

Supervised local error estimation for nonlinear image registration using convolutional neural networks

Koen A.J. Eppenhof^a and Josien P.W. Pluim^{a, b}

^aMedical Image Analysis (IMAG/e), Department of Biomedical Engineering,
Eindhoven University of Technology, Eindhoven, The Netherlands

^bImage Sciences Institute, University Medical Center Utrecht, Utrecht, The Netherlands

ABSTRACT

Error estimation in medical image registration is valuable when validating, comparing, or combining registration methods. To validate a nonlinear image registration method, ideally the registration error should be known for the entire image domain. We propose a supervised method for the estimation of a registration error map for nonlinear image registration. The method is based on a convolutional neural network that estimates the norm of the residual deformation from patches around each pixel in two registered images. This norm is interpreted as the registration error, and is defined for every pixel in the image domain. The network is trained using a set of artificially deformed images. Each training example is a pair of images: the original image, and a random deformation of that image. No manually labeled ground truth error is required. At test time, only the two registered images are required as input. We train and validate the network on registrations in a set of 2D digital subtraction angiography sequences, such that errors up to eight pixels can be estimated. We show that for this range of errors the convolutional network is able to learn the registration error in pairs of 2D registered images at subpixel precision. Finally, we present a proof of principle for the extension to 3D registration problems in chest CTs, showing that the method has the potential to estimate errors in 3D registration problems.

Keywords: nonlinear image registration, registration validation, registration error estimation, convolutional networks

1. INTRODUCTION

Assessing the accuracy of nonlinear image registration is a difficult problem. Often surrogate metrics are used to estimate the registration quality. These metrics measure the registration error indirectly using scores based on tissue overlap (e.g. Dice or Jaccard scores), image similarity metrics (e.g. mutual information and normalized correlation), and consistency metrics (e.g. inverse consistency error). Because these metrics do not necessarily correlate with the true registration error, they have to be used with caution.¹ An alternative is to use a dense set of corresponding landmarks to measure the registration error between the registered images. Generally, such dense sets of landmarks are hard to obtain, and require considerable manual labor.

For true quantification of registration accuracy it is useful to have quantitative estimates of the registration error for every pixel, i.e. an error map for the entire image domain. Such an error map can be used in practice to assess the quality or confidence of the registration locally, or to boost the registration by combining an ensemble of registration algorithms based on their local error.² It is desirable that an error map estimation method is independent of the registration method. This is especially important when the error estimates are used to compare multiple registration methods.

We propose a supervised method that maps a pair of registered images to a continuous image registration error map. The method uses a convolutional neural network that is trained to estimate the norm of the residual displacement in the center pixel of a pair of image patches. When trained, the network can be applied directly on a pair of registered images. In this paper, we evaluate this method on registration of frames from 2D digital subtraction angiography (DSA) sequences. Additionally, we show a proof of concept on 3D images using thoracic CT data.

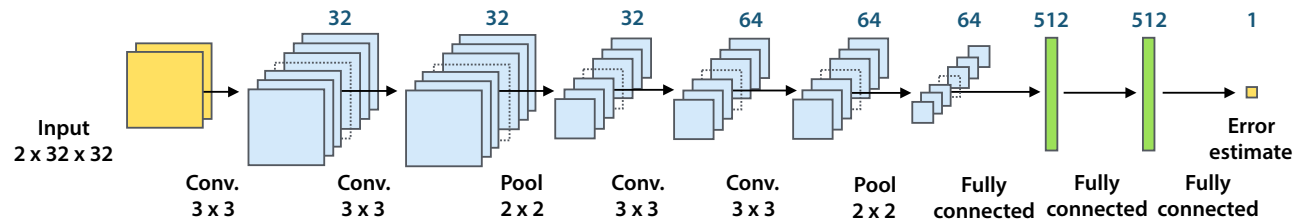


Figure 1. Convolutional network architecture. The network takes two 2D image patches of the same size as input and outputs the norm of the deformation at the patches' center pixel. The convolutional and fully connected layers all have ReLU activation functions. Dropout was applied to the first two fully connected layers.

2. METHOD

Let $I_F(\mathbf{x})$ and $I_M(\mathbf{x})$ be two images. Registration of the moving image $I_M(\mathbf{x})$ to the fixed image $I_F(\mathbf{x})$ aims to estimate the displacement $\mathbf{u}(\mathbf{x})$ that is needed to spatially align $I_M(\mathbf{x} + \mathbf{u}(\mathbf{x}))$ to $I_F(\mathbf{x})$. The estimated displacement $\hat{\mathbf{u}}(\mathbf{x})$ found in the registration will generally differ from the true displacement $\mathbf{u}(\mathbf{x})$. We define the registration error map as the L_2 -norm of this difference $\varepsilon(\mathbf{x}) = \|\mathbf{u}(\mathbf{x}) - \hat{\mathbf{u}}(\mathbf{x})\|$ for every pixel \mathbf{x} . We train the convolutional neural network to estimate the error map based on 32×32 image patches around each pixel in a pair of fixed and registered images. During training, these images are artificially deformed images, while testing is performed on simulated registration problems.

2.1 Convolutional network architecture

The network architecture is illustrated in Figure 1. The two input patches are treated as two channels of the same image. The network consists of two convolutional layers followed by a max-pooling layer, followed by two more convolutional layers, and another max-pooling layer. The convolutional layers have 3×3 kernels. The pooling layers have 2×2 windows, and use a stride of 2×2 , subsampling their inputs by a factor two along every axis. A trio of fully connected layers produces the final result: a single value for the estimate of the registration error between the two input patches. All convolutional and fully connected layers are equipped with a rectified linear unit (ReLU) activation function. Additionally, the first two fully connected layers have a dropout probability of 50% at training time to avoid overfitting.³

2.2 Training set

The training set was created by randomly deforming an image, and extracting pairs of patches from the deformed and original images. The deformations were obtained by sampling a 10×10 grid of random two-dimensional vectors sampled from a uniform distribution. Deformations for the full image were obtained by bicubic interpolation of this grid. The deformations were scaled to have a maximal norm of eight pixels, which we assume is close to the maximal deformation norm in image registration for the set of 2D DSA images we used. An example is shown in Figure 2.

2.3 Training

The network was trained using the open-source deep learning framework Caffe.⁴ Stochastic gradient descent (SGD) with mini-batches of 64 samples was used to minimize the Euclidean loss function

$$L = \frac{1}{2N} \sum_{k=1}^N (\varepsilon(\mathbf{x}_k) - \hat{\varepsilon}(\mathbf{x}_k))^2 \quad (1)$$

where $\hat{\varepsilon}_k(\mathbf{x})$ is the network output for example k , $\varepsilon_k(\mathbf{x})$ the displacement norm for example k , and $N = 64$ is the batch size. We use a momentum of 0.5 such that the previously seen training samples have an influence on the current weight update, which reduces oscillations in the weights.

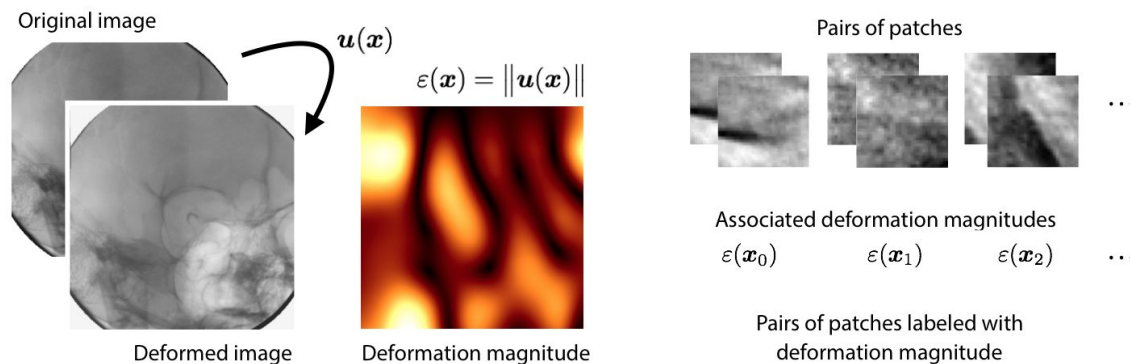


Figure 2. The training set consists of patches from an image, and a randomly deformed version of that image. Using the images and the deformation magnitude, a training set of example patches with associated deformation is created.

2.4 Testing data

To validate the method, we construct gold standard error maps in the following way: We registered two images using the registration toolbox `elastix`.⁵ The images were registered using a cubic B-spline transform and the mutual information similarity metric regularized with a bending energy potential term, which was optimized using adaptive stochastic gradient descent.⁶ We used a multiresolution registration algorithm with four resolutions and a B-spline grid spacing of 8×8 pixels at the finest resolution, with each coarser resolution doubling the grid spacing. The registered image and associated displacements we call $I^*(x)$ and $u(x)$ respectively. Using the three lowest resolutions of this registration, we created a less accurately registered image $I^\dagger(x)$ and displacements $u^\dagger(x)$. This less accurate registration omits registration at the final, finest resolution. The registered images $I^*(x)$ and $I^\dagger(x)$ are the input to the neural network. The norm of the difference between the displacements $I^*(x)$ and $I^\dagger(x)$, i.e. $\epsilon^*(x) = \|u^*(x) - u^\dagger(x)\|$, is the simulated gold standard error which the network's output will be compared to.

3. EXPERIMENTS

3.1 Materials

To train, test, and validate the method, we used 22 sequences of cerebral digital subtraction angiography (DSA) images.⁷ Each sequence shows the passage of a bolus of contrast agent through a part of the vasculature. In the center image of the sequences, bolus passage is maximal. By subtracting an image prior to the arrival of the bolus from an image in which contrast agent is present, an angiogram can be produced, provided these images are registered. Each of the images in the set has a resolution of 512×512 pixels. For training and testing, the first, center, and last frames of eleven sequences were used, ensuring inclusion of images with and without contrast agent. The remaining sequences were used for validation, where from each sequence the first frame and the center frame of the sequence were selected. The training images were randomly deformed as described in Section 2.2. To increase the number of training examples, this was repeated for ten random deformations of each image. The 330 pairs of images were divided into 4096 pairs of patches each, resulting in 1,351,680 examples in total. This set was randomly shuffled and further divided into a training set containing 90% of the data, and a validation set containing the remaining 10%. The other eleven sequences were used for testing the method, and used to create gold standard error maps as described in Section 2.4.

3.2 Results

To compare the distributions of the gold standard and estimated errors, we used violin plots, which show the distributions on either side of a vertical line for every registration pair in the test set (Figure 3). For pairs 7 and 11 we can see that larger errors are underestimated, which results in a bump in the distribution in the six-to-eight pixel range. Underestimation of these errors is also shown when we make a correlation plot of the estimates and the gold standard (Figure 5).

Table 1. RMSD, normalized RMSD, and Pearson's correlation coefficient for the registration pairs in the test set.

Registration pair	1	2	3	4	5	6	7	8	9	10	11
RMSD (pixels)	0.148	0.754	0.587	0.193	0.139	1.013	5.850	0.368	0.183	0.401	5.963
NMRSD	0.071	0.042	0.051	0.042	0.049	0.212	0.152	0.074	0.074	0.037	0.169
Pearson's r	0.940	0.954	0.967	0.936	0.977	0.550	0.863	0.915	0.940	0.982	0.823

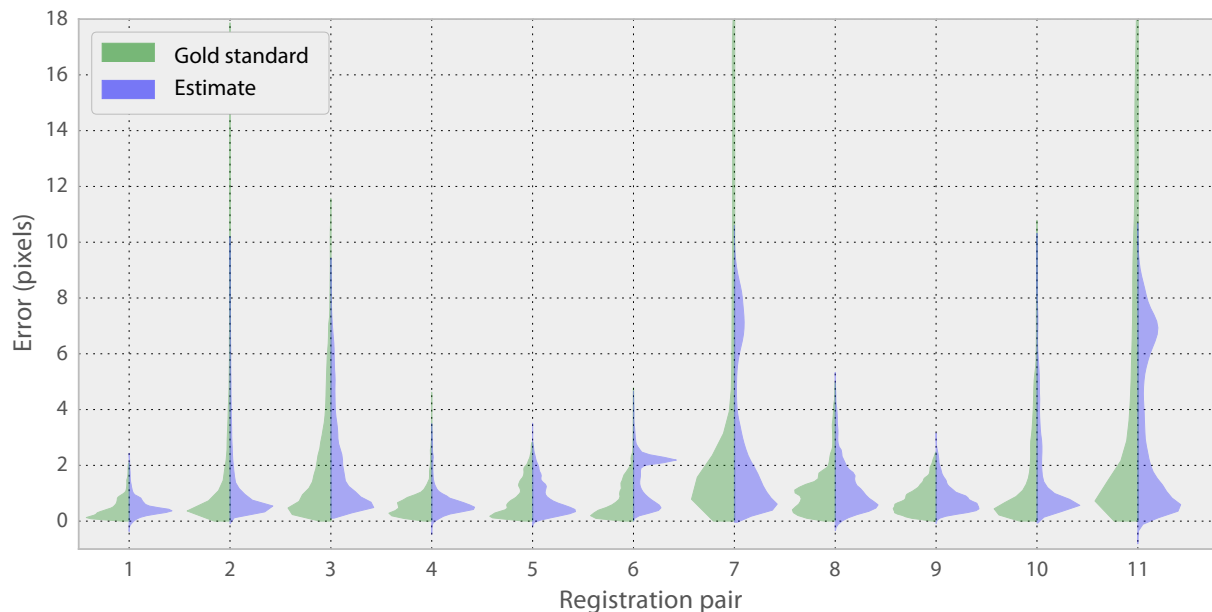


Figure 3. Violin plots showing the probability densities of the gold standard and estimated errors.

We calculated Pearson's r correlation coefficient and the root mean square differences (RMSD) between the gold standard error and the estimated error. To be able to compare registration pairs that have different ranges of registration errors, we also calculate the normalized RMSD (NRMSD), defined as $\text{NRMSD} = \frac{\text{RMSD}}{\varepsilon_{\max}^* - \varepsilon_{\min}^*}$. These values are displayed in Table 1. The only pairs having a larger NRMSD than 10% are pairs 6, 7, and 11. Figure 4 shows estimates for the two best results in terms of NRMSD (registration pairs 2 and 10) and the two poorest estimates (pairs 7 and 6).

3.3 Application to 3D images

The presented method has been trained and evaluated on 2D images. We also include a proof of principle for the application of the method to 3D images. This requires upgrading the convolutional network with 3D convolution and pooling operations. Future work will involve training and evaluating the method on chest 3D CT scans in the DURLAB data set.⁸ A preliminary result is shown in Figure 6.

4. DISCUSSION

The violin plots show that the network distribution of the error estimates is similar to the gold standard error distribution, except when the error range exceeds eight pixels. This can be explained by the absence of larger errors in the training set. The correlation plot (Figure 5) shows that the network is not able to extrapolate to larger values, but instead returns underestimates in the range of six-to-eight pixels. Consequently, the high precision for error estimates smaller than six pixels cannot be guaranteed for larger errors.

Qualitative comparisons of gold standard and estimated error maps in Figure 4 show good agreement. The poor performance for registration pair 6 (Figure 4D) can be attributed to the large homogeneous areas in the

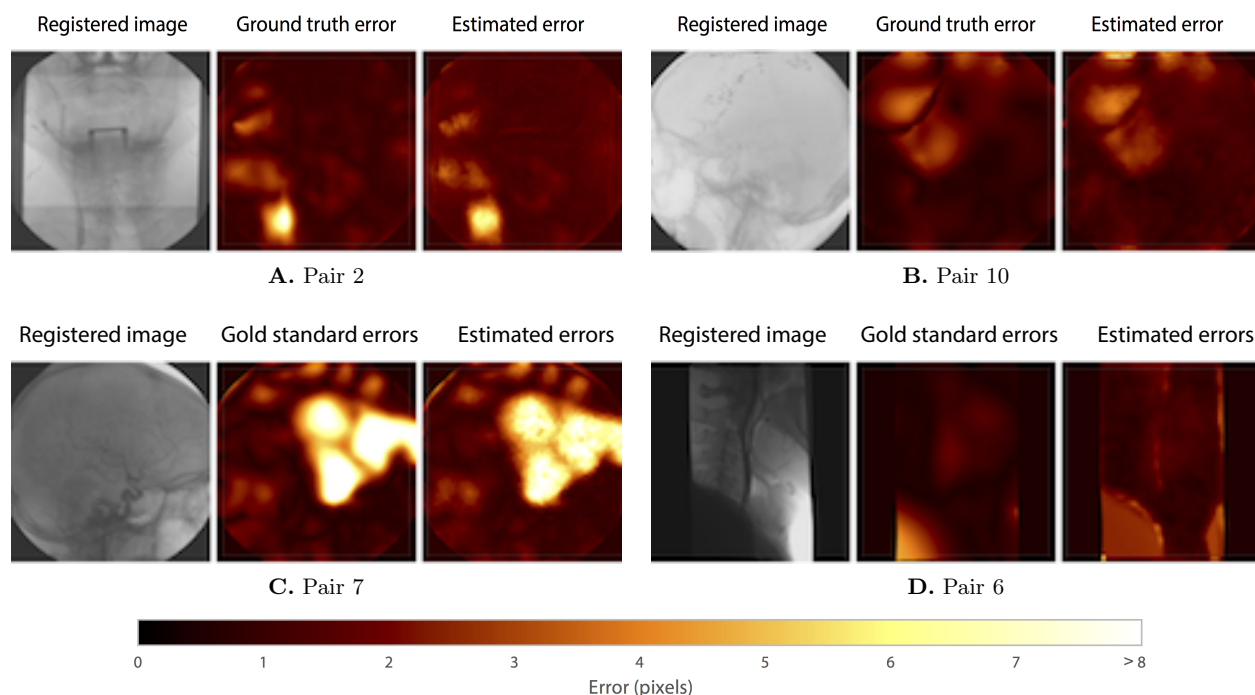


Figure 4. The registered image, gold standard error map, and the estimated error map for the two best results (A, B) and the two poorest results (C, D) in terms of NRMSD.

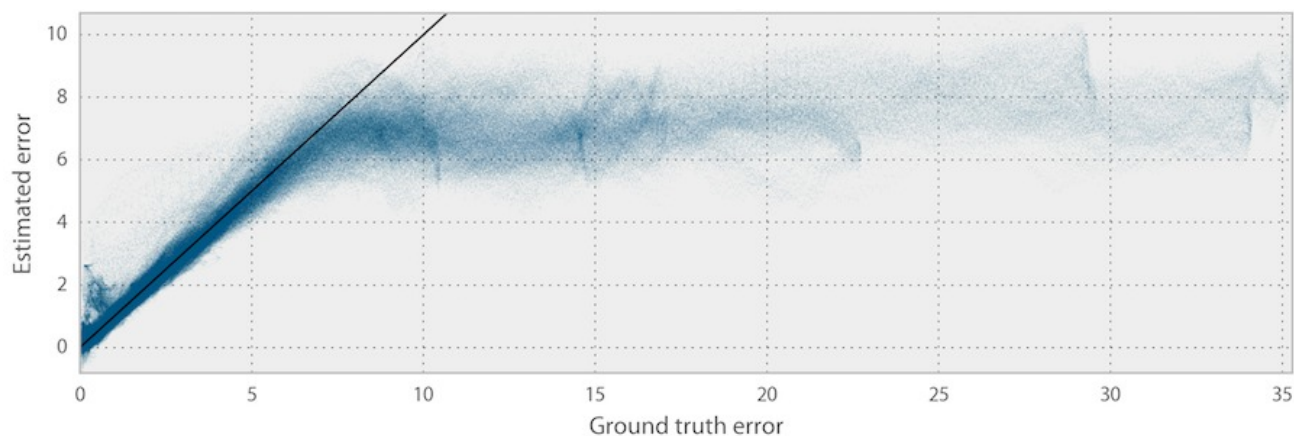


Figure 5. Correlation plot for all 4.8 million pixels in the 11 testing registration pairs. The error estimate has high correlation with the gold standard for errors below eight pixels. The black line shows the identity function.

image that coincide with a large error in the gold standard error map. It is suspected that the network relies on texture to find the displacement magnitude between image patches, and is unable to do so for areas with low contrast. In the current implementation, a maximal error of eight pixels is used in the training set. Although this maximum can easily be increased, in our experiments we found that this will decrease the performance on the smaller deformations and comes at the expense of longer training and testing times.

We show that a set of synthetically deformed images can adequately train the convolutional network. A manually labeled ground truth is therefore not required. The use of data augmentation permits training on only a small set of representative images – in this case eleven images.

Intensity-based registration methods and DSA images were used to evaluate the method. Because only the registered images are required as input, the methodology is independent of the registration method or image

modality, and application to other kinds of images would only require a small set of representative images to train the network on.

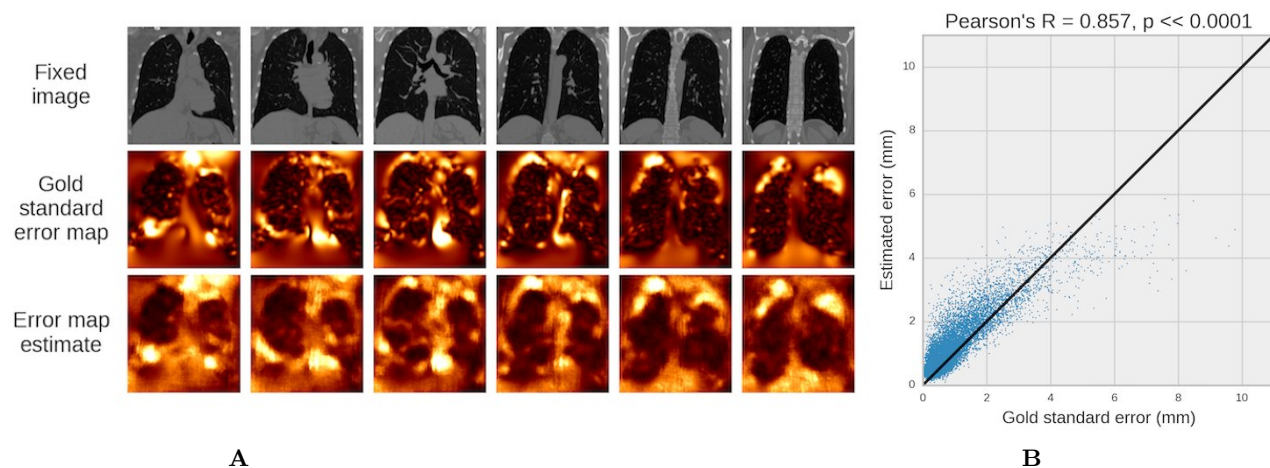


Figure 6. **A.** Example of DIRLAB 3D CT with gold standard error map and error map estimate. **B.** Correlation plot of registration error estimates for voxels inside the lungs.

5. CONCLUSION

In this paper we propose a supervised method for quantitative estimation of nonlinear registration error maps. The method uses a convolutional neural network that is trained using synthetically deformed images, and does not require a manually labeled ground truth. The network learns deformation magnitudes in a patch-wise way, and can be efficiently applied to a pair of registered images at test time. Evaluation of the method shows that this approach is very accurate for the range of errors that are present in the training set. The methodology can in principle be applied to any mono-modal registration method. We have shown the methodology can be extended to 3D images using 3D convolutional neural networks. Further development and validation of this 3D method will be the topic of future research.

REFERENCES

- [1] Rohlfing, T., "Image similarity and tissue overlaps as surrogates for image registration accuracy: Widely used but unreliable," *IEEE Transactions on Medical Imaging* **31**(2), 153–163 (2012).
- [2] Muenzing, S., van Ginneken, B., Viergever, M., and Pluim, J., "Dirboost—an algorithm for boosting deformable image registration: Application to lung ct intra-subject registration," *Medical Image Analysis* **18**(3), 449–459 (2014).
- [3] Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., and Salakhutdinov, R., "Dropout: A simple way to prevent neural networks from overfitting," *Journal of Machine Learning Research* **15**, 1929–1958 (2014).
- [4] Jia, Y., Shelhamer, E., Donahue, J., Karayev, S., Long, J., Girshick, R., Guadarrama, S., and Darrell, T., "Caffe: Convolutional architecture for fast feature embedding," *arXiv:1408.5093* (2014).
- [5] Klein, S., Staring, M., Murphy, K., Viergever, M., and Pluim, J., "elastix: A toolbox for intensity-based medical image registration," *IEEE Transactions on Medical Imaging* **29**(1), 196–205 (2010).
- [6] Klein, S., Pluim, J. P. W., Staring, M., and Viergever, M. A., "Adaptive stochastic gradient descent optimisation for image registration," *International Journal of Computer Vision* **81**(3), 227–239 (2008).
- [7] Meijering, E. H., Zuiderveld, K. J., and Viergever, M. A., "Image registration for digital subtraction angiography," *International Journal of Computer Vision* **31**(2/3), 227–246 (1999).
- [8] Castillo, R., Castillo, E., Fuentes, D., Ahmad, M., Wood, A. M., Ludwig, M. S., and Guerrero, T., "A reference dataset for deformable image registration spatial accuracy evaluation using the copdgene study archive," *Physics in Medicine and Biology* **58**(9), 2861–2877 (2013).