# A Probabilistic, Non-parametric Framework for Inter-modality Label Fusion

Juan Eugenio Iglesias[1], Mert Rory Sabuncu[1] and Koen Van Leemput[1,2,3]

[1] Martinos Center for Biomedical Imaging, MGH, Harvard Medical School, USA
[2] Department of Applied Mathematics and Computer Science, DTU, Denmark
[3] Departments of Information and Computer Science and of Biomedical Engineering and Computational Science, Aalto University, Finland

**Abstract.** Multi-atlas techniques are commonplace in medical image segmentation due to their high performance and ease of implementation. Locally weighting the contributions from the different atlases in the label fusion process can improve the quality of the segmentation. However, how to define these weights in a principled way in inter-modality scenarios remains an open problem. Here we propose a label fusion scheme that does not require voxel intensity consistency between the atlases and the target image to segment. The method is based on a generative model of image data in which each intensity in the atlases has an associated conditional distribution of corresponding intensities in the target. The segmentation is computed using variational expectation maximization (VEM) in a Bayesian framework. The method was evaluated with a dataset of eight proton density weighted brain MRI scans with nine labeled structures of interest. The results show that the algorithm outperforms majority voting and a recently published inter-modality label fusion algorithm.

## 1 Introduction

Automated segmentation of brain MRI scans is a key step in most neuroimaging pipelines. Manual delineation of structures of interest is time consuming and rater dependent, making automated approaches desirable. Some of the most popular segmentation methods in the recent literature are based on the multi-atlas paradigm, in which a set of training images with manual annotations (henceforth "atlases") are deformed to the image to analyze. The deformations are used to propagate the annotations to *target* space, where they are finally merged into an estimate of the segmentation; this step is known as label fusion. Multi-atlas techniques overcome the main limitation of using a single atlas: the fact that a single template can seldom cover all the anatomical variability within a population.

The simplest form of label fusion is majority voting, in which the most frequent label is assigned to each voxel [1]. Better results can be achieved by locally weighting the contribution of the atlases by their similarity to the target scan after registration. In [2], Isgum et al. define weights by inverting the absolute difference in image intensities. A more principled framework based on a generative model was proposed by Sabuncu et al. [3]: a smooth, discrete, latent membership

field determines from which deformed atlas the intensity and label are taken at each spatial location; Bayesian inference is used to compute the segmentation.

These methods rely on the consistency of voxel intensities across the atlases and the target scan. This assumption falters in inter-modality scenarios, including MRI when the atlases and the target have been acquired with different hardware or pulse sequences. This is often the case when analyzing clinical or legacy data. Another example of application that could benefit from inter-modality label fusion is the analysis of infant brain MRI, in which the intensities are very different than in scans from adults due to ongoing myelination.

While the image registration literature has managed inter-modality scenarios using metrics based on mutual information [4], label fusion across modalities remains an open problem. One possible approach is to arbitrarily define weights based on the mutual information or cross-correlation computed in a window around each voxel [5], yet the optimality of such an approach is unclear.

A principled way of carrying out label fusion across modalities was presented in [6]. This method is based on a generative model in which the intensity of the voxels corresponding to each label follows a Gaussian distribution. The parameters of the Gaussian are estimated from the data, making the fusion robust against changes in modality or MRI contrast. While this approach outperforms heuristic schemes based on local cross-correlations [7], it still has two disadvantages. First, the performance is poor when the Gaussian assumption is violated, such as in the thalamus or the putamen in brain MRI scans. And second, since the fusion only considers the deformed labels from the atlases, it ignores potentially valuable information from their intensities.

Here we propose a generalization of Sabuncu et al.'s model to inter-modality scenarios. The generative model is essentially the same; however, we do not hypothesize a Gaussian relation between the intensities of the atlases and the target scan. Instead, we assume a more flexible model based on a semi-parametric or non-parametric conditional distribution of the intensities of the target given the intensities of the atlases. Using a Bayesian framework, we first estimate this conditional distribution and also a multiplicative bias field. Then, the segmentation is computed as the most likely labels given these estimates and the input image.

## 2    Methods

### 2.1    Generative Model

The proposed generative model is shown in Fig. 1a, and the corresponding equations in Fig. 1b. We assume that registration is a preprocessing step: the intensities $\{I_n\}$ and corresponding labels $\{L_n\}$ of the $N_{atl}$ deformed atlases are constant during the fusion. $M(\mathbf{x}) \in \{1, \ldots, N_{atl}\}$ is a discrete, latent membership field that indexes which atlas generated the label and intensity of the voxel at spatial location $\mathbf{x}$. $M(\mathbf{x})$ is smooth thanks to a Markov random field (MRF) prior (Eq. 1 in Fig. 1b, where $\mathcal{V}_x$ is the 6-neighborhood of $\mathbf{x}$ and $\delta(\cdot)$ is Kronecker's delta).

Given $M(\mathbf{x})$, the label of a voxel $L(\mathbf{x})$ is sampled from a categorical distribution given by a logOdds model [8] defined by the warped labels of atlas $M(\mathbf{x})$

(Eq. 2 in Fig. 1b, where $\rho$ is the slope of the model and $D_n^l$ is the signed distance transform for atlas $n$ and label $l \in \{1, \ldots, \mathcal{L}\}$). This prior has been shown to outperform taking the label $L_{M(\mathbf{x})}(\mathbf{x})$ directly [3]. Given $L(\mathbf{x})$, the uncorrupted (i.e., bias field corrected) intensity of the voxel $\tilde{I}(\mathbf{x})$ is sampled from the conditional distribution $p(\tilde{I}(\mathbf{x})|I_{M(\mathbf{x})}(\mathbf{x}), \Theta)$, which depends on the model parameters $\Theta$ (Eq. 3 in Fig. 1b). Here we have made two assumptions. First, that the conditional PDF is stationary in space. Second, we assume that the image intensities are consistent across atlases, which is fair when the atlases are from the same modality. This allows us to represent the conditional probability density function (PDF) with a single, atlas-independent PDF.

For the conditional PDF of intensities, we consider two different models. First, a semi-parametric (SP) model which describes $\tilde{I}(\mathbf{x})$ with a Gaussian PDF with mean and variance depending on $I_{M(\mathbf{x})}(\mathbf{x})$. And second, a non-parametric (NP) model based on a collection of conditional 1D histograms ($h$) of $\tilde{I}(\mathbf{x})$ depending on $I_{M(\mathbf{x})}(\mathbf{x})$. These two forms of the algorithm allow us to isolate and understand the effects of the two contributions of this model with respect to [6]: the departure from the Gaussian model (NP) and using the intensities of the atlases in the fusion (NP and SP). The models are given by:

$$\text{Semi-parametric (SP): } p(\tilde{I}(\mathbf{x})|I_{M(\mathbf{x})}, \Theta) = \mathcal{N}(\tilde{I}(\mathbf{x}); \mu_{Q[I_{M(\mathbf{x})}]}, \sigma^2_{Q[I_{M(\mathbf{x})}]}). \quad (1)$$

$$\text{Non-parametric (NP): } p(\tilde{I}(\mathbf{x})|I_{M(\mathbf{x})}(\mathbf{x}), \Theta) = h(Q[\tilde{I}(\mathbf{x})]; Q[I_{M(\mathbf{x})}(\mathbf{x})]), \quad (2)$$

where $Q[\cdot]$ is a nearest neighbor interpolator that quantizes the intensities of the target scan and the atlases into discrete sets $A$ and $B$, respectively.

Finally, $\tilde{I}(\mathbf{x})$ is corrupted by a low-frequency, multiplicative bias field to yield the observed intensities $I(\mathbf{x})$ (Eq. 4 in Fig. 1). The bias field is modeled as the exponential (to ensure non-negativity) of a linear combination of smooth basis functions $\{\psi_k\}$. The linear coefficients $\mathbf{b} = \{b_k\}$ are grouped with the parameters of the conditional PDF of intensities into the model parameters $\theta = (\mathbf{b}, \{h(a; b)\})$ (NP) or $\theta = (\mathbf{b}, \{\mu_b, \sigma_b^2\})$ (SP). A flat prior distribution $p(\theta) \propto 1$ completes the model. Note that the denominator in Eq. 5 in Fig. 1b ensures integration to one.
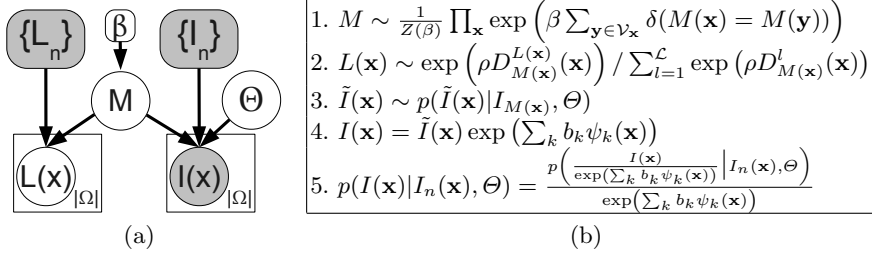
## 2.2 Inference

The segmentation $\hat{L}$ of image $I$ is estimated by maximizing the posterior probability: $p(L|I, \{I_n\}, \{L_n\})$. This leads to an intractable integral for which we can use the approximation that the posterior distribution of the model parameters is a Dirac's delta, i.e., $p(\theta|I, \{I_n\}) \approx \delta(\theta = \hat{\theta})$:

$$\hat{L} = \operatorname*{argmax}_L \int p(L|\theta, I, \{I_n\}, \{L_n\}) p(\theta|I, \{I_n\}) d\theta \approx \operatorname*{argmax}_L p(L|\hat{\theta}, I, \{I_n\}, \{L_n\}), \quad (3)$$

Thus, we first need to estimate $\hat{\theta}$ and then use it to compute $\hat{L}$ with Eq. 3.

**Estimating $\hat{\theta}$:** the problem here is

$$\hat{\theta} = \operatorname*{argmax}_\theta \; p(\theta|I, \{I_n\}) = \operatorname*{argmax}_\theta \; \log p(I|\{I_n\}, \theta). \quad (4)$$

Fig. 1. Graphical model (a) and corresponding equations (b). Random variables are in circles, constants in boxes, observed variables shaded and plates indicate replication.

Note that $\theta$ does not depend on $\{L_n\}$. Optimizing Eq. 4 requires marginalizing over $M$, which is intractable due to the MRF. Therefore, we use VEM to compute an approximate solution by optimizing a lower bound J instead:

$$\log p(I|\{I_n\}, \theta) \geq J(q(M), \theta) = \log p(I|\{I_n\}, \theta) - KL[q(M)||p(M|I, \theta, \{I_n\})] \quad (5)$$

$$= H[q] + \sum_M q(M) \log p(M, I|\theta, \{I_n\}), \quad (6)$$

where $H[\cdot]$ is the entropy of a random variable, $KL$ represents the (non-negative) Kullback-Leibler divergence and $q$ is a distribution over $M$ which is restricted to having a simpler form than $p(M|I, \theta, \{I_n\})$. We alternately optimize $J$ with respect to $q$ (E step) and $\Theta$ (M step).

In the **E step**, we work with Eq. 5. Maximizing $J$ amounts to minimizing the KL divergence. The standard "mean field" approximation assumes that $q$ factorizes as $q(M) = \prod_x q_x(M(\mathbf{x}))$, where $q_x(m)$ is a categorical distribution over the atlas indices $m = 1, \ldots, N_{atl}$, and $\sum_m q_x(m) = 1$. This yields:

$$\underset{q}{\operatorname{argmin}} \sum_{\mathbf{x}} \sum_m q_x(m) \log \frac{q_x(m)}{p(I(\mathbf{x})|I_m(\mathbf{x}), \theta)} - \beta \sum_{\mathbf{x}} E_{q_x} \left[ \sum_{y \in \mathcal{V}_x} q_y(M(\mathbf{x})) \right]. \quad (7)$$

Building the Lagrangian and setting derivatives to zero gives:

$$q_x(m) \propto p(I(x)|I_m(x), \theta) \exp[\beta \sum_{y \in \mathcal{V}_x} q_y(m)], \quad (8)$$

such that $\sum_m q_x(m) = 1$. We can solve this equation with fixed point iterations.

In the **M step**, it is more convenient to work with Eq. 6: since we are optimizing for $\theta$, we can disregard $H(q)$. Because of the structure of $q$, we have:

$$\underset{\Theta}{\operatorname{argmax}} \sum_{\mathbf{x}} \sum_m q_x(m) \left[ \log p(\tilde{I}(\mathbf{x})|I_m(\mathbf{x}), \Theta) - \sum_k b_k \psi_k(\mathbf{x}) \right], \quad (9)$$

with $\tilde{I}(\mathbf{x}) = I(\mathbf{x})e^{-\sum_k b_k \psi_k(\mathbf{x})}$. The solution depends on whether we consider the SP or the NP model. In the first case, replacing $p(\tilde{I}(\mathbf{x})|I_m(\mathbf{x}), \Theta)$ by the

corresponding Gaussian (Eq. 1), taking derivatives with respect to the means and variances, and setting them to zero gives:

$$\mu_b = \sum_{\mathbf{x}} \sum_{m} w_x^b(m) I(\mathbf{x}) e^{-\sum_k b_k \psi_k(\mathbf{x})} \Big/ \sum_{\mathbf{x}} \sum_{m} w_x^b(m), \qquad (10)$$

$$\sigma_b^2 = \sum_{\mathbf{x}} \sum_{m} w_x^b(m) (I(\mathbf{x}) e^{-\sum_k b_k \psi_k(\mathbf{x})} - \mu_b)^2 \Big/ \sum_{\mathbf{x}} \sum_{m} w_x^b(m), \qquad (11)$$

where $w_x^b(m) = q_x(m) \delta (Q[I_m(\mathbf{x})] = b)$. Eqns. 10 and 11 are weighted means and variances depending on the (approximate) membership posteriors $q_x(m)$.

For the NP model, we substitute Eq. 2 into Eq. 9 and build a Lagrangian to ensure integration to one. Taking derivatives and setting them to zero yields:

$$h(a; b) \propto \sum_{\mathbf{x}} \sum_{m} w_x^b(m) \delta \left( Q[I(\mathbf{x}) e^{-\sum_k b_k \psi_k(\mathbf{x})}] = a \right). \qquad (12)$$

such that $\sum_{a \in A} h(a; b) = 1/\Delta$, where $\Delta$ is the quantization interval. Again, Eq. 12 is simply a weighted histogram.

Finally, we use a quasi Newton solver with an explicit line search to optimize Eq. 9 for the bias field coefficients, both in the SP and BP case.

**Computing the final segmentation:** given $\theta$, estimating the final segmentation with Eq. 3 still requires an intractable sum over $M$. However, since $q(M)$ minimizes the KL divergence with $p(M|\hat{\theta}, I, \{I_n\})$, we approximate:

$$\hat{L} = \arg\max_L \sum_M p(L|M, \{L_n\}) p(M|\hat{\theta}, I, \{I_n\}) \approx \arg\max_L \sum_M p(L|M, \{L_n\}) q(M)$$

$$= \arg\max_L \prod_{\mathbf{x}} \sum_m q_x(m) p(L(\mathbf{x})|L_m) \Rightarrow \hat{L}(\mathbf{x}) = \arg\max_{L(\mathbf{x})} \sum_m q_x(m) p(L(\mathbf{x})|L_m) \quad (13)$$

**Summary of the algorithm:** we initialize the bias field coefficients $b_k = 0$, and the distribution $q_x(m) = 1/N_{atl}$. Next, we alternate the E and M steps until convergence. The E step updates $q$ with fixed point iterations of Eq. 8. The M step first updates the bias field by numerically optimizing Eq. 9 with respect to $\{b_k\}$, and then the parameters of the conditional PDF with Eqns. 10, 11 (SP model) or Eq. 12 (NP). Upon convergence, the final segmentation is computed with Eq. 13. The method is illustrated with a simple example in Fig. 2.

## 3   Experiments and results

We used 39 manually delineated (see protocol in [9]) T1 MRI scans as atlases to segment 36 brain structures in eight proton-density (PD) scans. The annotations of the PD scans were made on co-registered T1 data, allowing consistent annotations across the two datasets. FreeSurfer [10] was used to skull-strip the volumes and intensity-normalize the atlases, since consistent intensities are assumed.
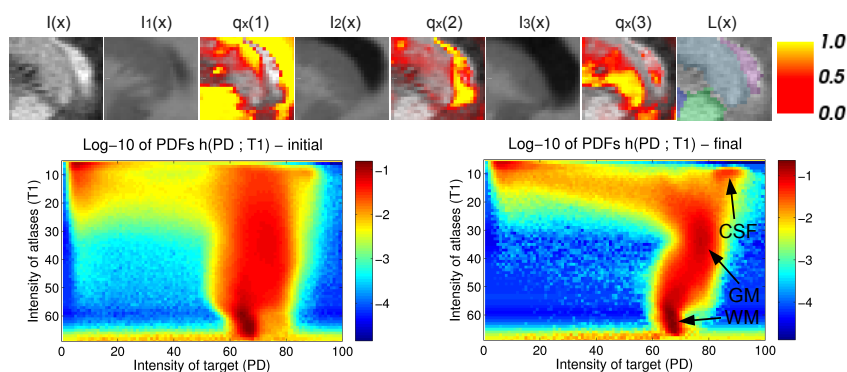
We used Elastix [11] to register the T1 to the PD scans (b-spline transform, mutual information). We compared four label fusion methods: majority voting, the Gaussian model from [6] and the proposed approach (SP and NP). We set $\rho = 1$, $\beta = 0.75$ (as in [6]), $\{\psi_k\}$ to a fourth order polynomial and the number of bins $|A| = |B| = 64$. The Dice overlap between the manual and automatically generated labels was used as measure of performance. Statistical significance was assessed with paired Wilcoxon signed rank tests. For simpler presentation of results, we merged the label of each left structure with its right counterpart, and used a representative subset of structures in the evaluation (as in [3]): white matter (WM), cortex (CT), lateral ventricle (LV), thalamus (TH), caudate (CA), putamen (PT), pallidum (PA), hippocampus (HP) and amygdala (AM).

Box plots of the structure-wise Dice scores are shown in Fig. 3a, whereas segmentations for a sample axial slice are shown in Fig. 3b. The p-values for the statistical tests comparing the NP method (the top-performing algorithm) with the other competing approaches are shown in Tab. 1. Majority voting produces decent outputs for the subcortical structures, but fails to extract the convoluted white matter surface, which is very difficult to register (see Fig. 3b-iii). It also produces bad results for the ventricles, as illustrated in the same figure. The Gaussian model gives excellent results for the cortex, but falters when the normality assumption does not hold. This is often the case for the thalamus and the putamen. For instance, the latter leaks into the white matter in Fig. 3b-iv.
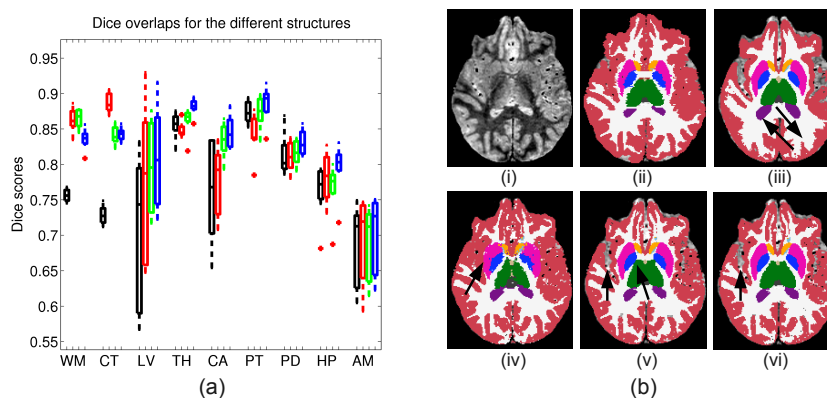
The NP version of the proposed approach significantly outperforms majority voting for every brain structure (Tab. 1). It also yields Dice scores significantly higher than those from the Gaussian model for all subcortical structures. Only in the cortex and the white matter the performance is inferior; see for instance the mistake marked by the arrow in Fig. 3b-vi. This is because the registration is poor for these convoluted structures, making a simple Gaussian intensity model more suitable. Overall, the mean improvement in Dice score is $\sim 2\%$ over the Gaussian model and $\sim 5\%$ over majority voting. The SP version also beats majority voting and the Gaussian model. However, it performs slightly worse than the NP method, likely due to its inability to describe multimodal shapes in the conditional intensity PDF (e.g., see atlas intensity range 10-20 in Fig. 2).

## 4    Discussion

We presented a cross-modality label fusion method based on a generative model that describes the relationship between image intensities in a SP or NP manner. The algorithm often converges in less than 15 iterations (about 20 minutes on a modern PC). The results show that using the intensities of the atlases in the fusion allows the SP algorithm to outperform previously proposed inter-modality label fusion techniques. Moreover, departure from the Gaussian model allows the NP model to further improve the results. Exploring more flexible SP models (such as those based on mixtures of Gaussians), incorporating the registration step into the framework and using more accurate approximations than nearest neighbors when estimating the conditional intensity PDF remain as future work.

**Fig. 2.** Intermediate outputs for a sagittal patch of a PD scan segmented with three T1 atlases and the NP model. Top row: $I$, $\{I_n\}$, $q_x(m)$ overlaid on $I$, and segmentation $L$ (green = putamen, purple = ventricle, blue = caudate); $q_x(m)$ highlights where each atlas contributed to generating $I, L$. Bottom row: initial and final estimates of the conditional PDFs $h(a; b)$; the latter are much sharper and emphasize three modes corresponding to gray matter (GM), white matter (WM) and cerebrospinal fluid (CSF).



**Fig. 3.** (a) Box plot of Dice scores for majority voting (black), the Gaussian model [6] (red) and the proposed SP (green) and NP (blue) methods. Lines are at the three quartile values, whiskers extend to 1.5 times the interquartile range from the box, and dots mark outliers. (b-i) Axial slice of a PD scan, (ii) manual labels, (iii) segmentation from majority voting, (iv) the Gaussian model, (v) the SP model, (vi) the NP model. Arrows point at mistakes. The color code is: red = CT, white = WM, pink = PT, dark blue = PD, light blue = CA, green = TH, purple = LV, orange = accumbens.

## Acknowledgements

| Method | WM | CT | LV | TH | CA | PT | PD | HP | AM | All combined |
|---|---|---|---|---|---|---|---|---|---|---|
| Maj. Vot. | 75.6 | 72.8 | 70.7 | 85.5 | 76.2 | 87.4 | 81.2 | 76.3 | 68.6 | 77.2 |
| p-value | 0.008 | 0.008 | 0.008 | 0.008 | 0.016 | 0.078 | 0.023 | 0.008 | 0.008 | $6.8\ 10^{-13}$ |
| Gaussian | 86.3 | **88.7** | 77.3 | 84.9 | 77.7 | 84.2 | 81.1 | 78.2 | 69.3 | 80.8 |
| p-value | 0.008 | 0.008 | 0.11 | 0.008 | 0.008 | 0.008 | 0.016 | 0.078 | 0.008 | $8.2\ 10^{-4}$ |
| SP | **86.5** | 84.1 | 79.6 | 86.2 | 83.4 | 87.4 | 81.7 | 76.6 | 68.9 | 81.6 |
| p-value | 0.008 | 0.46 | 0.016 | 0.008 | 0.055 | 0.016 | 0.023 | 0.008 | 0.008 | $2.8\ 10^{-5}$ |
| NP | 83.5 | 84.3 | **80.9** | **88.1** | **84.6** | **88.6** | **83.1** | **79.6** | **70.1** | **82.5** |

**Table 1.** Mean Dice scores (in %, highest in bold) and p-values with respect to NP.

# References

1. Heckemann, R., Hajnal, J., Aljabar, P., Rueckert, D., Hammers, A.: Automatic anatomical brain mri segmentation combining label propagation and decision fusion. NeuroImage **33**(1) (2006) 115–126
2. Isgum, I., Staring, M., Rutten, A., Prokop, M., Viergever, M., van Ginneken, B.: Multi-atlas-based segmentation with local decision fusionapplication to cardiac and aortic segmentation in CT scans. IEEE Trans. Med. Im. **28**(7) (2009) 1000–1010
3. Sabuncu, M., Yeo, B., Van Leemput, K., Fischl, B., Golland, P.: A generative model for image segmentation based on label fusion. IEEE Trans. Med. Im. **29** (2010) 1714–1729
4. Maes, F., Collignon, A., Vandermeulen, D., Marchal, G., Suetens, P.: Multimodality image registration by maximization of mutual information. IEEE Trans. Med. Im. **16**(2) (1997) 187–198
5. Yushkevich, P.A., Wang, H., Pluta, J., Das, S.R., Craige, C., Avants, B.B., Weiner, M.W., Mueller, S.: Nearly automatic segmentation of hippocampal subfields in in vivo focal T2-weighted MRI. NeuroImage **53**(4) (2010) 1208–1224
6. Iglesias, J., Sabuncu, M., Van Leemput, K.: A generative model for multi-atlas segmentation across modalities. In: 9th IEEE International Symposium on Biomedical Imaging (ISBI). (2012) 888–891
7. Iglesias, J., Sabuncu, M., Van Leemput, K.: A generative model for probabilistic label fusion of multimodal data. In Yap, P.T., Liu, T., Shen, D., Westin, C.F., Shen, L., eds.: Multimodal Brain Image Analysis. Volume 7509 of Lecture Notes in Computer Science. Springer Berlin Heidelberg (2012) 115–133
8. Pohl, K., Fisher, J., Shenton, M., McCarley, R., Grimson, W., Kikinis, R., Wells, W.: Logarithm odds maps for shape representation. In Larsen, R., Nielsen, M., Sporring, J., eds.: Medical Image Computing and Computer-Assisted Intervention MICCAI 2006. Volume 4191 of Lecture Notes in Computer Science. Springer Berlin Heidelberg (2006) 955–963
9. Caviness Jr, V., Filipek, P., Kennedy, D.: Magnetic resonance technology in human brain science: blueprint for a program based upon morphometry. Brain Dev. **11**(1) (1989) 1–13
10. Fischl, B., Salat, D., Busa, E., Albert, M., Dieterich, M., Haselgrove, C., van der Kouwe, A., Killiany, R., Kennedy, D., Klaveness, S., Montillo, A., Makris, N., Rosen, B., Dale, A.: Whole brain segmentation: Automated labeling of neuroanatomical structures in the human brain. Neuron **33** (2002) 341–355
11. Klein, S., Staring, M., Murphy, K., Viergever, M., Pluim, J.: Elastix: a toolbox for intensity-based medical image registration. IEEE Trans. Med. Im. **29**(1) (2010) 196–205