# ASSESSING THE ACCURACY OF NON-RIGID REGISTRATION WITH AND WITHOUT GROUND TRUTH

*Some people*

Some places

## ABSTRACT

We compare two methods for assessing the performance of groupwise non-rigid registration algorithms. One approach, which has been described previously, utilizes a measure of overlap between data labels. Our new approach exploits the fact that, given a set of non-rigidly registered images, a generative statistical appearance model can be constructed. We observe that the quality of the model depends on the quality of the registration, and can be evaluated by comparing synthetic images sampled from the model with the original image set. We derive indices of model specificity and generalisation that can be used to assess model/registration quality. We show that both approaches demonstrate the loss of registration as a set of correctly registered MR images of the brain is progressively perturbed. We compare the sensitivities of the different methods and show that, as well as requiring no ground truth, our new specificity measure provides the most sensitive approach to detecting misregistration.

## 1. INTRODUCTION

Non-rigid registration (NRR) of both pairs and groups of images has been used increasingly in recent years, as a basis for medical image analysis. Applications include structural analysis, atlas matching and change analysis [5]. The problem is highly under-constrained and the plethora of different algorithms that have been proposed generally produce different results for a given set of images [4, 19].

Various methods have been proposed for assessing the results of NRR [9, 11, 16, 15]. Most of these require access to some form of ground truth. One approach involves the construction of artificial test data, which limits application to 'off-line' evaluation. Other methods can be applied directly to real data, but require that anatomical ground truth be provided, typically involving annotation by an expert. This makes validation expensive and prone to subjective error.

We present two methods for assessing the performance of non-rigid registration algorithms; one requires ground truth to be provided *a priori*, whereas the other does not. We compare the two approaches by systematically varying the quality of registration of a set of MR images of the brain.
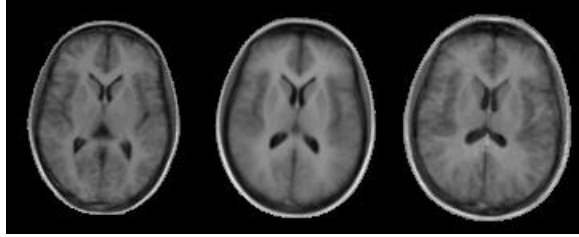
## 2. METHOD

The first of the proposed methods for assessing registration uses a generalisation of Tanimoto's spatial overlap measure [7]. We start with a manual mark-up of each image, providing an anatomical/tissue label for each voxel, and measure the overlap of corresponding labels following registration. We represent each label using a binary image, but after warping and interpolation into a common reference frame, based on the results of NRR, we obtain a set of fuzzy label images. These are used to compute the generalised overlap score [1] as follows:
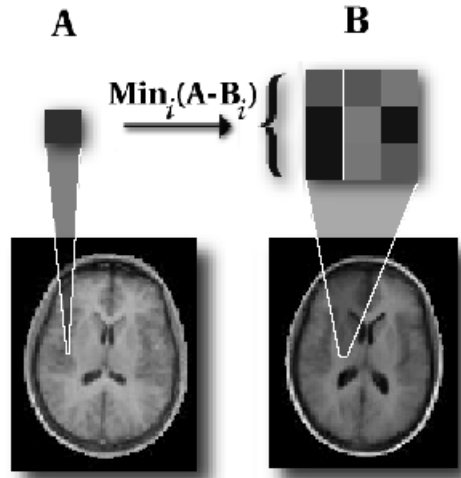
$$O = \frac{\sum\limits_{\text{pairs},k} \sum\limits_{\text{labels},l} \alpha_l \sum\limits_{\text{voxels},i} MIN(A_{kli}, B_{kli})}{\sum\limits_{\text{pairs},k} \sum\limits_{\text{labels},l} \alpha_l \sum\limits_{\text{voxels},i} MAX(A_{kli}, B_{kli})} \quad (1)$$

where $i$ indexes voxels in the registered images, $l$ indexes the label and $k$ indexes image pairs. $A_{kli}$ and $B_{kli}$ represent voxel label values in a pair of registered images and are in the range [0, 1]. The $MIN()$ and $MAX()$ operators are standard results for the intersection and union of a fuzzy set. This generalised overlap measures the consistency with which each set of labels partitions the image volume. The parameter $\alpha_l$ affects the relative weighting of different labels. With $\alpha_l = 1$, label contributions are implicitly volume weighted with respect to one another. We have also considered the cases where $\alpha_l$ weights for the inverse label volume (which makes the relative weighting of different labels equal), where $\alpha_l$ weights for the inverse label volume squared (which gives labels of smaller volume higher weighting) and where $\alpha_l$ weights for a measure of label complexity (which we define arbitrarily as the mean absolute voxel intensity gradient in the label).
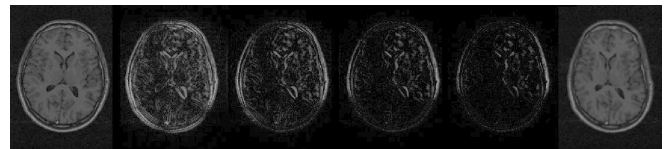
The second method assesses registration in terms of the quality of a generative statistical appearance model, constructed from the registered images. The idea is that a correct registration produces an anatomically meaningful dense correspondence between the set of images, resulting in a better statistical appearance model. We define model quality using two measures – specificity and generalisation [18]. Both are measures of overlap between
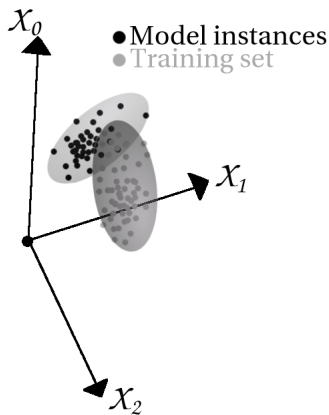
**Fig. 1**. The effect of varying the first model parameter of a brain appearance model by $\pm 2.5$ standard deviations.



**Fig. 2**. Training set and model in hyperspace



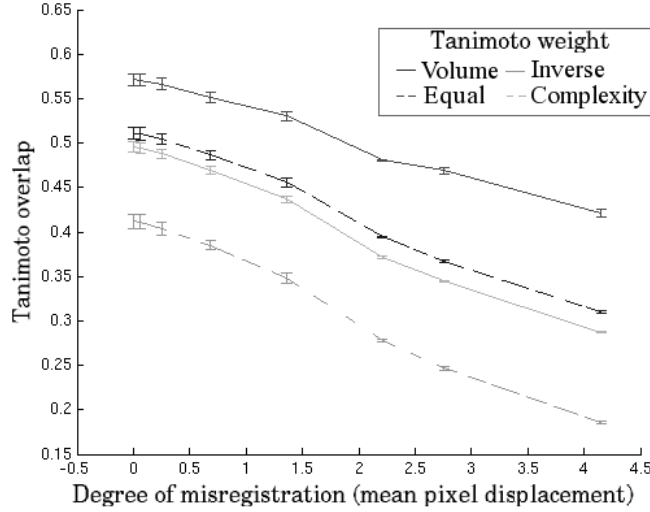**Fig. 3**. The calculation of a shuffle difference image



**Fig. 4**. Shuffle distance evaluation: **Left:** one image, **Right:** another image, **Centre, from left to right:** images showing contributions to shuffle distance, for $r = 0$ (abs. diff.), 1.5, 2.1 & 3.7 respectively.
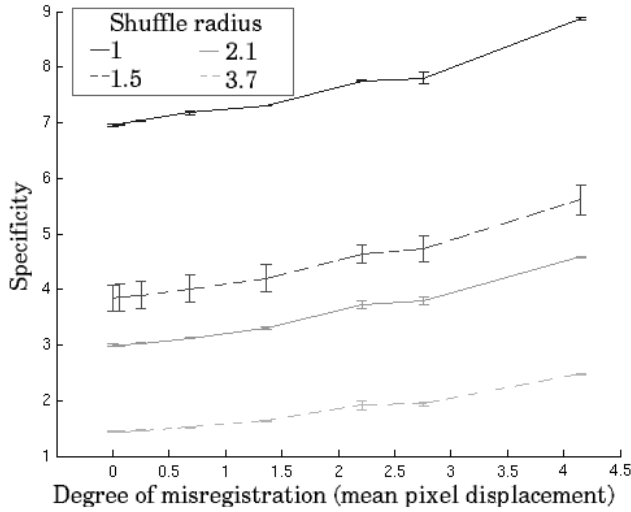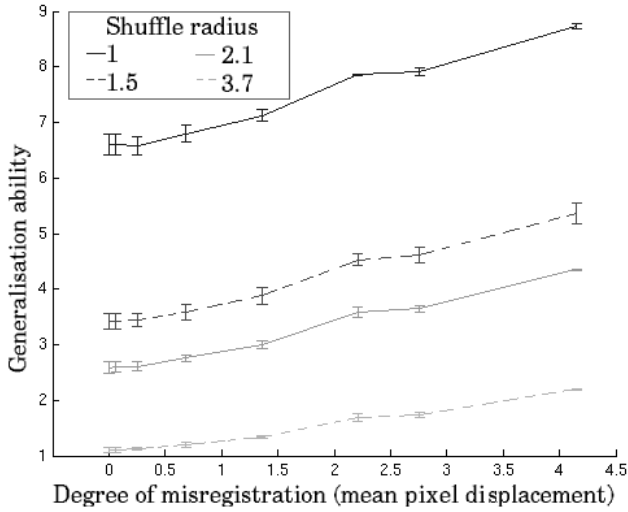
Registration is then evaluated through specificity and generalisation ability [18] of the model, or the ability of the model to i) generate realistic examples of the modelled entity and ii) represent well both seen and unseen examples of the modelled class. In practice, these are evaluated by using generative properties of the model to produce a large number of synthetic examples (in this case brain images) that are then compared to real examples in the original set using some pre-defined image distance measure. Minimum distances of synthetic examples to examples in the original set and vice versa, give model specificity and generalisation respectively. Image distance is measured as a mean shuffle distance, or minimum Euclidean distance between a pixel in one image and a corresponding neighbourhood of pixels in the other.

To test the validity of the proposed methods, the brain images were annotated with 6 tissue classes including gray, white matter and CSF that provided the ground truth for image correspondence. Initially, the images were brought into alignment using an NRR algorithm based on the MDL optimisation. A test set of different registrations was then created by applying random perturbations to each image in the registered set using diffeomorphic clamped-plate splines. By choosing

**Fig. 6**. Appearance model which was built automatically by group-wise registration. First mode is shown, ±2.5 standard deviations.

a different perturbation seed for each image and gradually increasing the magnitude of the perturbations, a series of image sets of progressively worse spatial correspondence and thus registration quality were obtained. By measuring the quality of the registration at each step, the proposed registration assessment measures can be validated.

## 3. RESULTS

Overall, the above approach was applied 10 times using 10 different perturbation seeds to ensure that both methods are consistent and results unbiased. Results of the proposed measures for increasing registration perturbation are shown in Figure 5. Note that Generalisation and Specificity plotted for different shuffle neighbourhood radius are in error form, i.e. they increase with decreasing performance. All metrics are generally well-behaved and show a monotonic decrease in registration performance. Such results directly validate the model-based metrics, which are shown be in agreement with the ground truth embodied in the region overlap based measure.
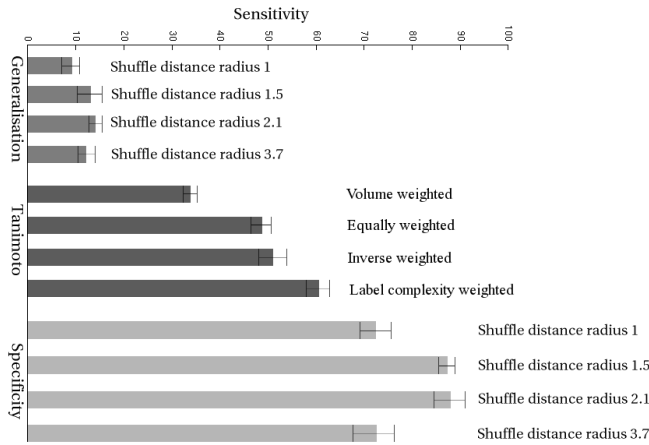
These results also demonstrate that, for all sizes of shuffle neighbourhood, the specificity and generalisation values increase (get worse) with increasing mis-registration. The results for different sizes of shuffle neighbourhood demonstrate that the range of mis-registration over which distinct values of specificity and generalisation are obtained increases as the neighbourhood size increases. We observe similar behaviour as the value of $\alpha_l$ is altered.

Finally, in order to obtain a quantitative comparison of the proposed algorithms we explore sensitivity of the proposed metrics, where the slighter the difference which can be detected reliably, the more sensitive the method. Sensitivity is in this case defined as the rate of change in the measure for a given perturbation range, normalised by the average uncertainty in the measurement over that range. More formally, sensitivity can be defined thus:

$$\frac{m - m_0}{d} / \overline{\sigma} \qquad (2)$$

where $m$ is the quality measured for a given value of displacement, $m_0$ is the measured quality at registration, $d$ is the degree of deformation and $\overline{\sigma}$ is the mean over the error

**Fig. 5**. Behaviour of proposed metrics with increasing registration perturbation: a) Generalisation, b) Specificity and c) Tantimoto overlap

bars. Sensitivity is evaluated for all three of the proposed metrics and shown in Figure 3 with errors bars based on both an inter-instantiation error and a measure-specific error. The Specificity measure is the most sensitive for any radius of the shuffle distance followed by the overlap metric and Generalisation, with shuffle radii of 1.5 and 2.1 (equivalent to 3x3 and 5x5 neighbourhoods) giving optimal sensitivity.



**Figure 2.** The sensitivity of the different registration assessment methods.

## 4. CONCLUSIONS

We have introduced a model-based approach to assessing the accuracy of non-rigid registration, without the need for ground truth. The validation experiments, based on perturbing correspondences obtained using ground truth, show that we are able to detect increasing mis-registration using just the registered image data. The results obtained for different sizes of shuffle neighbourhood show that the use of shuffle distance rather than Euclidean distance improves the range of mis-registration over which we can detect significant changes in registration accuracy.

More broadly, registration performance can be evaluated reliably both in the cases when ground truth information is available and when it is not. In particular, the methods based on generative statistical model evaluation are shown to be in agreement with the ground truth expressed through the true image region overlap metric based on the Tantimoto formulation. Proposed metrics are also shown to have sufficient sensitivity to detect very subtle changes in registration performance, on the level of perturbations measured in fractions of a pixel.

We believe that this represents an important advance in the assessment of NRR, because it establishes an entirely objective basis for evaluating the reliability of NRR-based experiments, and for comparing the performance of different methods of NRR. The fact that no ground truth data is required means that the method can be applied routinely. Further work is needed to compare the results obtained using our new approach with those obtained using more sophisticated segmentation-based methods of evaluation.

## 5. ACKNOWLEDGEMENT

## 6. REFERENCES

[1] W. R. Crum, O. Camara, D. Rueckert, K. Bhatia, M. Jenkinson, and D. L. G. Hill. Generalised overlap measures for assessment of pairwise and groupwise image registration and segmentation. In *Proceedings of MICCAI*, 3749:99-106, 2005.

[2] T. F. Cootes, C. Beeston, G. J. Edwards, and C. J. Taylor. A unified framework for atlas matching using active appearance models. In *Information Processing in Medical Imaging*, 1613:322-333, 1999.

[3] T.F. Cootes, G.J. Edwards and C.J.Taylor. Active appearance models. In *European Conference on Computer Vision*, 2:484-498, 1998.

[4] T. F. Cootes, S. Marsland, C. J. Twining, K. Smith, and C. J. Taylor. Groupwise diffeomorphic non-rigid registration for automatic model building. In *European Conference on Computer Vision*, 2034:316-27, 2004.

[5] W. R. Crum, T. Hartkens, and D. L. G. Hill. Non-rigid image registration: theory and practice. *British Journal of Radiology*, 77:140-153, 2004.

[6] R. H. Davies, C. J. Twining, T. F. Cootes, J. C. Waterton, and C. J. Taylor. A minimum description length approach to statistical shape modeling. *IEEE Transactions on Medical Imaging*, 21(5):525-537, 2002.

[7] M. Beauchemin and K. P. B. Thomson. The evaluation of segmentation results and the overlapping area matrix. *International Journal of Remote Sensing*, 18(18):3895-3899, 1997.

[8] G. J. Edwards, T. F. Cootes, and C. J. Taylor. Face recognition using active appearance models. In *European Conference on Computer Vision*, 2:581-595, 1998.

[9] J. M. Fitzpatrick and J. B. West. The distribution of target registration error in rigid-body point-based registration. *IEEE Transaction Medical Imaging,* 20:917-27, 2001.

[10] A. F. Frangi, D. Rueckert, J. A. Schnabel, and W. J. Niessen. Automatic construction of multiple-object three-dimensional statistical shape models: application to cardiac modelling. *IEEE Transactions on Medical Imaging*, 21:1151-66, 2002.

[11] P. Hellier, C. Barillot, I. Corouge, B. Giraud, G. Le Goualher, L. Collins, A. Evans, G. Malandain, and N. Ayache. Retrospective evaluation of inter-subject brain registration. In *Medical Image Computing and Computer-Assisted Intervention*, 2208:258-265, 2001.

[12] K. N. Kutulakos. Approximate N-view stereo. In *European Conference on Computer Vision*, 1:67-83, 2000.

[13] D. Rueckert, A. F. Frangi, and J. A. Schnabel. Automatic construction of 3-D statistical deformation models of the brain using nonrigid registration. *IEEE Transactions on Medical Imaging*, 22(8)1014-1025, 2003.

[14] D. Rueckert, L. I. Sonoda, C. Hayes, D. L. G. Hill, M. O. Leach, D. J. Hawkes. Non-rigid registration using free-form deformations: application to breast MR images. *IEEE Transactions on Medical Imaging*, 18(8):712-721, 1999.

[15] P. Rogelj, S. Kovacic, and J. C. Gee. Validation of a nonrigid registration algorithm for multimodal data. *Medical Imaging*, volume 4684, 2002.

[16] J. A. Schnabel, C. Tanner, A. Castellano-Smith, M. O. Leach, C. Hayes, A. Degenhard, R Hose, D. L. G. Hill, and D. J. Hawkes. Validation of non-rigid registration using finite element methods. In *Information Processing in Medical Imaging,* 2082:344-357, 2001.

[17] M. B. Stegmann, B. K. Ersboll, and R. Larsen. FAME - a flexible appearance modeling environment. *IEEE Transactions on Medical Imaging*, 22(10):1319-1331, 2003.

[18] C. J. Twining, T.F. Cootes, S. Marsland, S. V. Petrovic, R. S. Schestowitz, and C. J. Taylor. A unified information-theoretic approach to groupwise non-rigid registration and model building. In *Information Processing in Medical Imaging*, 3565:1-14, 2005.

[19] B. Zitova and J. Flusser. Image registration methods: a survey. *Image Vision Computing*, 21:977-1000, 2003.