evaluation of an algorithm. However, since in most cases these datasets have been anonymized using conventional techniques, it is not clear whether certain patients appear in more than one dataset. Therefore, in this context, our trained networks could be applied to check for mutual exclusiveness with respect to included patients between multiple datasets.

The publication of medical image datasets is an area of conflict. While, on the one hand, many patients may benefit from recent advances (e. g., the development of diagnostic algorithms), there are, on the other hand, patients who may be seriously harmed by the fact that their data is publicly available. With our work, we focus on providing empirical evidence for this issue and draw attention to the risks. The legal situation for the publication of medical data is currently regulated by the HIPAA (in the United States) and the GDPR (in Europe). We therefore contend that the corresponding ethics committees are responsible for weighing the benefits and the risks as well as for assessing the appropriateness of current regulations.

Potential solutions to the problems addressed in our work may be found in privacy-preserving approaches such as differential privacy (DP)[36,37] which guarantees that the global statistical distribution of a dataset is retained while individually recognizable information is reduced[38]. This means that an outside observer is unable to draw any conclusions about the presence or absence of a particular individual. Consequently, algorithms trained with DP are able to withstand linkage attacks attempting to reveal the identities of patients in the dataset used to train the algorithm. One commonly-used technique to achieve DP is to modify the input by adding noise to the dataset (local DP)[38,39]. Furthermore, DP can be applied to the computation results (global DP) or to algorithm updates[38]. However, training models with DP degrades the quality of the model (privacy-utility trade-off) which is problematic in medicine where high diagnostic utility is required. Therefore, further exploration on these topics is necessary before general conclusions can be made.

Aside from perturbation-based privacy approaches, we want to mention that the use of collaborative decentralized learning protocols such as federated learning (FL)[40] can significantly contribute to a safer use of medical data. By training a machine learning model collaboratively without centralizing the data, the need of raw data sharing or dataset release is eliminated[41]. Thus, the medical data is able to reside with its owner, e. g., the healthcare institution where the data was acquired, which resolves data governance and ownership issues[38,42]. However, FL itself does not provide full data security and privacy, meaning that some risks remain unless combined with other privacy-preserving methods.

## Methods

The research was carried out in accordance with the relevant guidelines and regulations of the institution conducting the experiments.

### Siamese neural networks

To re-identify patients from their chest radiographs, we employ SNN architectures for both the classification and the retrieval tasks. A SNN receives two input images which are processed by two identical feature extraction blocks sharing the same set of network parameters. The resulting feature representations can then be used to compare the inputs. The concept of a SNN was initially introduced by Bromley et al.[43] for handwritten signature verification. Taigman et al.[44] applied this idea in the field of face verification and proposed the *DeepFace* system. Moreover, Koch et al.[45] presented an approach for one-shot learning on the Omniglot[46] and MNIST[47] datasets.

### NIH ChestX-ray14 dataset

With a total of 112,120 frontal-view chest radiographs from 30,805 unique patients, the NIH ChestX-ray14[16] dataset counts to one of the largest publicly available chest radiography datasets in the scientific community. Due to follow-up scans, the image collection provides an average of 3-4 images per patient. The originally acquired radiographs were published as 8-bit gray-scale PNG images with a size of $1024 \times 1024$ pixels. Associated metadata is available for all images in the dataset. The additional data comprises information about the underlying disease patterns (either no finding or a combination of up to 14 common thoracic pathologies), the number of follow-up images taken, the patients' age and gender, and the projection view (anterior-posterior or posterior-anterior) used for radiography acquisition. According to the publisher, the dataset was carefully screened to remove all personally identifiable information before release[48]. Therefore, the patient names were replaced by integer IDs. Moreover, personal data in the image domain itself has been made unrecognizable by placing black boxes over the corresponding image areas.

### CheXpert dataset

The CheXpert[14] dataset contains 224,316 frontal and lateral chest radiographs of 65,240 patients, who underwent a radiographic examination from Standford University Medical Center between October 2002 and July 2017. The originally acquired

radiographs were published as 8-bit gray-scale JPG images with varying image resolutions. Note that only frontal chest radiographs were used in our work, whereas lateral images were excluded.

## COVID-19 Image Data Collection

The COVID-19 Image Data Collection[19] is a dataset that was created and published as an initiative to provide COVID-19 related chest radiographs and CT scans for machine learning tasks. It comprises data of 448 unique patients and a total of around 950 images with different image resolutions. In this work, only the available frontal radiographs were utilized, whereas the lateral images and CT scans were discarded.

## Data preparation

Since SNN architectures require pairs of images for training and evaluation, we construct both positive and negative image pairs from the images contained in the ChestX-ray14 dataset. In this context, a positive pair consists of two images belonging to the same patient, whereas a negative pair comprises two images that belong to different patients. Mathematically, the constructed dataset can be described according to

$$\mathscr{S} = \{(\boldsymbol{x}_{11}, \boldsymbol{x}_{12}, y_1), ..., (\boldsymbol{x}_{m1}, \boldsymbol{x}_{m2}, y_m), ..., (\boldsymbol{x}_{M1}, \boldsymbol{x}_{M2}, y_M)\} \ , \ \text{with } y_m \in \{0, 1\} \ , \tag{1}$$

where the triplet $(\boldsymbol{x}_{m1}, \boldsymbol{x}_{m2}, y_m)$ represents one sample consisting of two images $\boldsymbol{x}_{m1}$ and $\boldsymbol{x}_{m2}$, and the corresponding label $y_m$. $M$ refers to the total number of samples and $m$ denotes an iterator variable in the range of $m \in [1, M]$. The class label $y_m$ symbolizes a binary variable that takes the value 0 for negative image pairs and 1 for positive image pairs.

To ensure that images from one patient only appear either in the training, validation, or testing set, we use the patient-wise splitting strategy. According to the official split provided with the ChestX-ray14 dataset, the data is randomly divided into 70 % training, 10 % validation, and 20 % testing. Based on this split, we construct the actual image pairs for each subset.

### Offline mining

For patient verification, we follow an offline mining approach, meaning that the positive and negative image pairs are generated once before conducting the experiments. First, the positive pairs are generated by only considering the patients for whom multiple images exist in the respective subset. For each patient with follow-up images, we produce all possible tuple combinations assuming the images to be unique. By following this approach, we are able to construct a total of around 400,000 positive image pairs for our training set. The negative pairs in each subset are randomly generated and concatenated with the respective positive pairs afterwards.

### Online mining

For the patient re-identification experiments, we choose an online mining approach, meaning that image pairs are formed in each batch during the training procedure. This means that the embeddings of all batch images are first computed and then subsequently used in all possible combinations as input for the loss function. Moreover, all patients with only one available image were discarded from the training set.

## Patient verification

### Deep learning architecture

For patient verification, the used SNN architecture (see Figure 8) receives two images $\boldsymbol{x}_1$ and $\boldsymbol{x}_2$ of size $3 \times 256 \times 256$. Both inputs are processed by a pre-trained ResNet-50 incorporated in each network branch. In its original version, the ResNet-50 was designed to classify images into 1,000 object categories trained on the ImageNet[49] dataset. To adapt the ResNet-50 to our specific needs, we replace its classification layer with a layer consisting of 128 output neurons producing the feature representations $\boldsymbol{z}_1$ and $\boldsymbol{z}_2$, respectively. To merge both network branches, the absolute difference of the sigmoid activations of the two feature vectors is computed. We add a fully-connected (FC) layer to reduce the dimensionality to one neuron, followed by another sigmoid activation function $\sigma$ which yields the final output score $\hat{y} \in [0, 1]$.

### Training strategy

The verification model is trained using the binary cross-entropy (BCE) loss. The network parameters are optimized by combining mini-batch stochastic gradient descent (SGD)[50,51] with the adaptive moment estimation (Adam)[52] method. The batch size $N_b$ is set to 32 in all our experiments. We use different LRs to investigate their effect on the model's performance. Furthermore, we include an early stopping criterion with a patience $p = 5$, which means that the training procedure stops as soon as the network does not improve for 5 epochs. We train the architecture using input dimensions of $3 \times 256 \times 256$.
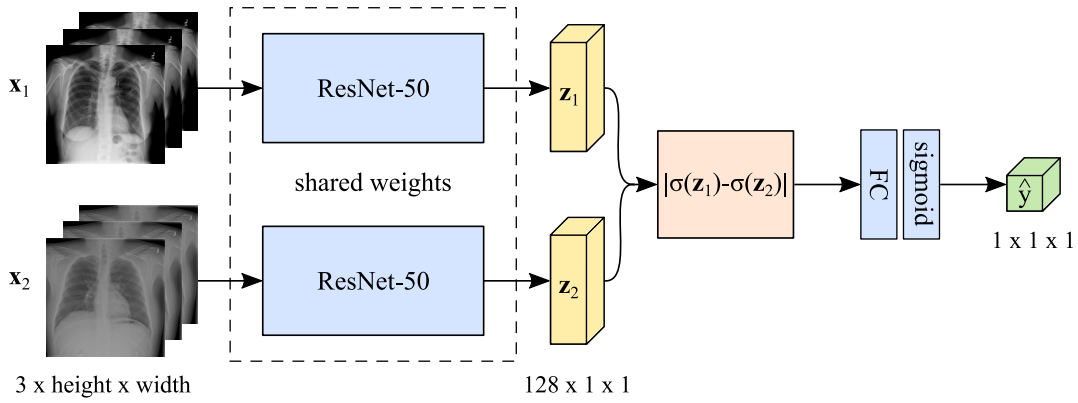
**Figure 8.** SNN architecture used for patient verification on the ChestX-ray14[16] dataset. The feature extraction blocks (light blue) share the same set of network parameters and produce the feature representations $z_1$ and $z_2$ (yellow). After merging (orange) and an additional FC and sigmoid layer $\sigma$ (light blue), the network yields the final output score $\hat{y}$ (green). For our patient re-identification experiments, we used the same architecture but ejected all the layers from the merging layer onwards.

### Evaluation techniques

We utilize ROC curves to visualize the trained verification models based on their performance. A ROC curve represents a two-dimensional graph in which the TPR is plotted against the false positive rate (FPR) at various threshold settings[53], thus indicating how many true positive classifications can be gained as an increasing number of false positive classifications is allowed. Additionally, we calculate the AUC which reflects a proportion of the area of the unit square and will always range from 0 to 1[53]. The higher the AUC score, the better the model's average performance. Nevertheless, it has to be mentioned that a classifier with a high AUC might perform worse in a specific region of ROC space than a classifier with a low AUC value. Moreover, we evaluate the performance by computing the accuracy, specificity, recall, precision and F1-score. Therefore, the threshold at the output neuron is set to $t = 0.5$.

## Patient re-identification
### Deep learning architecture

For patient re-identification, we train a SNN architecture which receives two images $x_1$ and $x_2$ of size $3 \times 1024 \times 1024$. Both inputs are processed by a pre-trained ResNet-50 incorporated in each network branch. However, the network head of the used ResNet-50 is slightly modified. The average pooling layer is replaced by an adaptive average pooling layer producing feature maps of size $5 \times 5$. In addition to the adaptive average pooling layer, an adaptive max-pooling layer is applied which also yields feature maps of size $5 \times 5$. The outputs of the pooling layers are concatenated and processed by a $1 \times 1$ convolutional layer reducing the number of feature maps from 2048 to 100. The feature maps are then flattened, followed by two successive FC layers resulting in 128-dimensional feature representations $z_1$ and $z_2$ for the first and the second network branch.

### Training strategy

The re-identification model is trained using the contrastive loss function[54] which is typically utilized to achieve a meaningful mapping $F$ from high to low dimensional space. By using the contrastive loss, the network learns to map similar inputs to nearby points on the output manifold while dissimilar inputs are mapped to distant points. Negative pairs contribute to the loss only if their distance is smaller than a certain margin $m$. In this work, the margin is set to $m = 1$.

For our image retrieval experiments, the SNN architecture is optimized using the SGD algorithm in combination with the 1cycle learning policy[55,56]. When using the 1cycle LR schedule, the LR $\eta$ steadily increases until it reaches a chosen maximum value and gradually decreases again thereafter. This schedule changes the LR after every single batch and is pursued a pre-defined number of epochs. The upper bound is chosen at 0.1584 with the help of a LR finder. The lower bound is set to 0.0063. The L2 regularization technique is used with a decay factor of $10^{-5}$. Moreover, the batch size is adjusted to 32. We optimize the SNN architecture by first training the adapted network head of the incorporated ResNet-50 for 30 epochs with all other parameters being frozen. Then, the complete architecture is trained for another cycle, this time consisting of 50 epochs.

Since the batch size limits the task of constructing informative positive and negative pairs in the online mining strategy, the concept of cross-batch memory[57] is utilized to generate sufficient pairs across multiple mini-batches. This concept is based upon the observation that the embedding features generally tend to change slowly over time. This "slow drift" phenomenon allows the use of embeddings of previous iterations that would normally be considered out-dated and discarded. For our experiments, a memory size of 128 is chosen, meaning that the last 4 batches are considered for mining.

### Evaluation techniques

To evaluate the re-identification performance of our trained model, several metrics are computed. *R*-Precision represents the precision at *R*, where *R* denotes the number of relevant images for a given query image. In other words, if the top-*R* retrieved images show *r* relevant images, then *R*-Precision can be calculated from Equation (2). Note that this value is then averaged over all query samples. Precision@1 constitutes a special case and evaluates how many times the top-1 images in the retrieved lists are relevant.

$$R\text{-}Precision = \frac{r}{R} \tag{2}$$

To further consider the order of the relevant images within the retrieved list, the mean average precision at *R* (*mAP@R*) is computed according to Equation (3). The *mAP@R* denotes the mean of the average precision scores at *R* (*AP@R*) over all *Q* query images. The *AP@R* (see Equation (4)) is the average of the precision values over all *R* relevant samples, where *P@i* refers to the precision at rank *i* and *rel@i* is an indicator function which equals 1 if the sample is relevant at rank *i* and 0 if it is not relevant.

$$mAP@R = \frac{1}{Q}\sum_{i=1}^{Q} AP_i@R \tag{3} \qquad AP@R = \frac{1}{R}\sum_{i=1}^{R} P@i \times rel@i \tag{4}$$

### Data availability

The NIH ChestX-ray14 dataset used throughout the current study is available via Box at `https://nihcc.app.box.com/v/ChestXray-NIHCC`. The COVID-19 Image Data Collection is available on GitHub at `https://github.com/ieee8023/covid-chestxray-dataset`. The CheXpert dataset can be requested at `https://stanfordmlgroup.github.io/competitions/chexpert`.

### Code availability

The code used to train and evaluate both the patient verification and the patient re-identification models is available at `https://github.com/kaipackhaeuser/CXR-Patient-ReID`. Correspondence and requests for materials should be addressed to K.P.

### Author contributions

A.M., V.C., S.G., C.S., and K.P. conceived the main idea. K.P. and N.M. performed the experiments and the evaluation. K.P. wrote the main part of the manuscript. S.G. offered continuous support during the experiments and the writing process. S.G., C.S., V.C., and A.M. provided expertise through intense discussions. All authors reviewed the manuscript.

### Competing interests

The authors declare no competing interests.

## References

1. Maier, A., Steidl, S., Christlein, V. & Hornegger, J. *Medical Imaging Systems: An Introductory Guide*, vol. 11111 (Springer, 2018).

2. Raoof, S. *et al.* Interpretation of Plain Chest Roentgenogram. *Chest* **141**, 545–558 (2012).

3. Gündel, S. *et al.* Learning to recognize Abnormalities in Chest X-Rays with Location-Aware Dense Networks. In *Iberoamerican Congress on Pattern Recognition*, 757–765 (Springer, 2018).

4. World Health Organization (WHO). Coronavirus. https://www.who.int/health-topics/coronavirus (2020). [Online; accessed: 21.12.2020].

5. Wang, L., Lin, Z. Q. & Wong, A. COVID-Net: a tailored deep convolutional neural network design for detection of COVID-19 cases from chest X-ray images. *Sci. Reports* **10**, 1–12 (2020).

6. Lee, C. S., Nagy, P. G., Weaver, S. J. & Newman-Toker, D. E. Cognitive and System Factors Contributing to Diagnostic Errors in Radiology. *Am. J. Roentgenol.* **201**, 611–617 (2013).

7. LeCun, Y., Bengio, Y. & Hinton, G. Deep learning. *Nature* **521**, 436–444 (2015).

8. Gündel, S. *et al.* Multi-task Learning for Chest X-ray Abnormality Classification on Noisy Labels. *arXiv: 1905.06362* (2019).

9. Akselrod-Ballin, A. *et al.* A region based convolutional network for tumor detection and classification in breast mammography. In *Deep learning and data labeling for medical applications*, 197–205 (Springer, 2016).

10. Rajpurkar, P. *et al.* CheXNet: Radiologist-Level Pneumonia Detection on Chest X-Rays with Deep Learning. *arXiv: 1711.05225* (2017).

11. Roh, Y., Heo, G. & Whang, S. E. A Survey on Data Collection for Machine Learning: A Big Data - AI Integration Perspective. *IEEE Transactions on Knowl. Data Eng.* (2019).

12. Maier, A., Syben, C., Lasser, T. & Riess, C. A gentle introduction to deep learning in medical image processing. *Zeitschrift für Medizinische Physik* **29**, 86–101 (2019).

13. Oakden-Rayner, L. Exploring Large-scale Public Medical Image Datasets. *Acad. Radiol.* **27**, 106–112 (2020).

14. Irvin, J. *et al.* CheXpert: A Large Chest Radiograph Dataset with Uncertainty Labels and Expert Comparison. In *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, 590–597 (2019).

15. Gohagan, J. K., Prorok, P. C., Hayes, R. B. & Kramer, B.-S. The Prostate, Lung, Colorectal and Ovarian (PLCO) Cancer Screening Trial of the National Cancer Institute: History, organization, and status. In *Controlled Clinical Trials*, vol. 21, 251S–272S (2000).

16. Wang, X. *et al.* ChestX-ray8: Hospital-scale Chest X-ray Database and Benchmarks on Weakly-Supervised Classification and Localization of Common Thorax Diseases. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2097–2106 (2017).

17. Bandyopadhyay, D. *et al.* Covid-19 pandemic: cardiovascular complications and future implications. *Am. J. Cardiovasc. Drugs* 1–14 (2020).

18. Spinelli, A. & Pellino, G. Covid-19 pandemic: perspectives on an unfolding crisis. *The Br. journal surgery* (2020).

19. Cohen, J. P. *et al.* COVID-19 Image Data Collection: Prospective Predictions are the Future. *arXiv: 2006.11988* (2020).

20. Chung, A. Figure 1 COVID-19 Chest X-ray Dataset Initiative. https://github.com/agchung/Figure1-COVID-chestxray-dataset (2020).

21. Chung, A. ActualMed COVID-19 Chest X-ray Dataset Initiative. https://github.com/agchung/Actualmed-COVID-chestxray-dataset (2020).

22. Rahman, T., Chowdhury, M. & Khandakar, A. COVID-19 Radiography Database. https://www.kaggle.com/tawsifurrahman/covid19-radiography-database (2020).

23. Willemink, M. J. *et al.* Preparing Medical Imaging Data for Machine Learning. *Radiology* **295**, 4–15 (2020).

24. Centers for Disease Control and Prevention. Health Insurance Portability and Accountability Act of 1996 (HIPAA). https://www.cdc.gov/phlp/publications/topic/hipaa.html (2018). [Online; accessed: 23.12.2020].

25. European Union. Complete guide to GDPR compliance. https://gdpr.eu/ (2020). [Online; accessed: 23.12.2020].

26. O'Connor, M. Google axed release of vast x-ray dataset following NIH privacy concerns. https://www.healthimaging.com/topics/imaging-informatics/google-axed-release-x-ray-dataset-nih-concerns (2019). [Online; accessed: 17.12.2020].

27. Vincent, J. Google scrapped the publication of 100,000 chest x-rays due to last-minute privacy problems. https://www.theverge.com/2019/11/15/20966460/google-scrapped-publication-100000-chest-x-rays-nih-project-2017 (2019). [Online; accessed: 17.12.2020].

28. Noumeir, R., Lemay, A. & Lina, J.-M. Pseudonymization of Radiology Data for Research Purposes. *J. Digit. Imaging* **20**, 284–295 (2007).

29. Sweeney, L. k-anonymity: A model for protecting privacy. *Int. J. Uncertainty, Fuzziness Knowledge-Based Syst.* **10**, 557–570 (2002).

30. Gkoulalas-Divanis, A. & Loukides, G. *Medical Data Privacy Handbook* (Springer, 2015).

31. Zheng, M., Karanam, S., Wu, Z. & Radke, R. J. Re-Identification with Consistent Attentive Siamese Networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 5735–5744 (2019).

32. Selvaraju, R. R. *et al.* Grad-CAM: Visual Explanations from Deep Networks via Gradient-based Localization. In *Proceedings of the IEEE International Conference on Computer Vision*, 618–626 (2017).

33. Gündel, S. *et al.* Robust Classification from Noisy Labels: Integrating Additional Knowledge for Chest Radiography Abnormality Assessment. *Med. Image Analysis* 102087 (2021).

34. Nautsch, A. *et al.* Preserving privacy in speaker and speech characterisation. *Comput. Speech & Lang.* **58**, 441–480 (2019).

35. Tomashenko, N. *et al.* The VoicePrivacy 2022 Challenge Evaluation Plan. https://www.voiceprivacychallenge.org/docs/VoicePrivacy_2022_Eval_Plan_v1.0.pdf (2022).

36. Dwork, C. A Firm Foundation for Private Data Analysis. *Commun. ACM* **54**, 86–95 (2011).

37. Dwork, C., Roth, A. *et al.* The Algorithmic Foundations of Differential Privacy. *Foundations Trends Theor. Comput. Sci.* **9**, 211–407 (2014).

38. Kaissis, G. A., Makowski, M. R., Rückert, D. & Braren, R. F. Secure, privacy-preserving and federated machine learning in medical imaging. *Nat. Mach. Intell.* **2**, 305–311 (2020).

39. Sarwate, A. D. & Chaudhuri, K. Signal Processing and Machine Learning with Differential Privacy: Algorithms and challenges for continuous data. *IEEE signal processing magazine* **30**, 86–94 (2013).

40. Konečnỳ, J., McMahan, H. B., Ramage, D. & Richtárik, P. Federated Optimization: Distributed Machine Learning for On-Device Intelligence. *arXiv preprint arXiv:1610.02527* (2016).

41. Rieke, N. *et al.* The future of digital health with federated learning. *NPJ digital medicine* **3**, 1–7 (2020).

42. Kaissis, G. *et al.* End-to-end privacy preserving deep learning on multi-institutional medical imaging. *Nat. Mach. Intell.* **3**, 473–484 (2021).

43. Bromley, J. *et al.* Signature Verification using a "Siamese" Time Delay Neural Network. *Int. J. Pattern Recognit. Artif. Intell.* **7**, 669–688 (1993).

44. Taigman, Y., Yang, M., Ranzato, M. & Wolf, L. DeepFace: Closing the Gap to Human-Level Performance in Face Verification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 1701–1708 (2014).

45. Koch, G., Zemel, R. & Salakhutdinov, R. Siamese Neural Networks for One-shot Image Recognition. In *ICML Deep Learning Workshop* (2015).

46. Lake, B. M., Salakhutdinov, R. & Tenenbaum, J. B. Human-level concept learning through probabilistic program induction. *Science* **350**, 1332–1338 (2015).

47. LeCun, Y., Bottou, L., Bengio, Y. & Haffner, P. Gradient-Based Learning Applied to Document Recognition. *Proc. IEEE* **86**, 2278–2324 (1998).

48. National Institutes of Health (NIH). NIH Clinical Center provides one of the largest publicly available chest x-ray datasets to scientific community. https://www.nih.gov/news-events/news-releases/nih-clinical-center-provides-one-largest-publicly-available-chest-x-ray-datasets-scientific-community (2017). [Online; accessed: 05.01.2021].

49. Deng, J. *et al.* ImageNet: A Large-Scale Hierarchical Image Database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, 248–255 (2009).

50. LeCun, Y. A., Bottou, L., Orr, G. B. & Müller, K.-R. Efficient BackProp. In *Neural Networks: Tricks of the Trade*, 9–48 (Springer, 2012).

51. Goodfellow, I., Bengio, Y. & Courville, A. *Deep Learning* (MIT Press, 2016). URL: http://www.deeplearningbook.org.

52. Kingma, D. P. & Ba, J. Adam: A Method for Stochastic Optimization. *arXiv preprint arXiv:1412.6980* (2014).

53. Fawcett, T. An introduction to ROC analysis. *Pattern Recognit. Lett.* **27**, 861–874 (2006).

54. Hadsell, R., Chopra, S. & LeCun, Y. Dimensionality Reduction by Learning an Invariant Mapping. In *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06)*, vol. 2, 1735–1742 (IEEE, 2006).

55. Smith, L. N. & Topin, N. Super-Convergence: Very Fast Training of Neural Networks Using Large Learning Rates. In *Artificial Intelligence and Machine Learning for Multi-Domain Operations Applications*, vol. 11006, 1100612 (International Society for Optics and Photonics, 2019).

56. Smith, L. N. Cyclical Learning Rates for Training Neural Networks. In *2017 IEEE Winter Conference on Applications of Computer Vision (WACV)*, 464–472 (IEEE, 2017).

57. Wang, X., Zhang, H., Huang, W. & Scott, M. R. Cross-Batch Memory for Embedding Learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 6388–6397 (2020).

**Supplementary Table 1.** Additional patient verification results: Comparison of different data handling techniques (FTS and RNP), training set sizes $N_s$ and learning rates $\eta$. We present the AUC (together with the lower and upper bounds of the 95 % confidence intervals from 10,000 bootstrap runs), the accuracy, the specificity, the recall, the precision, and the F1-score.

| Data handling | $N_s$ | $\eta$ | AUC + 95 % CI | Accuracy ($\frac{TP+TN}{P+N}$) | Specificity ($\frac{TN}{N}$) | Recall ($\frac{TP}{P}$) | Precision ($\frac{TP}{TP+FP}$) | F1-score |
|---|---|---|---|---|---|---|---|---|
| FTS | 100,000 | $10^{-4}$ | $0.8509\ ^{0.8532}_{0.8485}$ | $0.7461\ \left(\frac{74,605}{100,000}\right)$ | $0.8414\ \left(\frac{42,071}{50,000}\right)$ | $0.6507\ \left(\frac{32,534}{50,000}\right)$ | $0.8040\ \left(\frac{32,534}{40,463}\right)$ | 0.7193 |
| | | $10^{-5}$ | $0.7429\ ^{0.7459}_{0.7399}$ | $0.6322\ \left(\frac{63,223}{100,000}\right)$ | $0.7924\ \left(\frac{39,619}{50,000}\right)$ | $0.4721\ \left(\frac{23,604}{50,000}\right)$ | $0.6945\ \left(\frac{23,604}{33,985}\right)$ | 0.5621 |
| | | $10^{-6}$ | $0.7257\ ^{0.7289}_{0.7225}$ | $0.6702\ \left(\frac{67,020}{100,000}\right)$ | $0.6719\ \left(\frac{33,594}{50,000}\right)$ | $0.6685\ \left(\frac{33,426}{50,000}\right)$ | $0.6708\ \left(\frac{33,426}{49,832}\right)$ | 0.6696 |
| | | $10^{-7}$ | $0.7440\ ^{0.7470}_{0.7409}$ | $0.6741\ \left(\frac{67,410}{100,000}\right)$ | $0.7002\ \left(\frac{35,008}{50,000}\right)$ | $0.6480\ \left(\frac{32,402}{50,000}\right)$ | $0.6837\ \left(\frac{32,402}{47,394}\right)$ | 0.6654 |
| | 200,000 | $10^{-4}$ | $0.8886\ ^{0.8906}_{0.8866}$ | $0.7920\ \left(\frac{79,203}{100,000}\right)$ | $0.8547\ \left(\frac{42,736}{50,000}\right)$ | $0.7293\ \left(\frac{36,467}{50,000}\right)$ | $0.8339\ \left(\frac{36,467}{43,731}\right)$ | 0.7781 |
| | | $10^{-5}$ | $0.8158\ ^{0.8185}_{0.8132}$ | $0.6758\ \left(\frac{67,582}{100,000}\right)$ | $0.8574\ \left(\frac{42,869}{50,000}\right)$ | $0.4943\ \left(\frac{24,713}{50,000}\right)$ | $0.7761\ \left(\frac{24,713}{31,844}\right)$ | 0.6039 |
| | | $10^{-6}$ | $0.7615\ ^{0.7645}_{0.7586}$ | $0.6755\ \left(\frac{67,552}{100,000}\right)$ | $0.7448\ \left(\frac{37,242}{50,000}\right)$ | $0.6062\ \left(\frac{30,310}{50,000}\right)$ | $0.7038\ \left(\frac{30,310}{43,068}\right)$ | 0.6514 |
| | | $10^{-7}$ | $0.7526\ ^{0.7556}_{0.7495}$ | $0.6819\ \left(\frac{68,194}{100,000}\right)$ | $0.6791\ \left(\frac{33,957}{50,000}\right)$ | $0.6847\ \left(\frac{34,237}{50,000}\right)$ | $0.6809\ \left(\frac{34,237}{50,280}\right)$ | 0.6828 |
| | 400,000 | $10^{-3}$ | $0.9541\ ^{0.9552}_{0.9529}$ | $0.8852\ \left(\frac{88,519}{100,000}\right)$ | $0.8655\ \left(\frac{43,275}{50,000}\right)$ | $0.9049\ \left(\frac{45,244}{50,000}\right)$ | $0.8706\ \left(\frac{45,244}{51,969}\right)$ | 0.8874 |
| | | $10^{-5}$ | $0.8518\ ^{0.8542}_{0.8494}$ | $0.7450\ \left(\frac{74,498}{100,000}\right)$ | $0.8410\ \left(\frac{42,050}{50,000}\right)$ | $0.6490\ \left(\frac{32,448}{50,000}\right)$ | $0.8032\ \left(\frac{32,448}{40,398}\right)$ | 0.7179 |
| | | $10^{-6}$ | $0.7417\ ^{0.7449}_{0.7386}$ | $0.6572\ \left(\frac{65,715}{100,000}\right)$ | $0.7152\ \left(\frac{35,761}{50,000}\right)$ | $0.5991\ \left(\frac{29,954}{50,000}\right)$ | $0.6778\ \left(\frac{29,954}{44,193}\right)$ | 0.6360 |
| | | $10^{-7}$ | $0.7306\ ^{0.7337}_{0.7275}$ | $0.6567\ \left(\frac{65,670}{100,000}\right)$ | $0.6801\ \left(\frac{34,006}{50,000}\right)$ | $0.6333\ \left(\frac{31,664}{50,000}\right)$ | $0.6644\ \left(\frac{31,664}{47,658}\right)$ | 0.6485 |
| RNP | 800,000 | $10^{-3}$ | $0.9826\ ^{0.9833}_{0.9820}$ | $0.9324\ \left(\frac{93,238}{100,000}\right)$ | $0.9393\ \left(\frac{46,966}{50,000}\right)$ | $0.9254\ \left(\frac{46,272}{50,000}\right)$ | $0.9385\ \left(\frac{46,272}{49,306}\right)$ | 0.9319 |
| | | $10^{-4}$ | $\mathbf{0.9940}\ ^{0.9944}_{0.9937}$ | $0.9555\ \left(\frac{95,545}{100,000}\right)$ | $0.9822\ \left(\frac{49,111}{50,000}\right)$ | $0.9287\ \left(\frac{46,434}{50,000}\right)$ | $0.9812\ \left(\frac{46,434}{47,323}\right)$ | 0.9542 |
| | | $10^{-5}$ | $0.9278\ ^{0.9294}_{0.9262}$ | $0.8339\ \left(\frac{83,392}{100,000}\right)$ | $0.8946\ \left(\frac{44,732}{50,000}\right)$ | $0.7732\ \left(\frac{38,660}{50,000}\right)$ | $0.8801\ \left(\frac{38,660}{43,928}\right)$ | 0.8232 |
| FTS | | $10^{-5}$ | $0.9200\ ^{0.9217}_{0.9182}$ | $0.8215\ \left(\frac{82,153}{100,000}\right)$ | $0.8888\ \left(\frac{44,440}{50,000}\right)$ | $0.7543\ \left(\frac{37,713}{50,000}\right)$ | $0.8715\ \left(\frac{37,713}{43,273}\right)$ | 0.8087 |
| | | $10^{-6}$ | $0.8669\ ^{0.8692}_{0.8646}$ | $0.7752\ \left(\frac{77,519}{100,000}\right)$ | $0.8165\ \left(\frac{40,823}{50,000}\right)$ | $0.7339\ \left(\frac{36,696}{50,000}\right)$ | $0.7999\ \left(\frac{36,696}{45,873}\right)$ | 0.7655 |
| | | $10^{-7}$ | $0.8126\ ^{0.8152}_{0.8099}$ | $0.7247\ \left(\frac{72,468}{100,000}\right)$ | $0.7502\ \left(\frac{37,509}{50,000}\right)$ | $0.6992\ \left(\frac{34,959}{50,000}\right)$ | $0.7368\ \left(\frac{34,959}{47,450}\right)$ | 0.7175 |