

# Image Projective Transformation Rectification with Synthetic Data for Smartphone-captured Chest X-ray Photos Classification

Chak Fong Chong<sup>a</sup>, Yapeng Wang<sup>a</sup>, Benjamin Ng<sup>a</sup>, Wuman Luo<sup>a</sup>, Xu Yang<sup>a,\*</sup>

<sup>a</sup>*Macao Polytechnic University, Macao SAR, China*

---

## Abstract

Classification on smartphone-captured chest X-ray (CXR) photos to detect pathologies is challenging due to the projective transformation caused by the non-ideal camera position. Recently, various rectification methods have been proposed for different photo rectification tasks such as document photos, license plate photos, etc. Unfortunately, we found that none of them is suitable for CXR photos, due to their specific transformation type, image appearance, annotation type, etc. In this paper, we propose an innovative deep learning-based **Projective Transformation Rectification Network** (PTRN) to automatically rectify CXR photos by predicting the projective transformation matrix. To the best of our knowledge, it is the first work to predict the projective transformation matrix as the learning goal for photo rectification. Additionally, to avoid the expensive collection of natural data, synthetic CXR photos are generated under the consideration of natural perturbations, extra screens, etc. We evaluate the proposed Method in the CheXphoto smartphone-captured CXR photos classification competition hosted by the Stanford University Machine Learning Group, our approach won **first place** with a huge performance improvement (ours 0.850, second-best 0.762, in AUC). A deeper study demonstrates that the use of PTRN successfully achieves the classification performance on the spatially transformed CXR

---

\*Corresponding author

*Email addresses:* chakfong.chong@mpu.edu.mo (Chak Fong Chong), yapengwang@mpu.edu.mo (Yapeng Wang), bng@mpu.edu.mo (Benjamin Ng), wumanluo@mpu.edu.mo (Wuman Luo), xuyang@mpu.edu.mo (Xu Yang)

photos to the same level as on the high-quality digital CXR images, indicating PTRN can eliminate all negative impacts of projective transformation on the CXR photos.

*Keywords:* Chest X-ray, Image rectification, Projective transformation, Camera perspective, Deep learning, Medical image analysis

---

## 1. Introduction

Chest X-ray (CXR), also named chest radiograph, is one of the most ubiquitous medical imaging techniques for chest disease diagnosis. However, due to the lack of resources, many CXRs cannot be formally reviewed by radiologists on time [1, 2]. For instance, in only one hospital in the United Kingdom, over 23,000 CXRs were not formally reviewed on time in the past 12 months [3]. Such serious situations provide a strong motivation to develop computer-aid diagnosis (CAD) systems for CXR automatic interpretation.

Convolutional neural networks (CNN) have been used to perform CXR multi-label classification tasks to detect pathologies. These CNN classification models are usually trained on datasets of digital CXRs [4, 5, 6] produced by modern computer-based digital X-ray systems [7, 8, 9, 10, 11]. The classification performances are promising and even approach radiologist-level [7, 9].

Unfortunately, in many developing countries and regions, traditional X-ray film systems are still widely used [12, 13]. CXR films are hard copies, so cannot be interpreted directly by computer-based models. To address this problem, a quick and convenient solution is to use smartphones to capture photos of CXR films (*CXR film photos*), then a CXR classification model interprets the captured CXR photos [10, 14, 15, 16]. Benefitting from the low prices and popularity of smartphones, it is a cheap and effective solution for countries and regions which have very limited medical resources. Moreover, the use of smartphones can facilitate both real-time and store-and-forward medical consultation to provide remote medical care and remote diagnosis [17, 18]. Especially, the sudden outbreak of COVID-19 in 2020 has promoted the practical applications of telemedicine and zero-contact diagnosis. Furthermore, it can help patients who have privacy concerns that not willing to go to hospitals for radiology reports.

Additionally, to demonstrate the generalization of our proposed method on CXR photos captured in various unconstrained scenarios, we also include

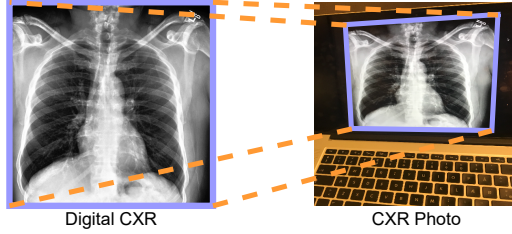


Figure 1: The CXR in a photo are warped by projective transformation due to non-ideal camera position.

smartphone-captured photos of digital CXR images displayed on the monitor screens (*CXR monitor photos*) to evaluate our proposed method in this paper.

The classification on smartphone-captured CXR photos is a very challenging task because CXR photos have very different appearances from digital CXRs, which leads to a significant performance drop. In a CXR photo, (1) the CXR is warped by projective transformation due to non-ideal camera position [19], as shown in Figure 1; (2) Natural perturbations such as environmental illuminations, camera out-focus, image noises, etc., exist in photos. Moreover, CXR photos are captured in various unconstrained scenarios, thus they have very different projective transformations and natural perturbations. Therefore, the classification performance drops significantly when a CNN model is trained on high-quality digital CXR images and then evaluated on CXR photos [10, 16, 15]. In the paper of [16], Chong et al. demonstrated that the projective transformation of CXR photos is the major reason for the significant classification performance drop.

Therefore, one obvious solution to improve classification performance on CXR photos is to automatically rectify the projective transformation of CXR photos. Chong et al. [16] proposed a method named generative adversarial network-based spatial transformation adversarial method (GAN-STAM) to automatically rectify CXR photos. However, the results showed that GAN-STAM cannot properly rectify CXR photos, probably due to the following reasons: the lack of training data, inappropriate rectification mechanism (affine transformation), instability of GAN training, etc.

Recently, various rectification methods have been proposed for different photo rectification tasks such as document photos [20, 21, 22, 23, 24, 25, 26], license plate photos [27, 28, 29, 30, 31], etc. Unfortunately, we found that

none of them is suitable for smartphone-captured CXR photos, due to their specific transformation type, image appearance, annotation type, etc. Normally, an existing method is particularly designed for one specific rectification task only. For example, in document photos unwarping, the documents in photos are folded and curved, which are not suitable for smartphone-captured CXR photos.

In this paper, we propose an innovative method to rectify the projective transformation of CXR photos caused by non-ideal camera position, named **Projective Transformation Rectification Network (PTRN)**. To the best of our knowledge, it is the first work to predict the projective transformation matrix as the learning goal for photo rectification. Additionally, synthetic CXR photos with transformation matrices as the ground truth annotations are generated for training PTRN. Moreover, the designs of both PTRN and the synthetic data generation framework are general to any camera-captured photos such as CCTV-captured license plate photos.

The primary contributions of this work are summarized:

1. We propose an innovative method named PTRN to rectify the projective transformation caused by non-ideal camera positions. The transformation matrix is set as the learning goal of PTRN for photo rectification. Additionally, the Intersection over Union (IoU), normally used to quantify the percent overlap of image segmentation, is first proposed as the performance metric to quantitatively evaluate the rectification performance, to the best of our knowledge.
2. We also propose the innovative synthetic CXR photos generation framework, which is designed under the consideration of the appearances of natural CXR photos including screen, background, and natural perturbations, to avoid the collection and annotation of expensive natural data. The framework produces projective transformation matrices as the ground truth annotations, instead of bounding boxes which are common in synthetic data generation frameworks for other object detection tasks. The proposed framework is fast yet efficient, as it only uses general image processing methods.
3. Specifically, we design a CXR photos classification pipeline in three steps: (1) PTRN predicts the projective transformation matrix of a CXR photo; (2) The photo is rectified using the predicted matrix; (3) A classifier trained on high-quality digital CXRs evaluates the rectified photo to produce the final result.



4. We train PTRN on synthetic data generated using the CheXpert digital CXR dataset [5] and the MS-COCO dataset [32]. We also train a classifier on the CheXpert dataset to build up the CXR photos classification pipeline. Then, the pipeline performance is evaluated on the CheXphoto smartphone-captured CXR photos classification dataset [14]. The pipeline achieves AUC 0.880/0.802 on the CXR monitor photos / CXR film photos, respectively, which is much higher than the performance without using PTRN to rectify photos (AUC 0.710/0.599).
5. The proposed PTRN with the classification pipeline is also evaluated on the CheXphoto competition hosted by Stanford and Vinbrain. Our pipeline achieves **first place** in the competition with AUC 0.850 which is much higher than the second-best approach (AUC 0.762) on the competition leaderboard. A deeper study demonstrates that the use of PTRN successfully achieves the classification performance on the spatially transformed CXR photos to the same level as on the high-quality digital CXR images, indicating PTRN can eliminate all negative impacts of projective transformation on the CXR photos.

The rest of the paper is organized as follows. In Section 2, work related to smartphone-captured CXR photos classification, image rectification, and deep learning with synthetic data is introduced. Section 3 introduces PTRN and the synthetic data framework. Section 4 shows the experiment results and discussion. Section 5 is the conclusion, limitations, and future work.

## 2. Related Work

### 2.1. Smartphone-captured CXR Photos Classification

Classification performance drop when current CXR classification models are evaluated on smartphone-captured CXR photos [10, 16, 7]. It is because CXR photos have very different appearances from digital CXRs that the models are trained on, as introduced in Section 1.

Several methods have been proposed to improve classification performance. One first attempt is Le et al. [33] proposed a method to crop the CXRs in photos using YOLO-v3 [34], but the classification performance is unacceptably low (AUC 0.684). Kuo et al. [10] developed a recalibration method that uses data augmentation methods to add additional noises into digital CXRs for training. The method achieves AUC 0.835 on CXR photos. However, this method does not tackle the projective transformation problem.

Chong et al. [16] found that once the projective transformation of CXR photos is perfectly rectified before classification, the classification performance significantly improves from AUC 0.801 to AUC 0.887. Therefore, they proposed GAN-STAM to automatically rectify the projective transformation. However, the classification performance attains AUC 0.865 only, which is still far away from the performance on perfectly rectified photos. It indicates GAN-STAM cannot properly rectify the projective transformation, probably due to these reasons: (1) GAN-STAM uses an affine transformation to rectify photos, but the CXRs in photos are warped by the projective transformation that is more complex than the affine transformation [35]. (2) the lack of data and annotations. (3) the training of GAN is unstable [36, 37]; and (4) no quantitative evaluation metric to measure the rectification performance.

## 2.2. Camera-captured Photos Rectification

**Traditional model-based methods.** Model-based methods have been proposed for document photos [38, 39, 40], QR code photos [41, 42], etc. These methods rely on image appearances. For document photos, horizontal sentences printed in the document are utilized by using image processing techniques like Radon transformation [43, 40] to calculate the distortion of the document such as the vanishing points and the level of curving. For a QR code photo, the QR code features (finder patterns and alignment patterns) are utilized to locate the QR codes for rectification. Overall, a traditional method is usually designed for rectifying a specific target. Additionally, these methods rely on low-level image patterns extracted by traditional digital image processing methods, so that unsuitable for photos captured in various unconstrained scenarios.

**Deep learning-based methods.** Deep learning has been widely used for such rectification tasks since deep learning models can learn from a large amount of data. Therefore, it is more robust to camera-captured natural photos and reduces the manual design of algorithms. Various rectification methods have been proposed for different photo rectification tasks in the literature such as document photos, license plate photos, etc. Unfortunately, we found that none of them is suitable for CXR photos. It is because a task has its distinctive properties such as the transformation type, image appearance, annotation type, etc. Hence, a rectification method is particularly designed for a specific task and not compatible with other tasks. For document photo unwarping, the documents are folded and curved [20, 21, 22, 23, 24, 25, 26], which is unsuitable for CXR photos. In license plate recognition, several

methods train models to crop/rectify the detected license plates in photos using affine transformation [27, 28, 29], but CXR photos are projective transformation which is more complex than affine transformation [35]. Some methods [30, 31] utilized the spatial transformer network (STN) [44] (a learnable neural network module that can perform spatial transformations to feature maps or images) to reduce the problem of transformation. In these methods, A STN-like module is inserted into the license plate recognition network to reduce the transformation problem for better recognition accuracy. The whole recognition network is trained in an end-to-end manner. Therefore, the STN-like module aims to select the region in the whole image that the recognition network is interested in (attention), instead of precisely rectifying the license plate. Besides, some methods were proposed for other rectification tasks like barcode photos [45] and text photos [46, 47]. The method for barcodes relies on the vertical lines of barcodes, while the transformation types of texts include line curvature which is very different from photos of CXR.

Overall, to the best of our knowledge, existing rectification methods are not compatible with the projective transformation of photos like CXR photos. Therefore, we propose PTRN which is the first method that predicts the projective transformation matrices to rectify the projective transformation of CXR photos caused by the non-ideal camera position.

### *2.3. Synthetic Data for Photo Rectification and Object Detection*

Synthetic data is one popular method to address the problem of the lack of data for training deep learning models, as it avoids expensive natural data collection. Synthetic data has also been used for object detection [48, 49, 50] and camera-captured photo rectification tasks [23, 24, 28].

A common strategy for generating a synthetic image is a composition of a foreground image (i.e., a transformed image) and a background image, plus data augmentation to simulate natural perturbations [48, 49, 50, 23, 24, 28]. However, the detailed generation steps are quite different across tasks, since a task has its distinctive properties such as natural appearance and annotation type. In text detection [48], a text usually appears in well-defined regions like a sign or a flat wall. Therefore, the generation step includes additional geometric estimation and segmentation to ensure the texts are properly placed in the correct positions. The annotation type of object detection is bounding box, while the annotation type of document photo unwarping is 3D coordinate maps [24, 23, 25]. Hence, the synthetic data generation steps are

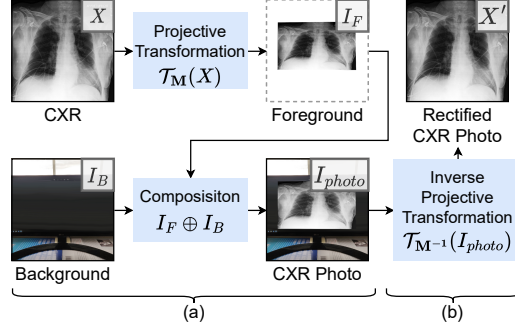


Figure 2: Formulation of projective transformation rectification for CXR photos.

different.

Our proposed framework for synthetic data generation has a couple of novelties in comparison to previous synthetic data frameworks. We propose the first CXR photos synthetic framework. The design of the framework is under consideration of the natural appearances of CXR photos. Furthermore, our framework is the first that uses the projective transformation matrix as the ground truth annotation to represent the transformation caused by the non-ideal camera position. Unlike object detection tasks that use bounding boxes that only locate the object, the projective transformation matrix can be used to precisely rectify the photo. Which is innovative and effective.

### 3. Methods

We propose Projective Transformation Rectification Network (PTRN) for rectifying the projective transformation of CXR photos caused by non-ideal camera position. PTRN is trained in an end-to-end manner on synthetic CXR photos training samples. The designs of both PTRN and the synthetic data generation framework are general such that PTRN can also be used for rectifying other camera-captured photos of images (e.g., CCTV-captured license plate photos), detailed in follows.

#### 3.1. Problem Formulation

This section formulates the problem of projective transformation rectification for CXR photos. The designs of PTRN and the synthetic data framework are based on the formulation.

A CXR photo  $I_{photo}$  can be considered as a composition of a transformed CXR  $I_F$  and a background image  $I_B$ , as in Figure 2 (a). We borrow a

notation  $\oplus$  from [51] to represent the composition of two images.  $I_{photo}$  can be formulated as:

$$I_{photo} = I_F \oplus I_B = \mathcal{T}_{\mathbf{M}}(X) \oplus I_B \quad (1)$$

where  $I_F$  is a CXR  $X$  warped by projective transformation with a 3-by-3 transform matrix  $\mathbf{M}$ :  $I_F = \mathcal{T}_{\mathbf{M}}(X)$ . A matrix  $\mathbf{M}$  has 8 parameters  $\theta_i$  that

represent a projective transformation:  $\mathbf{M} = \begin{bmatrix} \theta_1 & \theta_2 & \theta_3 \\ \theta_4 & \theta_5 & \theta_6 \\ \theta_7 & \theta_8 & 1 \end{bmatrix}$ . In this work,

all projective transformations are applied in the homogeneous coordinates and the coordinates of the 4 vertices of each image are assigned to be  $P_0 =$

$$\left\{ \begin{pmatrix} -1 \\ -1 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 \\ -1 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}, \begin{pmatrix} -1 \\ 1 \\ 1 \end{pmatrix} \right\}.$$

A CXR photo  $I_{photo}$  can be rectified using inverse projective transformation if  $\mathbf{M}$  is known (Figure 2 (b)):

$$X' = \mathcal{T}_{\mathbf{M}^{-1}}(I_{photo}) \quad (2)$$

where  $X'$  denotes the rectified CXR photos.

The design of PTRN is based on Equation 2: A CXR photo  $I_{photo}$  can be rectified by predicting  $\mathbf{M}$ . Therefore, the input of PTRN is a CXR photo  $I_{photo}$  and the output is the predicted matrix  $\hat{\mathbf{M}}$ . Based on this design, a training sample for PTRN should consist of a CXR photo  $I_{photo}$  and the matrix  $\mathbf{M}$ . To synthetically generate training samples, a generation framework is designed based on Equation 1: A training sample  $(I_{photo}, \mathbf{M})$  can be synthesized using a CXR  $X$ , a background image  $I_B$  and  $\mathbf{M}$ . The designs of PTRN and the synthetic data generation framework are detailed in the following sections. Besides, we note this formulation is general to other camera-captured photos of images, not CXRs only, since the CXR appearance is not considered in the formulation. E.g., the formulation is also suitable for CCTV-captured license plate photos.

### 3.2. Projective Transformation Rectification Network

PTRN is a deep neural network that receives a CXR photo  $I_{photo}$  and predicts the matrix  $\mathbf{M}$  for rectification. PTRN is particularly designed not to rely on any appearances of CXRs. Hence, PTRN is also suitable for rectifying projective transformation of camera-captured photos of different image types.

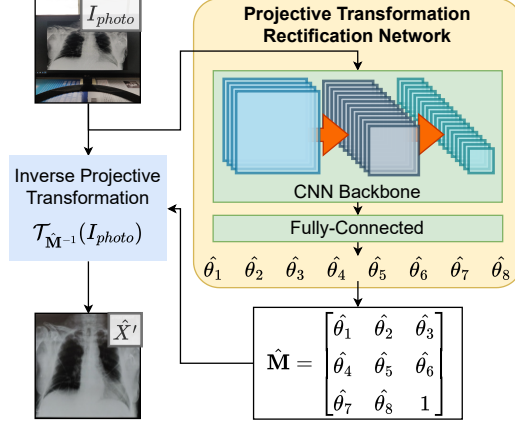


Figure 3: The architecture of PTRN and the steps of rectifying a CXR photo using PTRN.

### 3.2.1. Architecture

As shown in Figure 3 (right). The architecture of PTRN consists of two components: (1) a CNN backbone (e.g., ResNet-50 [52], DenseNet-121 [53]), followed by (2) a fully-connected layer to regress the 8 parameters  $\theta_i$  of a predicted matrix  $\hat{\mathbf{M}}$ .

Rectifying a CXR photo  $I_{photo}$  using PTRN has two steps, as shown in Figure 3: (1) PTRN predicts  $\hat{\mathbf{M}}$ ; (2) Apply the inverse projective transformation with  $\hat{\mathbf{M}}$  (Equation 2) to the CXR photo  $I_{photo}$ .

### 3.2.2. Training

PTRN is trained on synthetic training samples  $\{(I_{photo}^{(i)}, \mathbf{M}^{(i)})\}$  in an end-to-end manner. As PTRN predicts the parameters by regression, mean squared error (MSE) loss is used to calculate the loss between the model prediction and the ground truth. MSE loss is defined as  $L = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$  where  $y$  denotes the 8 parameters of the ground truth  $\mathbf{M}$ :  $y = [\theta_1 \dots \theta_8]$  and  $\hat{y}$  denotes the predicted 8 parameters  $\hat{y} = [\hat{\theta}_1 \dots \hat{\theta}_8]$ .

### 3.2.3. CXR Region Prediction

A matrix  $\mathbf{M}$  can be converted to a quadrilateral that represents a CXR region. The 4 vertices  $\begin{pmatrix} x_i \\ y_i \\ 1 \end{pmatrix}$  of the quadrilateral represented by  $\mathbf{M}$  are:

$$\begin{pmatrix} w_i x_i \\ w_i y_i \\ w_i \end{pmatrix} = \mathbf{M} \begin{pmatrix} u_i \\ v_i \\ 1 \end{pmatrix} \quad (3)$$

where  $\begin{pmatrix} u_i \\ v_i \\ 1 \end{pmatrix} \in P_0$ .

The quadrilateral can also be converted back to  $\mathbf{M}$  using Equation 3. This relation is used in the evaluation of the rectification performance of PTRN. For instance, a predicted matrix  $\hat{\mathbf{M}}$  can be converted to a predicted CXR region.

### 3.2.4. Evaluation and Rectification Performance Metric

PTRN is evaluated on natural test samples  $\left\{ \left( I_{\text{photo}}^{(i)}, P^{(i)} \right) \right\}$ . Each sample consists of a natural CXR photo  $I_{\text{photo}}$  and the 4 vertices  $\begin{pmatrix} x_i \\ y_i \\ 1 \end{pmatrix} \in P$  that are manually annotated to represent the ground-truth CXR region. The ground-truth label is the 4 vertices instead of the matrix  $\mathbf{M}$  since marking the 4 vertices is the simplest way to annotate the ground truth in practice.

We also propose using IoU between the ground truth CXR region and the predicted CXR region for each sample to measure the rectification performance, as the predicted region should be as overlapped as the ground truth region. The predicted region is calculated from  $\hat{\mathbf{M}}$  using Equation 3 and the ground truth region is given in the test sample. The IoU calculation of a sample is demonstrated in Figure 4.

### 3.3. Training Data Synthesis

PTRN requires a large amount of training data. Unfortunately, there are no suitable datasets. The collection of natural data samples is also expensive. Therefore, we propose a framework for the generation of synthetic training samples. To enable the generality of PTRN, the proposed framework is

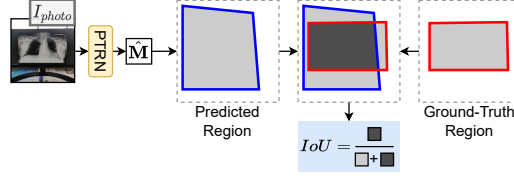


Figure 4: The rectification performance of PTRN is evaluated by the IoU between predicted regions and ground truth regions.

designed not to rely on the appearances of CXRs. Hence, the framework can also generate synthetic camera-captured photos of different image types with very few modifications. So that PTRN can be trained to rectify camera-captured photos of different images.

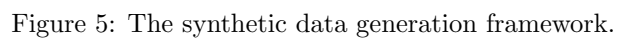
The generation of one synthetic training sample consists of 4 steps, as shown in Figure 5: (1) screen synthesis; (2) CXR projective transformation; (3) adding a background image; and (4) adding natural perturbations. The details of the 4 steps are described in the following subsections. A random CXR  $X$  and a random image  $R$  are required for generating a sample.

Note that the design of the framework does not exactly follow Equation 1, as various natural appearances of CXR photos are also considered such as screen and natural perturbations. Randomness is involved in the framework to increase the diversity of the synthesized CXR photos so that hopefully the trained PTRN can be generalized to photos captured in various scenarios. An ablation study is conducted in Section 4.6 to verify the framework. The framework only uses general image processing methods and general augmentation methods. Through experiment, the generation speed attains 102.89 samples/second in a PC with Intel i9-10900K CPU.

### 3.3.1. Screen Synthesis

In a natural CXR photo, the transformed CXR is often surrounded by dark padding, which is usually a monitor screen, as shown in Figure 6. As the dark padding looks like a part of the CXR. It may interfere with the prediction of PTRN. This step is to synthesize the screen to simulate such a situation to improve the robustness of PTRN in this situation. Considering that not all natural CXR photos have such a situation, the framework is designed to sometimes synthetic the screen. In experiments, the probability is set to be 0.3. If a screen is synthesized, the output of this step is a CXR with a synthetic screen  $X_{SC}$ , otherwise, the output is identical to the input





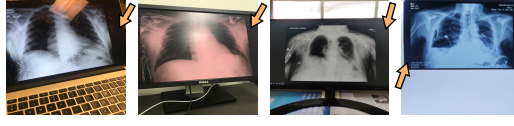


Figure 6: The CXRs in photos are usually surrounded with a dark area.

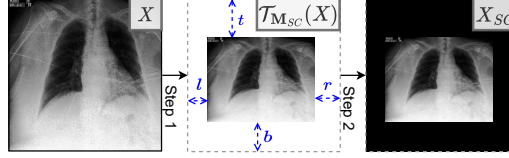


Figure 7: The two sub-steps of screen synthesis.

$X$ . The process of screen synthesis has two sub-steps, as shown in Figure 7.

**Sub-step 1:**  $X$  is scaled-down and translated, determined by the randomly generated padding sizes  $t, b, l, r$  (see Figure 7). In the experiments, we set:  $t, b, l, r \sim U_{[0,0.6]}$ . Subsequently, the operation of scaling down and translation is implemented by projective transformation. The transform matrix  $\mathbf{M}_{SC}$  is calculated by:

$$\mathbf{M}_{SC} = \begin{bmatrix} 1 - (l + r)/2 & 0 & (l - r)/2 \\ 0 & 1 - (t + b)/2 & (t - b)/2 \\ 0 & 0 & 1 \end{bmatrix} \quad (4)$$

. Finally,  $X$  is scaled-down and translated:  $\mathcal{T}_{\mathbf{M}_{SC}}(X)$ .

**Sub-step 2:**  $\mathcal{T}_{\mathbf{M}_{SC}}(X)$  and a plain dark color image  $I_{color}$  are composited to get the output  $X_{SC}$ :

$$X_{SC} = \mathcal{T}_{\mathbf{M}_{SC}}(X) \oplus I_{color} \quad (5)$$

. In the experiments, the RGB values of the  $I_{color}$  is randomly generated by  $R, G, B \sim U\{0, 19\}$ .

### 3.3.2. CXR Projective Transformation

This step is to simulate the projective transformation caused by the non-ideal camera position. Firstly, a projective transformation matrix  $\mathbf{M}$  is randomly generated. Then, projective transformation with  $\mathbf{M}$  is applied to the output of step 1 to get a transformed CXR  $I_F$ . The 8 parameters  $\theta_i$  of  $\mathbf{M}$  must be generated under some regulations since an inappropriate set of

| Parameters       | Generation   |
|------------------|--|
| Scaling          | $C_x, C_y \sim U_{[0.2, 0.8]},  C_x - C_y  \leq 0.2$ |
| Shearing         | $S_x, S_y \sim U_{[-0.1, 0.1]}$                      |
| Rotation         | $\alpha \sim U_{[-\pi, \pi]}$                        |
| Perspective Warp | $F_x, F_y \sim N(\mu = 0, \sigma^2 = 0.1^2)$         |
| Translation      | $T_x, T_y \sim N(\mu = 0, \sigma^2 = 0.25^2)$        |

Table 1: Setup in experiment for generating the parameters of the five individual actions. The setup of  $|C_x - C_y| \leq 0.2$  is to avoid the transformed CXRs be too narrow.

parameters results in a nonsensical transformed CXR, e.g., the shape is non-convex [54, 35]. To avoid such unexpected situations, we first equivalently represent a projective transformation by a sequential series of five individual actions: scaling, shearing, rotation, perspective warp, and translation. Then, the parameters of these actions are generated within an appropriate range. Specifically, when generating a matrix  $\mathbf{M}$ , the first step is to generate five 3-by-3 matrices that represent the five individual actions:

$$\mathbf{M}_C = \begin{bmatrix} C_x & 0 & 0 \\ 0 & C_y & 0 \\ 0 & 0 & 1 \end{bmatrix}, \mathbf{M}_S = \begin{bmatrix} 1 & S_x & 0 \\ S_y & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}, \mathbf{M}_R = \begin{bmatrix} \cos \alpha & \sin \alpha & 0 \\ -\sin \alpha & \cos \alpha & 0 \\ 0 & 0 & 1 \end{bmatrix}, \quad (6)$$

$$\mathbf{M}_P = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ P_x & P_y & 1 \end{bmatrix}, \mathbf{M}_T = \begin{bmatrix} 1 & 0 & T_x \\ 0 & 1 & T_y \\ 0 & 0 & 1 \end{bmatrix}, \quad (7)$$

where  $\mathbf{M}_C, \mathbf{M}_S, \mathbf{M}_R, \mathbf{M}_P, \mathbf{M}_T$  denote the transform matrix of scaling, shearing, rotation, perspective warp, and translation, respectively. The parameters in the five metrics are randomly generated in appropriate ranges. The setup in the experiment is shown in Table 1. Then, the generated  $\mathbf{M}$  is calculated by:

$$\mathbf{M} = \mathbf{M}_T \mathbf{M}_P \mathbf{M}_R \mathbf{M}_S \mathbf{M}_C \quad (8)$$

The generated  $\mathbf{M}$  is then used to transform the output of step 1 to calculate the transformed CXR  $I_F$ . The output of step 1 has two types:  $X$  and  $X_{SC}$ . In the case of  $X$ ,  $I_F$  is calculated by

$$I_F = \mathcal{T}_{\mathbf{M}}(X) \quad (9)$$

; In the case of  $X_{SC}$ , since the CXR region in  $X_{SC}$  has been transformed by  $\mathbf{M}_{SC}$ , the transformation by  $\mathbf{M}_{SC}$  must be inversed to ensure that  $\mathbf{M}$

| Perturbations     | Data Augmentation Methods              |
|-------------------|--|
| Illuminations     | Adding pixel values                    |
|                   | Multiplying pixel values               |
|                   | Adding hue and saturation value in HSV |
|                   | Color Enhancement                      |
|                   | Brightness Enhancement                 |
|                   | Sharpness enhancement                  |
| Out-focus         | Average blur                           |
| Image noises      | Adding Gaussian noise                  |
| Image compression | JPEG compression                       |

Table 2: Data augmentation methods for adding natural perturbations.

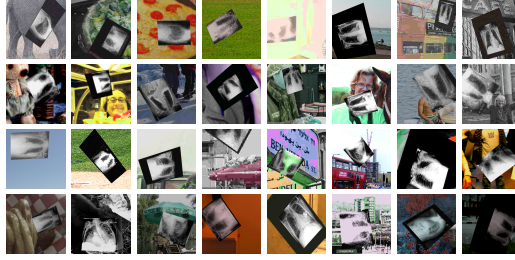


Figure 8: Synthetically generated CXR photos.

can correctly represent the CXR region in  $I_F$ . The transformed CXR  $I_F$  is calculated by

$$I_F = \mathcal{T}_{\mathbf{M}\mathbf{M}_{SC}^{-1}}(X_{SC}) \quad (10)$$

.

### 3.3.3. Adding Background Image

In practical situations, the background of CXR photos varies. We consider it as a random image  $R$ . This step is to composite the output  $I_F$  of step 2 and a background image. The output  $I_{FB}$  is calculated by  $I_{FB} = I_F \oplus R$ .

### 3.3.4. Adding Natural Perturbations

This step is to simulate the perturbations of natural photos. In the experiments, various data augmentation methods are used to simulate the perturbations, as listed in Table 2.

In the end,  $I_{photo}$  and matrix  $\mathbf{M}$  are composited to be a synthetic training sample  $(I_{photo}, \mathbf{M})$ . Figure 8 shows some samples of generated  $I_{photo}$ .

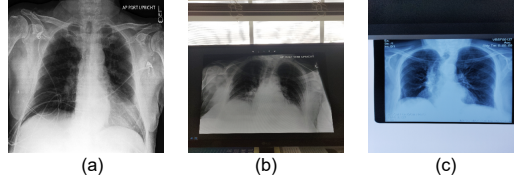


Figure 9: (a) CheXpert validation set. (b) CheXphoto-Monitor validation set. (c) CheXphoto-Film validation set.

## 4. Experimental Results and Discussion

All experiments are conducted on a desktop computer with hardware Intel i9-10900KF CPU, 128GB RAM, and RTX 3090 24G GPU. The deep learning platform is TensorFlow 2 [55] on Python 3, installed in a Linux Mint 20.1 OS. The data augmentation methods are implemented using imgaug library [56].

### 4.1. Datasets

#### 4.1.1. CheXpert

CheXpert [5] is a public digital CXR dataset for competition. It consists of 224,316 digital CXRs from 65,240 patients. Each digital CXR is labeled into at least one of 14 pathologies: No Finding, Enlarged Cardiomeastinum, Cardiomegaly, Lung Opacity, Lung Lesion, Edema, Consolidation, Pneumonia, Atelectasis, Pneumothorax, Pleural Effusion, Pleural Other, Fracture, and Support Devices. The validation set consists of 234 digital CXRs from 200 studies (Figure 9(a)).

#### 4.1.2. CheXphoto

CheXphoto [14] is a public smartphone-captured CXR photos dataset for competition. It consists of 10,507 CXR photos from 3,000 patients, which are sampled from digital CXRs in CheXpert. The labels are also inherited from CheXpert. The validation set consists of synthesis photos and natural photos. Only natural photos are used for evaluation in this paper. The validation set has two types of natural photos: (1) A total of 234 CXR monitor photos (Figure 9(b)). They are the photo version of the CheXpert validation set. Each photo is produced by using a smartphone camera to capture digital CXRs displayed on a monitor. (2) A total of 250 CXR film photos (Figure 9(c)). Each photo is produced by using a smartphone camera to capture

CXR film. In Section 4.4, PTRN is used for the CheXphoto competition. In the competition, the private test CXR photos were split into two private test sets by type (monitor/film). Therefore, we split the CheXphoto validation set into two sets by the type (monitor/film) for experiments, referred to as the CheXphoto-Monitor validation set and the CheXphoto-Film validation set, respectively. Note that the CheXphoto-Monitor validation set is the photo version of the CheXpert validation set (digital CXR), we use these two sets to demonstrate the effectiveness of PTRN in Section 4.5.

#### 4.2. Implementation Details

For PTRN, an ImageNet-pretrained DenseNet-201 [53] is used as the CNN backbone. Adam optimizer [57] with a learning rate of  $1 \times 10^{-5}$  and parameters ( $\beta_1 = 0.9, \beta_2 = 0.999, \epsilon = 1 \times 10^{-7}$ ) is used to update the weights. The batch size is set to 32. The input image is a  $224 \times 224$  pixel RGB image. Pixel values are linearly rescaled from  $[0, 255]$  to  $[0, 1]$ . Among these hyperparameters, only the learning rate was simply tuned while others remain default. The trained PTRN has sufficient rectification performance already (see Section 4.6), which demonstrates the easy training of PTRN.

For synthetic training data generation, we use the 224,316 digital CXRs from the CheXpert training set as the source of the CXR images, and the 41K images from the Microsoft COCO 2017 test set [32] as the source of the random images. When generating a sample, a CXR image and a random image are randomly picked from the sources. Benefitting from the fast generation speed of synthetic data, we dynamically generate the synthetic training data during training PTRN. Therefore, each training sample is used only once to avoid overfitting.

The performance of PTRN is validated on the CheXphoto-Monitor validation set per 100 weights updates. The checkpoint with the highest IoU is picked and furtherly tested on the CheXphoto-Film validation set. The evaluation of PTRN requires the 4 vertices of the ground truth region for each sample. Since the CheXphoto dataset did not provide these annotations, we manually annotated the 4 vertices of the CXR photos.

#### 4.3. CXR Photos Classification Pipeline

The pipeline to perform CXR photos classification (Figure 10) consists of three steps: (1) PTRN predicts the projective transformation matrix  $\hat{\mathbf{M}}$  of a CXR photo; (2) the photo is rectified using the predicted  $\hat{\mathbf{M}}$ ; (3) A

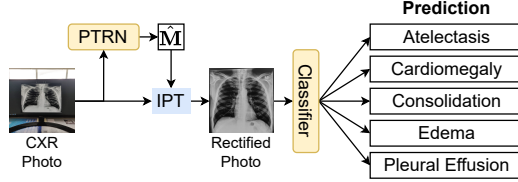


Figure 10: CXR photos classification pipeline.

classifier trained on high-quality digital CXRs evaluates the rectified CXR photo. This pipeline is used in the experiments in Sections 4.4 and 4.5.

#### 4.4. CheXphoto Competition

CheXphoto is a competition for smartphone-captured CXR photos classification hosted by Stanford and VinBrain [14]. A submitted model is tested to perform multi-label classification on the two private test sets, respectively:

1. CXR film photos. A total of 250 CXR films are captured as photos by a smartphone camera. The ranking of the leaderboard is sorted in terms of AUC on this private test set. We refer to this set as CheXphoto-Film private test set.
2. CXR monitor photos. A total of 668 digital CXRs from 500 studies are displayed on a monitor and captured as photos by an iPhone 8. We refer to this set as CheXphoto-Monitor private test set.

The performance of a model is measured by calculating the mean AUC-ROC score on 5 selected pathologies: Atelectasis, Cardiomegaly, Consolidation, Edema, and Pleural Effusion.

We build a classifier for this competition with the pipeline in Section 4.3. The classifier consists of 5 weighted average ensembles. Each ensemble outputs the predicted probability of a pathology. Each ensemble is composed of 4 to 6 trained single CNN classifiers for binary classification. A single classifier is a CNN model with an input of  $224 \times 224$  pixel grayscale image. This setup does not like other previous work, a single classifier only predicts one pathology instead of all 5 pathologies. It is because the numbers of training samples across the 5 pathologies are quite unbalanced. While training a multi-label classifier, the overfitting timings of the 5 pathologies are different. Therefore, we employ binary classifiers to obtain optimal performance for each pathology. The single classifiers are trained on the CheXpert training set with different configurations to improve the performance of the

| Method             | CXR Film Photos |              | CXR Monitor Photos |              |
|--------------------|-----------------|--------------|--------------------|--------------|
|                    | Validation      | Test         | Validation         | Test         |
| YOLOv3 [33]        | -               | -            | 0.684              | -            |
| GAN-STAM [16]      | -               | -            | 0.865              | -            |
| <b>PTRN (Ours)</b> | <b>0.868</b>    | <b>0.850</b> | <b>0.885</b>       | <b>0.891</b> |

Table 3: Quantitative results in AUC of PTRN on the CheXphoto validation sets and the CheXphoto competition private test sets.

| Rank     | Model                            | CXR Film Photos | CXR Monitor Photos |
|----------|----------------------------------|-----------------|--------------------|
| <b>1</b> | <b>LBC-v2</b>                    | <b>0.850</b>    | <b>0.89</b>        |
| <b>2</b> | <b>LBC-v0</b>                    | <b>0.820</b>    | <b>0.89</b>        |
| <b>3</b> | <b>Stellarium-CheXpert-Local</b> | <b>0.802</b>    | <b>0.88</b>        |
| 4        | MVD121                           | 0.762           | 0.83               |
| 5        | MVD121-320                       | 0.758           | 0.84               |

Table 4: The top 5 of the leaderboards of the CheXphoto competition (22 Nov 2022). The top 3 spots are all our PTRN.

constructed ensembles [58, 59]. The different configurations include different CNN models, learning rates, batch sizes, configurations of label smoothing regularization from [8] for handling uncertainty labels, and data augmentation methods such as adjusting brightness, contrast, and adding Gaussian noise. Cross entropy loss is used to calculate the loss. The validation dataset is the CheXphoto-Film validation set in which the projective transformation is perfectly rectified by manual operation.

The results are reported in Table 3. our pipeline outperforms all published work [33, 16] in the CheXphoto dataset. Moreover, our pipeline achieves **first place** on the CheXphoto competition leaderboard, as shown in Table 4. The first, second, and third places are all our work based on PTRN. We compare our pipeline with the fourth place, our approach yields a huge performance gap on CXR film photos (0.850 vs. 0.762, 0.088 higher in AUC).

#### 4.5. Classification Performance on Automatically Rectified CXR Photos

To verify the classification performance improvement in the CheXphoto competition is mostly contributed by PTRN instead of the ensembling, we compare the classification performance of a single CNN classifier with/without PTRN on CXR photos. The single CNN classifier is an Xception [60] trained on the CheXpert dataset. It achieves AUC 0.889/0.887 on the CheXpert validation/test set, respectively. The classification performance is close to the single model of #2 [8] in the CheXpert competition (validation AUC



| Method          | CheXpert     |              | CheXphoto-Monitor |              | CheXphoto-Film |              |
|-----------------|--------------|--------------|-------------------|--------------|----------------|--------------|
|                 | Validation   | Test         | Validation        | Test         | Validation     | Test         |
| CNN             | 0.889        | 0.887        | 0.821             | 0.710        | 0.722          | 0.599        |
| <b>PTRN+CNN</b> | <b>0.893</b> | <b>0.896</b> | <b>0.893</b>      | <b>0.880</b> | <b>0.791</b>   | <b>0.802</b> |

Table 5: Quantitative results in AUC of PTRN on the CheXphoto validation sets and the CheXphoto competition private test sets.

0.894). The CNN has been uploaded to the CheXphoto competition. In the leaderboard, the name of the CNN without PTRN is *Stellarium*, and the name of the CNN with PTRN is *Stellarium-CheXpert-local*. The results are shown in Table 5.

The first row reports the performance of the CNN without PTRN, huge performance drops of this CNN are observed on the CXR monitor photos and CXR film photos, since the photos experienced projective transformation with extra noises. In the second row, PTRN is used to rectify the projective transformation of CXR photos. The classification pipeline follows the one in Section 4.3.

For the results on the CXR monitor photos, it achieves AUC 0.893/0.880 on the validation/test set, respectively, which is far superior to the AUC scores before rectification (AUC 0.821/0.710). The data in the CheXphoto-Monitor validation set is the photos version of the digital CXRs in the CheXpert validation set. The performance on the CheXphoto-Monitor validation set (AUC 0.893) is the same as the performance on the CheXpert validation set (AUC 0.893). It indicates that PTRN can maintain the classification performance on CXR photos to the same level as on digital CXRs. It implies that PTRN is sufficient to eliminate all the negative impacts of projective transformation to classification performance.

In CXR film photos, huge performance improvements are observed after using PTRN (AUC 0.802/0.599, 0.203 improvement on the CheXphoto-Film private test set), which verifies the effectiveness of the PTRN on both types of photos.

Besides, we furtherly test the pipeline in the CheXpert dataset. Surprisingly, minor improvements are also observed. A further investigation is needed to verify the performance impact of PTRN on digital CXRs.

#### 4.6. Ablation Study on Training Sample Synthesis

The synthetic data framework consists of 4 steps. However, based on the formulation of Equation 1, a synthetic CXR photo can be composited by only

| Step |   |   |   | CheXphoto-Monitor Validation | CheXphoto-Film Validation   |
|------|---|---|---|------------------------------|-----------------------------|
| 1    | 2 | 3 | 4 | (for validation)             | (for test)                  |
| ✓    | ✓ | ✓ | ✓ | <b>0.942 (0.938, 0.946)</b>  | <b>0.892 (0.887, 0.898)</b> |
|      | ✓ | ✓ | ✓ | 0.911 (0.905, 0.916)         | 0.816 (0.808, 0.824)        |
| ✓    | ✓ | ✓ |   | 0.919 (0.914, 0.924)         | 0.800 (0.790, 0.809)        |
|      | ✓ | ✓ |   | 0.811 (0.875, 0.886)         | 0.815 (0.807, 0.824)        |

Table 6: Ablation study on the framework for generation of synthetic data (In mean IoU with 95% C.I.)

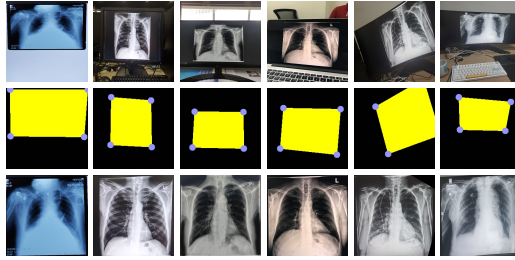


Figure 11: Top row: CXR photos. Mid row: Prediction of PTRN. Bot row: CXR photos that are rectified by using the prediction.

a transformed CXR and a background. It means in the framework, only step 2 (projective transformation) and 3 (adding background) are necessary. Step 1 (screen synthesis) and step 4 (adding natural perturbations) are additional to simulate the appearances of natural CXR photos. To study the impact of steps 1 and 4 to the rectification performance, we conducted an ablation study by removing step 1 or/and 4.

The results are shown in Table 5. Mean IoU (mIoU) with 95% confidence interval (C.I.) is reported. In the first row, the mIoU of PTRN that is trained with all four steps is reported. It achieves mIoU 0.942/0.892 on the validation/test set respectively. After removing the step(s), the mIoU on the validation set are dropped by approximately 0.02-0.07, and the mIoU on the test set is dropped much larger (approximately 0.07-0.10). It demonstrates the necessity and importance of step 1 and 4. It also verifies that the framework for the generation of synthetic CXR photos has a sufficient variance of the key variables controlling the transformed CXRs.

#### 4.7. Qualitative Evaluation

In Figure 11, we demonstrate using trained PTRN to automatically rectify six CXR photos captured in different scenarios. All six CXR photos are

properly rectified. These results demonstrate that PTRN is robust to CXR photos captured in different scenarios.

## 5. Conclusion

In this paper, we present a deep learning-based Projective Transformation Rectification Network (PTRN) that is trained with synthetic data for rectifying the projective transformation of CXR photos. We also propose a framework to generate synthetic data. Our pipeline achieves first place on the CheXphoto competition leaderboard, which has a significant improvement over [33] and [16]. This result verifies the performance of our PTRN and the framework for the generation of synthetic data. Additionally, the design of PTRN and the framework for the generation of synthetic data can be applied to other image recognition fields that encounter similar image distortion caused by the imperfect camera position.

This work has certain limitations. For example, in Section 4.2, the hyperparameters were not particularly tuned in the training of PTRN since it already produced satisfying rectification performance. The performance could be further improved by tuning hyperparameters or choosing other CNNs as the backbone. Moreover, the CXR photos classification pipeline consists of two CNN models, which requires a higher computation cost. The above limitations may lead to possible directions to extend or improve this work.

## Acknowledgment

This work is supported by Macao Polytechnic University under grant number RP/ESCA-01/2021.

## References

- [1] C. Q. Commission, et al., Radiology review: a national review of radiology reporting within the nhs in england, Care Quality Commission (2018).
- [2] R. C. of Radiologists, Unreported x-rays, computed tomography (ct) and magnetic resonance imaging (mri) scans: results of a snapshot survey of english national health service (nhs) trusts (2015).

- [3] I. M. Baltruschat, H. Nickisch, M. Grass, T. Knopp, A. Saalbach, Comparison of deep learning approaches for multi-label chest x-ray classification, *Scientific reports* 9 (1) (2019) 1–10.
- [4] X. Wang, Y. Peng, L. Lu, Z. Lu, M. Bagheri, R. M. Summers, Chestx-ray8: Hospital-scale chest x-ray database and benchmarks on weakly-supervised classification and localization of common thorax diseases, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 2097–2106.
- [5] J. Irvin, P. Rajpurkar, M. Ko, Y. Yu, S. Ciurea-Ilcus, C. Chute, H. Marklund, B. Haghighi, R. Ball, K. Shpanskaya, 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, 2019, pp. 590–597.
- [6] A. E. Johnson, T. J. Pollard, S. J. Berkowitz, N. R. Greenbaum, M. P. Lungren, C.-y. Deng, R. G. Mark, S. Horng, Mimic-cxr, a de-identified publicly available database of chest radiographs with free-text reports, *Scientific data* 6 (1) (2019) 1–8.
- [7] P. Rajpurkar, J. Irvin, K. Zhu, B. Yang, H. Mehta, T. Duan, D. Ding, A. Bagul, C. Langlotz, K. Shpanskaya, et al., Chexnet: Radiologist-level pneumonia detection on chest x-rays with deep learning, *arXiv preprint arXiv:1711.05225* (2017).
- [8] H. H. Pham, T. T. Le, D. Q. Tran, D. T. Ngo, H. Q. Nguyen, Interpreting chest x-rays via cnns that exploit hierarchical disease dependencies and uncertainty labels, *Neurocomputing* 437 (2021) 186–194.
- [9] Z. Yuan, Y. Yan, M. Sonka, T. Yang, Large-scale robust deep auc maximization: A new surrogate loss and empirical studies on medical image classification, in: *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 3040–3049.
- [10] P.-C. Kuo, C. C. Tsai, D. M. López, A. Karargyris, T. J. Pollard, A. E. Johnson, L. A. Celi, Recalibration of deep learning models for abnormality detection in smartphone-captured chest radiograph, *NPJ digital medicine* 4 (1) (2021) 1–10.

- [11] Q. Guan, Y. Huang, Y. Luo, P. Liu, M. Xu, Y. Yang, Discriminative feature learning for thorax disease classification in chest x-ray images, *IEEE Transactions on Image Processing* 30 (2021) 2476–2487.
- [12] A. B. Schwartz, G. Siddiqui, J. S. Barbieri, A. L. Akhtar, W. Kim, R. Littman-Quinn, E. F. Conant, N. K. Gupta, B. A. Pukenas, P. Ramchandani, et al., The accuracy of mobile teleradiology in the evaluation of chest x-rays, *Journal of telemedicine and telecare* 20 (8) (2014) 460–463.
- [13] S. Andronikou, K. McHugh, N. Abdurahman, B. Khoury, V. Mngomezulu, W. E. Brant, I. Cowan, M. McCulloch, N. Ford, Paediatric radiology seen from africa. part i: providing diagnostic imaging to a young population, *Pediatric radiology* 41 (7) (2011) 811–825.
- [14] N. A. Phillips, P. Rajpurkar, M. Sabini, R. Krishnan, S. Zhou, A. Pareek, N. M. Phu, C. Wang, M. Jain, N. D. Du, et al., Chexphoto: 10,000+ photos and transformations of chest x-rays for benchmarking deep learning robustness, in: *Machine Learning for Health*, PMLR, 2020, pp. 318–327.
- [15] P. Rajpurkar, A. Joshi, A. Pareek, A. Y. Ng, M. P. Lungren, Chexternal: Generalization of deep learning models for chest x-ray interpretation to photos of chest x-rays and external clinical settings, in: *Proceedings of the Conference on Health, Inference, and Learning*, 2021, pp. 125–132.
- [16] C. F. Chong, X. Yang, W. Ke, Y. Wang, Gan-based spatial transformation adversarial method for disease classification on cxr photographs by smartphones, in: *2021 Digital Image Computing: Techniques and Applications (DICTA)*, IEEE, 2021, pp. 01–08.
- [17] L. Liu, J. Gu, F. Shao, X. Liang, L. Yue, Q. Cheng, L. Zhang, et al., Application and preliminary outcomes of remote diagnosis and treatment during the covid-19 outbreak: retrospective cohort study, *JMIR mHealth and uHealth* 8 (7) (2020) e19417.
- [18] K. Karako, P. Song, Y. Chen, W. Tang, Realizing 5g-and ai-based doctor-to-doctor remote diagnosis: opportunities, challenges, and prospects, *BioScience Trends* (2020).

- [19] J. Liang, D. Doermann, H. Li, Camera-based analysis of text and documents: a survey, *International Journal of Document Analysis and Recognition (IJDAR)* 7 (2) (2005) 84–104.
- [20] C. Xue, Z. Tian, F. Zhan, S. Lu, S. Bai, Fourier document restoration for robust document dewarping and recognition, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 4573–4582.
- [21] X. Jiang, R. Long, N. Xue, Z. Yang, C. Yao, G.-S. Xia, Revisiting document image dewarping by grid regularization, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 4543–4552.
- [22] H. Bandyopadhyay, T. Dasgupta, N. Das, M. Nasipuri, Rectinet-v2: A stacked network architecture for document image dewarping, *Pattern Recognition Letters* 155 (2022) 41–47.
- [23] K. Ma, Z. Shu, X. Bai, J. Wang, D. Samaras, Docunet: Document image unwarping via a stacked u-net, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 4700–4709.
- [24] S. Das, K. Ma, Z. Shu, D. Samaras, R. Shilkrot, Dewarpnet: Single-image document unwarping with stacked 3d and 2d regression networks, in: *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 131–140.
- [25] H. Feng, W. Zhou, J. Deng, Q. Tian, H. Li, Docscanner: Robust document image rectification with progressive learning, *arXiv preprint arXiv:2110.14968* (2021).
- [26] Q. Quan, Q. Wang, Y. Du, L. Li, S. K. Zhou, Recovering medical images from ct film photos, *arXiv preprint arXiv:2203.05567* (2022).
- [27] S. M. Silva, C. R. Jung, A flexible approach for automatic license plate recognition in unconstrained scenarios, *IEEE Transactions on Intelligent Transportation Systems* (2021).
- [28] T. Björklund, A. Fiandrotti, M. Annarumma, G. Francini, E. Magli, Robust license plate recognition using neural networks trained on synthetic images, *Pattern Recognition* 93 (2019) 134–146.

- [29] S. M. Silva, C. R. Jung, License plate detection and recognition in unconstrained scenarios, in: Proceedings of the European conference on computer vision (ECCV), 2018, pp. 580–596.
- [30] H. Xu, Z.-H. Guo, D.-H. Wang, X.-D. Zhou, Y. Shi, 2d license plate recognition based on automatic perspective rectification, in: 2020 25th International Conference on Pattern Recognition (ICPR), IEEE, 2021, pp. 202–208.
- [31] H. Xu, X.-D. Zhou, Z. Li, L. Liu, C. Li, Y. Shi, Eilpr: Toward end-to-end irregular license plate recognition based on automatic perspective alignment, IEEE Transactions on Intelligent Transportation Systems 23 (3) (2021) 2586–2595.
- [32] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, C. L. Zitnick, Microsoft coco: Common objects in context, in: European conference on computer vision, Springer, 2014, pp. 740–755.
- [33] H. X. Le, P. D. Nguyen, T. H. Nguyen, K. N. Le, T. T. Nguyen, Interpretation of smartphone-captured radiographs utilizing a deep learning-based approach, arXiv preprint arXiv:2009.05951 (2020).
- [34] J. Redmon, A. Farhadi, Yolov3: An incremental improvement, arXiv preprint arXiv:1804.02767 (2018).
- [35] C. Solomon, T. Breckon, Fundamentals of Digital Image Processing: A practical approach with examples in Matlab, John Wiley & Sons, 2011.
- [36] A. Creswell, T. White, V. Dumoulin, K. Arulkumaran, B. Sengupta, A. A. Bharath, Generative adversarial networks: An overview, IEEE signal processing magazine 35 (1) (2018) 53–65.
- [37] K. Wang, C. Gou, Y. Duan, Y. Lin, X. Zheng, F.-Y. Wang, Generative adversarial networks: introduction and outlook, IEEE/CAA Journal of Automatica Sinica 4 (4) (2017) 588–598.
- [38] Y. Fang, C. Deyun, W. Rui, A distortion correction approach on natural scene text image, in: Proceedings of 2011 6th International Forum on Strategic Technology, Vol. 2, IEEE, 2011, pp. 1058–1061.

- [39] J. Liang, D. DeMenthon, D. Doermann, Geometric rectification of camera-captured document images, *IEEE transactions on pattern analysis and machine intelligence* 30 (4) (2008) 591–605.
- [40] Y. Takezawa, M. Hasegawa, S. Tabbone, Robust perspective rectification of camera-captured document images, in: *2017 14th IAPR International Conference on Document Analysis and Recognition (ICDAR)*, Vol. 6, IEEE, 2017, pp. 27–32.
- [41] Y.-H. Chang, C.-H. Chu, M.-S. Chen, A general scheme for extracting qr code from a non-uniform background in camera phones and applications, in: *Ninth IEEE International Symposium on Multimedia (ISM 2007)*, IEEE, 2007, pp. 123–130.
- [42] H. Tribak, Y. Zaz, Qr code recognition based on principal components analysis method, *International Journal of Advanced Computer Science and Applications* 8 (4) (2017).
- [43] S. R. Deans, *The Radon transform and some of its applications*, Courier Corporation, 2007.
- [44] M. Jaderberg, K. Simonyan, A. Zisserman, et al., Spatial transformer networks, *Advances in neural information processing systems* 28 (2015).
- [45] Y. Xiao, Z. Ming, 1d barcode detection via integrated deep-learning and geometric approach, *Applied Sciences* 9 (16) (2019) 3268.
- [46] F. Zhan, S. Lu, Esir: End-to-end scene text recognition via iterative image rectification, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 2059–2068.
- [47] C. Xue, S. Lu, S. Hoi, Detection and rectification of arbitrary shaped scene texts by using text keypoints and links, *Pattern Recognition* 124 (2022) 108494.
- [48] A. Gupta, A. Vedaldi, A. Zisserman, Synthetic data for text localisation in natural images, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 2315–2324.
- [49] D. Dwibedi, I. Misra, M. Hebert, Cut, paste and learn: Surprisingly easy synthesis for instance detection, in: *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 1301–1310.



- [50] F. Zhan, S. Lu, C. Xue, Verisimilar image synthesis for accurate detection and recognition of texts in scenes, in: Proceedings of the European Conference on Computer Vision (ECCV), 2018, pp. 249–266.
- [51] C.-H. Lin, E. Yumer, O. Wang, E. Shechtman, S. Lucey, St-gan: Spatial transformer generative adversarial networks for image compositing, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2018, pp. 9455–9464.
- [52] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 770–778.
- [53] G. Huang, Z. Liu, L. Van Der Maaten, K. Q. Weinberger, Densely connected convolutional networks, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2017, pp. 4700–4708.
- [54] R. I. Hartley, Theory and practice of projective rectification, International Journal of Computer Vision 35 (2) (1999) 115–127.
- [55] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin, S. Ghemawat, I. Goodfellow, A. Harp, G. Irving, M. Isard, Y. Jia, R. Jozefowicz, L. Kaiser, M. Kudlur, J. Levenberg, D. Mané, R. Monga, S. Moore, D. Murray, C. Olah, M. Schuster, J. Shlens, B. Steiner, I. Sutskever, K. Talwar, P. Tucker, V. Vanhoucke, V. Vasudevan, F. Viégas, O. Vinyals, P. Warden, M. Wattemberg, M. Wicke, Y. Yu, X. Zheng, TensorFlow: Large-scale machine learning on heterogeneous systems, software available from tensorflow.org (2015).  
URL <https://www.tensorflow.org/>
- [56] A. B. Jung, K. Wada, J. Crall, S. Tanaka, J. Graving, C. Reinders, S. Yadav, J. Banerjee, G. Vecsei, A. Kraft, Z. Rui, J. Borovec, C. Vallentin, S. Zhydenko, K. Pfeiffer, B. Cook, I. Fernández, F.-M. De Rainville, C.-H. Weng, A. Ayala-Acevedo, R. Meudec, M. Laporte, et al., imgaug, <https://github.com/aleju/imgaug>, online; accessed 01-Feb-2020 (2020).
- [57] D. P. Kingma, J. Ba, Adam: A method for stochastic optimization, arXiv preprint arXiv:1412.6980 (2014).

- [58] Z.-H. Zhou, Ensemble methods: foundations and algorithms, CRC press, 2012.
- [59] L. K. Hansen, P. Salamon, Neural network ensembles, IEEE transactions on pattern analysis and machine intelligence 12 (10) (1990) 993–1001.
- [60] F. Chollet, Xception: Deep learning with depthwise separable convolutions, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2017, pp. 1251–1258.