

FAIM – A ConvNet Method for Unsupervised 3D Medical Image Registration

Dongyang Kuang¹ and Tanya Schmah²

University of Ottawa, Ottawa, Canada.

dykuangii@gmail.com

University of Ottawa, Ottawa, Canada.

tschmah@uottawa.ca

Abstract. We present a new unsupervised learning algorithm, “FAIM”, for 3D medical image registration. With a different architecture than the popular “U-net” [9], the network takes a pair of full image volumes and predicts the displacement fields needed to register source to target. Compared with “U-net” based registration networks such as VoxelMorph [2], FAIM has fewer trainable parameters but can achieve higher registration accuracy as judged by Dice score on region labels in the Mindboggle-101 dataset. Moreover, with the proposed penalty loss on negative Jacobian determinants, FAIM produces deformations with many fewer “foldings”, i.e. regions of non-invertibility where the surface folds over itself. In our experiment, we varied the strength of this penalty and investigated changes in registration accuracy and non-invertibility in terms of number of “folding” locations. We found that FAIM is able to maintain both the advantages of higher accuracy and fewer “folding” locations over VoxelMorph, over a range of hyper-parameters (with the same values used for both networks). Further, when trading off registration accuracy for better invertibility, FAIM required less sacrifice of registration accuracy. Codes for this paper will be released upon publication.

Keywords: Image registration · Convolutional neural network · Unsupervised registration · Folding penalization

1 Introduction

Image registration is a key element of medical image analysis. The spatial deformations required to optimally register images are highly non-linear, especially for regions such as the cerebral cortex, the folding patterns of which can vary significantly between individuals. Most state-of-the-art registration algorithms, such as ANTs [1], use geometric or variational methods that are guaranteed to produce diffeomorphisms, i.e. smooth invertible deformations with a smooth inverse. These algorithms are very computationally intensive and still do not generally find optimal deformations. One general problem is that the optimization problems solved by these algorithms are highly nonconvex. Another is that they treat each pair of images to be registered *de novo*, without any learning.

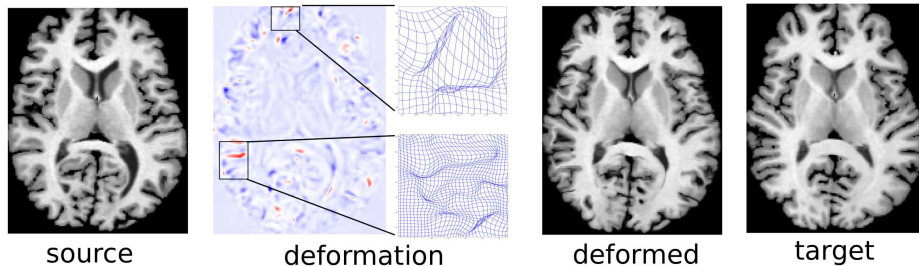


Fig. 1: An axial slice of a deformation produced by a CNN method: VoxelMorph-1, with its default L_2 regularization parameter $\lambda = 1$ on spatial gradients. The first and last images in the row are the source and target images, while the third one is the deformed source image produced by the method. The second image in the row shows values of the Jacobian determinant of the predicted deformation, with “folding” locations (negative determinant) marked in red. The deformed grids illustrate parts of the deformation.

A revolution is taking place in the last few years in the application of machine learning methods to medical image processing, including registration tasks. Supervised methods for registration, as in [14,11,8], learn from known reference deformations for training data – either actual “ground truth” in the case of synthetic image pairs, or deformations computed by other automatic or semiautomatic methods. Unsupervised methods, as in [7,13,10,2], do not require reference deformations, but instead minimize some cost function modeling the goodness of registration, optionally regularized by a term constraining the deformation. These methods have properties complementary to the standard geometric methods: they are very fast (at test time) and have the ability to learn automatically from data; however the predicted deformations are not guaranteed to be diffeomorphisms. In particular, there are often many regions where one image has been “folded” over itself by a non-invertible transformation. In these regions the Jacobian matrix of the deformation has negative determinant, as shown in Figure 1. These spatial foldings are not physically possible and thus constitute registration errors when used in clinical applications. The frequency of this kind of error has limited the adoption of neural network methods in medical image registration.

To address this problem, we propose a new unsupervised image registration algorithm, FAIM (for FAsT IMage registration) with an explicit anti-folding regularization. Using the MindBoggle101 dataset [6], we compared FAIM’s response on both registration accuracy and anti-folding performance with an U-net[9] based network VoxelMorph [2]. We also examined the trade-off behavior on accuracy and number of foldings on both networks.

2 Methods

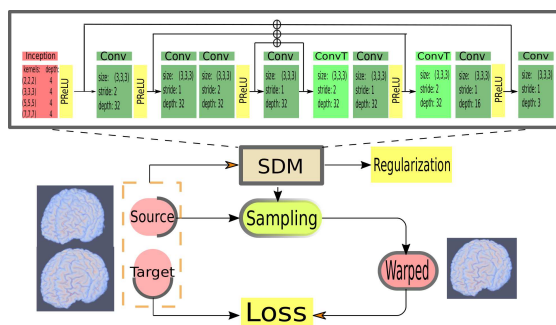


Fig. 2: FAIM network architecture.

Our architecture is directly inspired by the spatial transformer network (STN) of Jaderberg et al. [4], which is used to learn the proper parametrized transformation of the input feature so that later tasks such as classifications can be better performed. This kind of module is originally developed for 2D images and only affine and thin

plate spline transformation were implemented. In some very recent research [10,2], this framework begins to appear in 3D medical image registration. All these works aim to find an optimal parametrized transformation $\phi : \Omega \rightarrow \mathbf{R}^3$, for image domain $\Omega \subset \mathbf{R}^3$, such that the warped volume $S \circ \phi^{-1}(x)$ from a moving/source volume $S(x)$ is well aligned to the fixed/target volume $T(x)$. In our network, we use displacement field $\mathbf{u}(x)$ to parametrize the deformation ϕ by $S \circ \phi^{-1}(x) = S(x + \mathbf{u}(x))$, which is learned through a spatial deformation module (SDM). Figure 2 shows the flow chart when the network is in training and a closer look at the SDM.

During training, the moving volume and the target volume are stacked together as the input feeding into SDM. The first layer is inspired by Google’s Inception module [12]. The purpose of this layer is trying to compare and capture information at different spatial scales for later registration. PReLU [3] activations are used at the end of each convolutional block except the last layer which uses linear activation to produce displacement fields. The sampling module then takes the displacements and generates a deformed grid and use it to sample the source image to produce the warped image. We use kernel stride > 1 to reduce the volume size instead of inserting max pooling layers. Transposed convolutional layers are used for upsampling. There are three “add” skip connections between the downsampling and upsampling path to help the gradient flow.

The total training loss is the sum of an image dissimilarity term L_{image} and regularization terms $L_{total} = L_{image}(S, T) + \alpha R_1(\mathbf{u}) + \beta R_2(\mathbf{u})$, defined in Table 1. The main loss L with cross correlation (CC) in this paper is for the similarity between the warped source and target, while the first regularization term R_1 regularizes the overall smoothness of the predicted displacements. The second regularization aims specifically at penalizing transformations that have many negative Jacobian determinants. Transformations that have all non-negative Jacobian determinants will not be penalized.

$L_{image}(S, T): 1 - CC(S \circ \phi^{-1}, T)$
Regularization: $R_1(\mathbf{u}) = \ D\mathbf{u}\ _2$
Regularization: $R_2(\mathbf{u}) = 0.5 (\det(D\phi^{-1}) - \det(D\phi^{-1}))$

Table 1: Loss and regularization functions used.

3 Experiments

3.1 Mindboggle101 dataset

This dataset, created by Klein et al. [6], is based on a collection of 101 T1-weighted MRIs from healthy subjects. The Freesurfer package (<http://www.martinos.org/freesurfer>) was used to preprocess all images, and then automatically label the cortex using its DK cortical parcellation atlas. For 54 of the images, including the OASIS-TRT-20 subset, these automatic parcellations were manually edited to follow a custom labeling protocol, DKT. We use the variant DKT25, with 25 cortical regions per hemisphere. Details of data collection and processing, including atlas creation, are described in [6].

In the present paper, we used brain volumes from the following three named subsets of Mindboggle101, for a total of 62 volumes: NKI-RS-22, NKI-TRT-20 and OASIS-TRT-20. These images are already warped to MNI152 space. We normalized the intensity of each brain volume by its maximum voxel intensity. Each image has dimensions $182 \times 218 \times 182$, which we truncated to $144 \times 180 \times 144$. With this resolution, FAIM has 179,787 trainable parameters, which is about only 70% of VoxelMorph’s 259,675 trainable parameters.

Figure 4 shows the region corresponding to one label in the parcellation. The geometrical complexity of this cortical surface parcellation leads to very challenging registration tasks.

3.2 Evaluation Methods

We divide each dataset into sets of training and test images, and use these to form training and test sets of *pairs* of images. The training set consists of all ordered brain volume pairs¹ from the union of the NKI-RS-22 and NKI-TRT-20 subsets (1722 pairs in total), and the test set consists of all ordered pairs from the OASIS-TRT-20 subset (380 pairs in total). We train FAIM and VoxelMorph on all pairs of images from the training set, and then examine their predicted deformations with pairs of images from the test set. The Adam optimizer [5] is used. When not otherwise specified, both networks are trained on our training set with the same hyperparameters: learning rate = 10^{-4} , epochs = 10, $\alpha = 1$.

¹ Their corresponding labels are not used in training.

We use predicted deformations to warp corresponding ROI labels from source to target per pair. Registration accuracy is primarily evaluated using the Dice score.

$$\mathbf{Dice}(X, Y) = 2 \frac{|X \cap Y|}{|X| + |Y|}. \quad (1)$$

It measures the degree of overlap between corresponding regions in the parcellations associated with each image. The quality of the predicted deformations ϕ is assessed by the total number of locations where Jacobian determinant $\det(\nabla\phi^{-1}(x))$ are negative,

$$\mathcal{N} := \sum \delta(\det(D\phi^{-1}) < 0).$$

3.3 Results

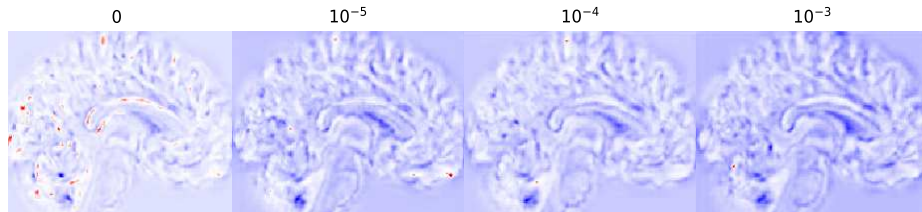


Fig. 3: Locations where $\det(D\phi^{-1}) < 0$ (marked in red) with different β shown on one slice. Predictions here are done using FAIM.

Figure 3 visualizes the effect of the second regularization term $R_2(\mathbf{u})$ that penalizes “foldings” directly during training. When the regularization is not used, $\beta = 0$, there are multiple locations visible in the transformation whose Jacobian determinant are negative. The number is greatly reduced with $\beta = 10^{-5}$, and almost eliminated at higher β values. Numerical results are given in Table 2. Figure 4 provides a visualization of one predicted label rendered in 3D. AntsSynQuick, Voxelmorph and FAIM appear to produce quite similar results on this label, but the underlying transformations are different as shown in later detailed comparisons.

We selected 5 scales of regularization strength β from 0 to 10^{-2} and trained both FAIM and VoxelMorph under the same hyper-parameters. We summarize the mean Dice score across all predicted ROI labels with their corresponding target labels in the test set and mean \mathcal{N} (i.e. $\overline{\mathcal{N}}$) of all predicted deformations in the test set in Table 2. As one can see in the table, FAIM has higher registration accuracy under all the considered β values and lower number of “foldings” in the predicted deformations when increasing β . A more detailed comparison on accuracy in terms of Dice score with $\beta = 10^{-3}$ and relations among Dice score, $\overline{\mathcal{N}}$ and β are listed together in Figure 5.

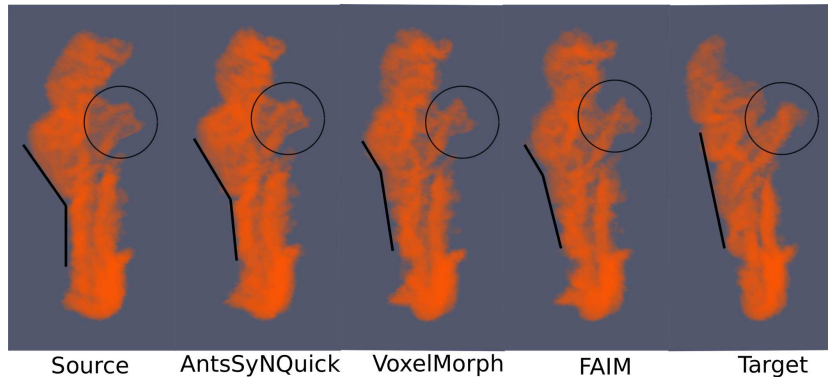
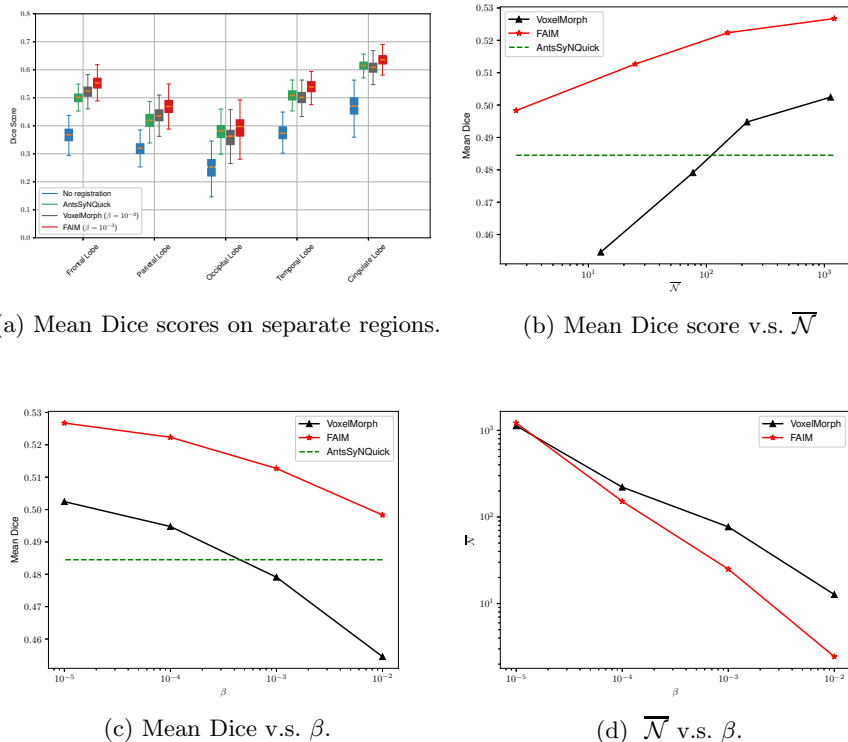


Fig. 4: One label (left superior parietal) for one source-target image pair, and the warped source labels produced by different methods. Notice that all three methods are aware of correct regions needed to deform such as thinning the part marked by the circled and straightening the region indicated by the black lines.

Mean Dice	$\beta = 0$	10^{-5}	10^{-4}	10^{-3}	10^{-2}
VoxelMorph	0.5066	0.5024	0.4948	0.4791	0.4545
FAIM	0.5330	0.5267	0.5230	0.5126	0.4983
Mean \mathcal{N}	$\beta = 0$	10^{-5}	10^{-4}	10^{-3}	10^{-2}
VoxelMorph	49406	1129	221	77	13
FAIM	59115	1215	151	25	2

Table 2: Mean Dice scores and mean number of “folding” locations with different β values. For comparison, the mean Dice score for ANTs SyNQuick is 0.4845.

From Figure 5 (a), FAIM has higher registration accuracy among the three compared methods in all the five regions on the brain. Figure 5 (c) suggests this advantage of FAIM in accuracy is consistent across different values of β , with a mean improvement of approximately 3%. To investigate how the networks balance the two competing tasks of high accuracy and low number of “folding” locations, we plotted mean Dice score against mean number of “folding” locations in Figure 5 (b). In this figure, the flatter curve from FAIM suggests the accuracy of it is more robust with respect to numbers of “folding” locations in its predictions when compared with VoxelMorph. In other words, the higher slope for VoxelMorph shows that to achieve the same gain in reducing number of “foldings”, U-net based VoxelMorph has to sacrifice more in registration accuracy. Finally, we check the sensitivity of the control of β over negative Jacobian determinants in Figure 5 (d) by visualizing \mathcal{N} against β . The sharper slope of FAIM in this log-log plot reveals that we will have more gain in reducing negative Jacobian determinant per unit increase of the regularization strength β when compared with VoxelMorph.



(a) Mean Dice scores on separate regions.

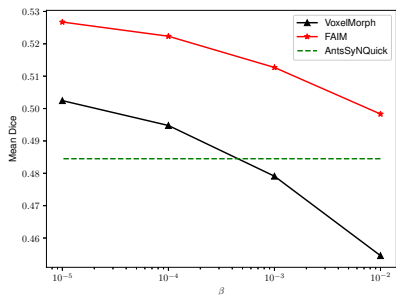
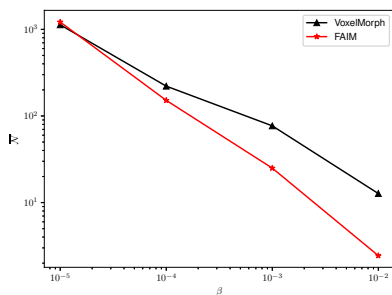
(b) Mean Dice score v.s. \bar{N} (c) Mean Dice v.s. β .(d) \bar{N} v.s. β .

Fig. 5: A summary plots of our experiments. In (b) and (c), mean Dice score of AntsSynQuick are also plotted as a horizontal dashed line free of the two parameters.

4 Discussion

We have developed an unsupervised learning algorithm, FAIM, for 3D medical image registration with an option to directly penalize “foldings”, which are spatial locations where the deformation is non-invertible, indicated by a negative determinant of the Jacobian matrix. Our algorithm is similar to the U-net based registration network VoxelMorph of Balakrishnan et al. [2], however our architecture design and loss functions are different. Our anti-folding penalty is similar to (but different from) the penalty used by Zhang et al. [15]. We compared FAIM experimentally to VoxelMorph on the Mindboggle101 dataset [6]. Our experiments showed that FAIM has advantages in several aspects including: fewer trainable parameters, higher registration accuracy as measured by Dice score, and less sacrifice needed when trading off registration accuracy for better invertibility (fewer “foldings”). In fact, as seen in Table 2, FAIM with regularization parameter $\beta = 10^{-2}$ produces deformations that are almost completely

invertible (foldings occurring at only 2 voxels per brain on average) while still having better registration accuracy than the ANTs SyNQuick method. While we recognise that ANTs is capable of producing more accurate registrations with well-chosen hyperparameters, our results suggest that NN methods may now be seriously considered for some applications where geometric and variational methods such as ANTs are currently used.

References

1. Avants, B.B., Tustison, N.J., Song, G., Cook, P.A., Klein, A., Gee, J.C.: A reproducible evaluation of ants similarity metric performance in brain image registration. *Neuroimage* **54**(3), 2033–2044 (2011)
2. Balakrishnan, G., Zhao, A., Sabuncu, M.R., Gutttag, J., Dalca, A.V.: An unsupervised learning model for deformable medical image registration. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 9252–9260 (2018)
3. He, K., Zhang, X., Ren, S., Sun, J.: Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In: *Proceedings of the IEEE international conference on computer vision*. pp. 1026–1034 (2015)
4. Jaderberg, M., Simonyan, K., Zisserman, A., et al.: Spatial transformer networks. In: *Advances in neural information processing systems*. pp. 2017–2025 (2015)
5. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014)
6. Klein, A., Tourville, J.: 101 labeled brain images and a consistent human cortical labeling protocol. *Frontiers in neuroscience* **6**, 171 (2012)
7. Li, H., Fan, Y.: Non-rigid image registration using fully convolutional networks with deep self-supervision. *arXiv preprint arXiv:1709.00799* (2017)
8. Rohé, M.M., Datar, M., Heimann, T., Sermesant, M., Pennec, X.: Svf-net: Learning deformable image registration using shape matching. In: *International Conference on Medical Image Computing and Computer-Assisted Intervention*. pp. 266–274. Springer (2017)
9. Ronneberger, O., Fischer, P., Brox, T.: U-net: Convolutional networks for biomedical image segmentation. In: *International Conference on Medical image computing and computer-assisted intervention*. pp. 234–241. Springer (2015)
10. Shan, S., Guo, X., Yan, W., Chang, E.I., Fan, Y., Xu, Y., et al.: Unsupervised end-to-end learning for deformable medical image registration. *arXiv preprint arXiv:1711.08608* (2017)
11. Sokooti, H., de Vos, B., Berendsen, F., Lelieveldt, B.P., Išgum, I., Staring, M.: Nonrigid image registration using multi-scale 3d convolutional neural networks. In: *International Conference on Medical Image Computing and Computer-Assisted Intervention*. pp. 232–239. Springer (2017)
12. Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., Erhan, D., Vanhoucke, V., Rabinovich, A.: Going deeper with convolutions. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 1–9 (2015)
13. Wang, S., Kim, M., Wu, G., Shen, D.: Scalable high performance image registration framework by unsupervised deep feature representations learning. In: *Deep Learning for Medical Image Analysis*, pp. 245–269. Elsevier (2017)
14. Yang, X., Kwitt, R., Niethammer, M.: Fast predictive image registration. In: *Deep Learning and Data Labeling for Medical Applications*, pp. 48–57. Springer (2016)

15. Zhang, J.: Inverse-consistent deep networks for unsupervised deformable image registration. arXiv preprint arXiv:1809.03443 (2018)