# Tissue Classification of Large-scale Multi-site MR Data Using Fuzzy k-Nearest Neighbor Method

Ali Ghayoor[a,e], Jane S. Paulsen[b,c], Regina E. Y. Kim[b], Hans J. Johnson[a,b,d,e],

[a]Dept. of Electrical and Computer Engineering, The University of Iowa, Iowa City, IA, USA
[b]Dept. of Psychiatry, University of Iowa Hospitals & Clinics, Iowa City, IA USA
[c]Dept. of Neurology, University of Iowa Hospitals & Clinics, Iowa City, IA, USA
[d]Iowa Informatics Institute, The University of Iowa, Iowa City, IA, USA
[e]Iowa Institute for Biomedical Imaging, The University of Iowa, Iowa City, IA, USA

## ABSTRACT

This paper describes enhancements to automate classification of brain tissues for multi-site degenerative magnetic resonance imaging (MRI) data analysis. Processing of large collections of MR images is a key research technique to advance our understanding of the human brain. Previous studies have developed a robust multi-modal tool for automated tissue classification of large-scale data based on expectation maximization (EM) method initialized by group-wise prior probability distributions. This work aims to augment the EM-based classification using a non-parametric fuzzy k-Nearest Neighbor (k-NN) classifier that can model the unique anatomical states of each subject in the study of degenerative diseases. The presented method is applicable to multi-center heterogeneous data analysis and is quantitatively validated on a set of 18 synthetic multi-modal MR datasets having six different levels of noise and three degrees of bias-field provided with known ground truth. Dice index and average Hausdorff distance are used to compare the accuracy and robustness of the proposed method to a state-of-the-art classification method implemented based on EM algorithm. Both evaluation measurements show that presented enhancements produce superior results as compared to the EM only classification.

**Keywords:** Tissue classification, segmentation, expectation maximization, fuzzy k-nearest neighborhood method, multi-site studies, neurodegenerative diseases.

## 1. INTRODUCTION

Brain tissue segmentation on structural magnetic resonance imaging (MRI) has received considerable attention; one of the classic neuroimaging challenges is the segmentation of MR images into white matter (WM), grey matter (GM) and cerebrospinal fluid (CSF). Volumetric measurements in different brain regions are important in studies on aging and neurodegenerative disorders[1] like Alzheimer's disease, Schizophrenia and Huntington's Disease (HD).

Given the relevance of brain tissue segmentation, different automated segmentation methods have been proposed over the years. Almost all of these methods rely on a supervised or unsupervised voxel classifier. Supervised methods use manually segmented training data to learn the typical distribution of intensity or appearance features for the tissue classes.[2] Unsupervised methods, particularly those based on expectation maximization (EM), do not require training data and are therefore more widely used than the supervised methods. EM-based methods start with an initial segmentation, which is often based on a probabilistic brain tissue atlas that is registered to the unlabeled target scans, and from this initialization, class-specific Gaussian intensity distributions are estimated. This intensity model can then be used to update the segmentation and this process is repeated until the segmentation converges.[1]

Kim and Johnson[3] implemented an iterative optimization framework between bias-correction, registration, and tissue classification using expectation maximization (EM) method for large-scale heterogeneous multi-site longitudinal MR data analysis. In this study, we propose to extend the $EM$-based classification using a non-parametric fuzzy k-Nearest Neighbor (k-NN) classifier that avoids biases inherent in $EM$ use of prior probability distributions that may not represent diseased anatomical states.

---

Send correspondence to Hans J. Johnson (hans-johnson@uiowa.edu)

## 2. METHODS

### 2.1 General Framework

This paper describes the algorithmic enhancements on the implementation of a framework developed by Kim and Johnson[3] that iteratively incorporates bias-field correction, image registration, and tissue classification. Enhancements applied for more accurate subject specific tissue classification in processing of heterogeneous multi-site degenerative MR data.

Our atlas based framework takes inputs of any combination of modalities with any number of scan repetitions if the input modalities have comparable resolution and voxel sizes. First, using a Rigid-type registration, all intra-subject scans are spatially normalized into a common subject-specific reference orientation defined by anterior commissure (AC), and posterior commissure (PC) landmarks, and mid-saggital plane.[4] Then, all intra-modal scan repetitions are averaged together to increase the signal-to-noise ratio for each modality. After that, all the atlas priors are placed into the subject space using an atlas to subject transformation that is derived from a high-deformable registration algorithm (SyN)[5,6] to enhance the accuracy of the subject-specific tissue priors.

The warped priors in the subject space are tissue probability maps giving the probability of a certain voxel belonging to a certain tissue, and they are used to initialize the Gaussian distribution parameters for the $EM$ algorithm. Finally, the processes of posterior estimation, bias field correction, and the registration are iteratively updated multiple times until convergence.

### 2.2 New classifier

A non-parametric subject specific fuzzy $k - NN$ classifier complements the $EM$ estimate of tissues using the information from multi-modal scans. The new classifier takes the output tissue probability maps (TPMs), $P_c(x)$, from the $EM$ algorithm, where "$x$" represents a voxel location, and "$c$" is a single tissue class:

$$\begin{aligned}
&\forall x \in \{voxel \quad locations\}, \\
&\quad \forall c \in \{1, \ldots, \mathbb{C}\}, \\
&\exists \quad 0 \leq P_c(x) \leq 1 \quad s.t. \quad \sum_{c=1}^{\mathbb{C}} P_c(x) = 1
\end{aligned} \tag{1}$$

Where $\mathbb{C}$ is the total number of tissue types, and $P_c(x)$ represent how likely voxel location $x$ belongs to tissue type $c$.

**Training sample set** To find the candidate training sample locations "$t$" for the $k - NN$ classifier, all $P_c(x)$ from $EM$ posterior $TPMs$ are thresholded in order to identify those sample locations that have a sufficient probability to belong to a single tissue type. Increasing the threshold leads to fewer but more reliable tissue samples. A threshold of 0.7 is chosen based on the results presented by Vrooman $et\ al.$[7] for brains tissue types.

$$t \in \{x \quad | \quad \exists c \quad s.t. \quad P_c(x) \geq 0.7\} \tag{2}$$

Where training sample location $t$ is assigned with a label pointing to tissue region $c$.

Chosen training samples are then represented in an $\mathbb{F}$-dimensional feature space with:

$$\mathbb{F} = \mathbb{M} + \mathbb{C} \tag{3}$$

Where $\mathbb{F}$ is the number of features; $\mathbb{M}$ is the number of input multi-modal scans, and $\mathbb{C}$ is the total number of tissue types. The feature vector corresponding to the training sample $t$ is created as:

$$\begin{bmatrix} I_1(t), & ..., & I_{\mathbb{M}}(t), & \breve{P}_1(t), & ..., & \breve{P}_{\mathbb{C}}(t) \end{bmatrix} \tag{4}$$

Where $I_m(t)$, $m \in \{1, \ldots, \mathbb{M}\}$ represents the intensity value of the $m^{th}$ input image scan at sample location $t$, and $\breve{P}_c(t)$, $c \in \{1, \ldots, \mathbb{C}\}$ is a binary value derived from the $c^{th}$ $EM$ posterior $TPM$ at sample location $t$, such that:

$$\breve{P}_c(t) = \begin{cases} 1 & if \quad P_c(t) \geq 0.01 \\ 0 & if \quad P_c(t) < 0.01 \end{cases} \tag{5}$$

In fact, our feature space defines all the candidate regions, suggested by $EM$ results, that the current sample location $t$ probably belongs to by more than one percent chance. In this way, the fuzzy $K - NN$ classifier is restricted to only biological plausible results, and it is not biased by the probability values estimated in $EM$ step.

Finally, each created feature vector is added to a *training sample set*, and its known label code is added to the corresponding row of a *labels vector*.

**Test sample set**   Test sample locations "$s$" are the center points of the voxel locations in the first input scan, and for each test location, a feature vector is created as shown in equation (4). All the test feature vectors are then added to a *test sample set*.

**Run the algorithm**   The training and test sample sets and the labels vector are passed to a fuzzy $k - NN$ algorithm where the following procedure is performed on each test sample location:

1. In the feature space, the Euclidean distances between each test sample and all the training samples are computed. Distances are calculated through a k-dimensional tree structure[8] that is a data structure for organizing points in a k-dimensional space using space partitioning.

2. The first $\mathbb{K}$ nearest neighbors are identified from the computed distance vector. $\mathbb{K}$ needs to be an odd number, and it was set to 45 as suggested by Vrooman *et al.*[7] and Cocosco *et al.*[9]

3. New probabilities, $P_c(s)$, are computed for the test location $s$ showing how likely the current test location belongs to each tissue type. If $\mathbb{N}$ out of $\mathbb{K}$ nearest neighbors belong to tissue class $c$, then:

$$\forall s \in \{test \quad sample \quad locations\},$$
$$\forall c \in \{1, \ldots, \mathbb{C}\},$$
$$P_c(s) = \frac{\sum_{o=1}^{\mathbb{N}} \frac{1}{d_{c,o}^2}}{\sum_{i=1}^{\mathbb{K}} \frac{1}{d_i^2}} \tag{6}$$

Where $d_{c,o}$ is the distance of the $o^{th}$ occurrence of class $c$ to the current test location $s$; $d_i$ is the distance to the $i^{th}$ neighbor of the current test sample, and $P_c(s)$ represents the probability that the current test location $s$ belongs to class $c$.

4. For the test location $s$, all computed $P_c(s)$, $c \in \{1, \ldots, \mathbb{C}\}$ are stored in one **row** of a $\mathbb{S} \times \mathbb{C}$ *likelihood* matrix, where $\mathbb{S}$ is the total number of test locations, and $\mathbb{C}$ is the total number of tissue types:

$$\text{likelihood matrix} = \begin{bmatrix} & & & \cdot & & & \\ & & & \cdot & & & \\ & & & \cdot & & & \\ P_1(s), & ..., & P_c(s), & ..., & P_{\mathbb{C}}(s) \\ & & & \cdot & & & \\ & & & \cdot & & & \\ & & & \cdot & & & \end{bmatrix}_{\mathbb{S} \times \mathbb{C}} \tag{7}$$

5. Finally, new tissue probability maps are created by rearranging each **column** of the likelihood matrix to an output probability image. There are $\mathbb{C}$ output probability maps corresponding to all interested tissue types.

# 3. EXPERIMENTAL METHODS

The accuracy and effectiveness of the proposed method is evaluated qualitatively and quantitatively.

3D Slicer[10] was used to visually compare the segmentation results of the proposed enhancements to the technique established by Kim and Johnson.[3] Qualitative investigation was done using a sample $T1$-weighted MR scan that was arbitrarily selected from our local University of Iowa SIEMENS Trio Tim 3 Tesla scan protocol. This protocol was used as part of the multi-site international PREDICT-HD[11] project.

Quantitative evaluation of the proposed enhancements used a set of 18 synthetic MR datasets of a brain subject from BrainWeb database.[12] The BrainWeb database provides a rich set of multi-spectral data as *input sources* to our algorithm that include both $T1$ and $T2$ modality scans. BrainWeb also provides a simulation of the heterogeneous nature of the multi-site real data with input variants that represent six levels of noise and three degrees of bias-field for each $T1$ and $T2$. Finally, the BrainWeb data provides a set of tissue segmentation *baselines* for comparison against each *output result* from our algorithm.

The accuracy and robustness of the proposed enhancements by a fuzzy $k-NN$ algorithm were then compared to the reported results by Kim and Johnson[3] derived from an $EM$-based only classification. For this purpose, the similarity of both methods were compared against the segmentation baseline provided by BrainWeb along with the evaluation datasets.

Our atlas based approach uses the atlas definitions from two $T1$ and $T2$ modalities with priors for 15 discrete region-specific tissue types listed in table 1. This is a slightly simplified approach to that taken in[3] where 17 regions were identified with the Basal region being subdivided into (Caudate, Putamen, Accumben) regions.

Table 1: Atlas definition of 15 region-specific intensity-context priors. Each tissue type is sub-divided into regions of interest with given names. (Gm = Grey matter, Wm = white matter, Csf = cerebrospinal fluid, Crbl = Cerebellum, Vb = venous blood)

| Tissue | Name |
|---|---|
| Grey matter | Basal |
| | Hippocampus |
| | Crbl Gm |
| | Surf Gm |
| White matter | Wm |
| | Crbl Wm |
| Csf | Csf |
| Wm & Gm | Thalamus |
| | Globus |
| Venous blood | Vb |
| Background | Not Gm |
| | Not Wm |
| | Not Csf |
| | Not Vb |
| | Air |

The software is implemented based on the *InsightToolkit* libraries[13, 14] and conforms to the coding style, testing, and software guidelines identified by the National Alliance for Medical Image Computing (NAMIC) group. Our implementation is publicly available via BRAINSTools package[15] and contributes to a fully automated processing pipeline for MR images.[16, 17]

# 4. RESULTS

Figure 1 shows the visual comparison of the results on a sample MR scan from the PREDICT-HD study. As shown by corresponding arrows in both images, the segmentation boundaries of GM, WM and CSF from our
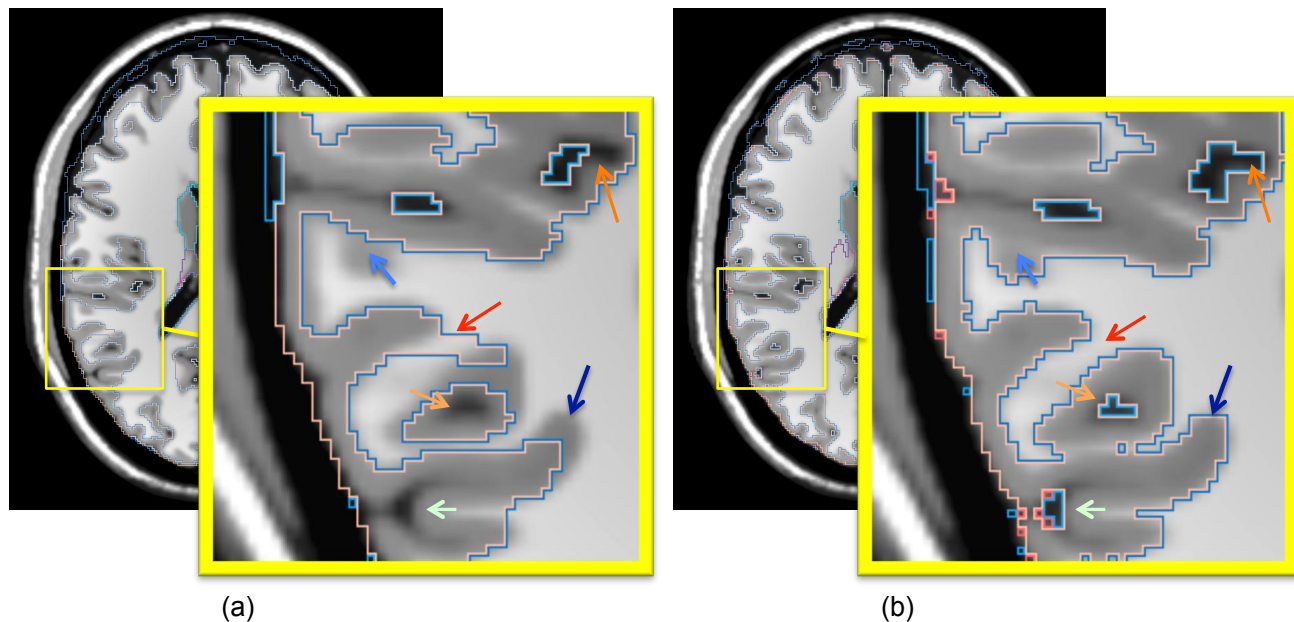
Figure 1: Visual comparison of segmentation results on a sample PREDICT-HD MR data between (a) $EM$ only based classification and (b) proposed enhancements using a fuzzy $k - NN$ classifier. More accurate delineation is achieved by the proposed approach.

proposed approach (using $K - NN$) (Fig. 1(b)) are more agreeable to real anatomical tissue boundaries than the results derived from $EM$-based only classification (Fig. 1(a)).

In order to compare the quantitative results, two independent measures, "Dice index" and "average Hausdorff distance[18]", are reported to compare the results of the automated delineations against the ground truth. Figure 2 shows the Dice index (larger is better) and average Hausdorff distance (smaller is better) evaluated along three degrees of bias-field ($rf = 0\%$, $rf = 20\%$ and $rf = 40\%$) and *six* levels of noise (0%, 1%, 3%, 5%, 7%, and 9%) for three tissue types (WM, GM and CSF). The results of $EM$ method are shown in *black* while the *blue* color is used for the results of proposed enhancements using a $k - NN$ classification.

## 5. DISCUSSION AND CONCLUSIONS

This work improved automated classification of brain tissues for multi-center $3D$ MRI data analysis. Previous studies have used expectation maximization (EM) based classification that is group specific and uses a *priori* knowledge for all the subjects in an atlas based approach. This paper, however, emphasized the importance of a non-parametric model's utility in neurodegenerative diseases, since each subject has unique anatomical states in longitudinal degenerative studies that may not be represented by prior probability distributions. Enhancements were suggested by augmenting the $EM$-based classification using a fuzzy k-nearest neighbor (k-NN) classifier that builds up a model for each individual subject and complements the classification results that $EM$ produces. A Fuzzy $k - NN$ method was selected, as this non-parametric classifier is subject specific; it is not biased by prior probability distributions, and it potentially can model more complex decision boundaries than a Gaussian distribution based mixture method.

The proposed implementation generates more precise results than $EM$ only classification, since both similarity measures, Dice index and the average Hausdorff distance, show improvement on the results of $k - NN$ classifier as demonstrated in Figure 2. Also, qualitative observations in Figure 1 show that our method especially provides more accuracy in delineation of sophisticated tissue boundaries where tissue regions are highly interleaved together like GM and WM boundaries.
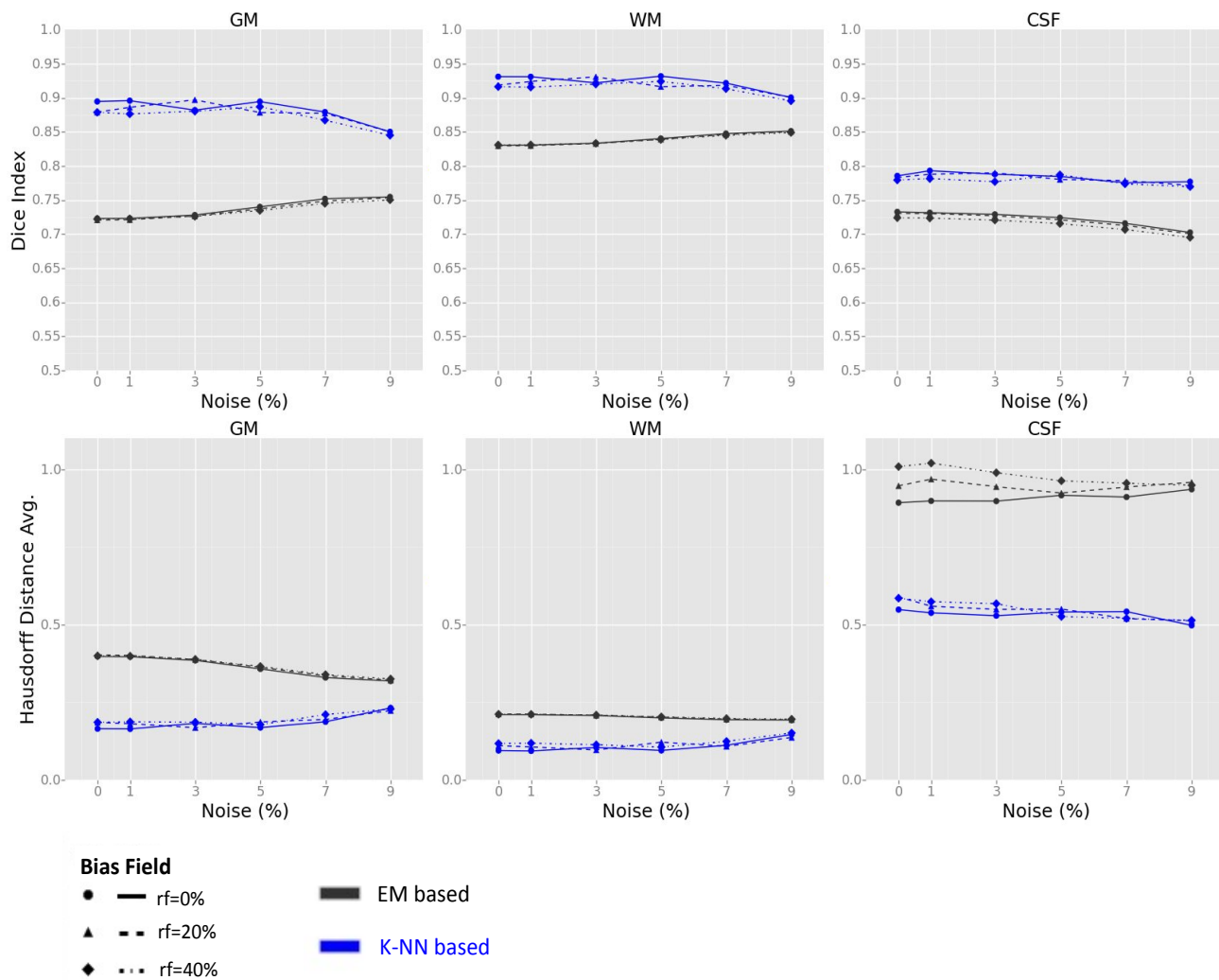
Figure 2: Comparison of the classification of cerebrospinal fluid (CSF), Grey matter (GM) and White matter (WM) tissues between $EM$ only based classification (black) and the extended method by a fuzzy $k-NN$ classifier (blue) using two independent measures, Dice index (larger is better) and average Hausdorff distance (smaller is better) . The evaluation is performed along three degrees of bias-field (rf=0, rf=20 and rf=40) and six levels of noise (0%, 1%, 3%, 5%, 7%, and 9%) along x-axis. Both similarity measures show improvement on the results of the proposed $k-NN$ enhancements.

# 6. ACKNOWLEDGMENTS

## REFERENCES

[1] Vrooman, H., van der Lijn, F., and Niessen, W., "Auto-knn: brain tissue segmentation using automatically trained knearest-neighbor classification," in [*Proceedings of the MICCAI WorkshopsThe MICCAI Grand Challenge on MR Brain Image Segmentation (MRBrainS13)*], (2013).

[2] Anbeek, P., Vincken, K. L., van Bochove, G. S., van Osch, M. J. P., and van der Grond, J., "Probabilistic segmentation of brain tissue in mr imaging," *NeuroImage* **27**(4), 795–804 (2005).

[3] Kim, E. and Johnson, H., "Robust multi-site mr data processing: iterative optimization of bias correction, tissue classification, and registration," *Frontiers in neuroinformatics* **7**, 1–11 (2013).

[4] Ghayoor, A., Vaidya, J. G., and Johnson, H. J., "Development of a novel constellation based landmark detection algorithm," in [*Proc. SPIE, Med. Img. 2013: Image Processing*], **8669**, 86693F (6pages) (2013).

[5] Avants, B. B., Epstein, C. L., Grossman, M., and Gee, J. C., "Symmetric diffeomorphic image registration with cross-correlation: evaluating automated labeling of elderly and neurodegenerative brain," *Medical image analysis* **12**(1), 26–41 (2008).

[6] Avants, B. B., Tustison, N., and Song, G., "Advanced normalization tools (ANTS)," *Insight Journal* (July 2009). Availabel online at: `http://hdl.handle.net/10380/3113`.

[7] Vrooman, H. A., Cocosco, C. A., van der Lijn, F., Stokking, R., Ikram, M. A., Vernooij, M. W., Breteler, M. M. B., and Niessen, W. J., "Multi-spectral brain tissue segmentation using automatically trained k-nearest-neighbor classification," *NeuroImage* **37**(1), 71–81 (2007).

[8] Bentley, J. L., "Multidimensional binary search trees used for associative searching," *Commun. ACM* **18**(9), 509–517 (1975).

[9] Cocosco, C. A., Zijdenbos, A. P., and Evans, A. C., "A fully automatic and robust brain mri tissue classification method," *Medical Image Analysis* **7**(4), 513–527 (2003).

[10] Fedorov, A., Beichel, R., Kalpathy-Cramer, J., Finet, J., Fillion-Robin, J.-C., Pujol, S., Bauer, C., Jennings, D., Fennessy, F., Sonka, M., Buatti, J., Aylward, S. R., j. V. Miller, Pieper, S., and Kikinis, R., "3D slicer as an image computing platform for the quantitative imaging network," *Magnetic Resonance Imaging* **30**(9), 1323–1341 (2012).

[11] "PREDICT-HD; an observational study of the earliest signs of huntington disease." `https://www.predict-hd.net`.

[12] Cocosco, C., Kollokian, V., Kwan, R., and Evans, A., "Brainweb: Online interface to a 3d mri simulated brain database," *NeuroImage* **5**, S425 (1997).

[13] Johnson, H. J., McCormick, M. M., and Ibanez, L., [*The ITK Software Guide Book 1: Introduction and Development Guidelines Fourth Edition Updated for ITK version 4.7*], Kitware, Inc. (2015).

[14] Johnson, H. J., McCormick, M. M., and Ibanez, L., [*The ITK Software Guide Book 2: Design and Functionality Fourth Edition Updated for ITK version 4.7*], Kitware, Inc. (2015).

[15] Johnson, H. J., Ghayoor, A., and Kim, R. E., "BRAINSTools; a suite of tools for medical image processing focused on brain analysis." `https://github.com/BRAINSia/BRAINSTools`. Online; accessed 25-July-2015.

[16] Kim, E. Y., Magnotta, V. a., Liu, D., and Johnson, H. J., "Stable Atlas-based Mapped Prior (STAMP) machine-learning segmentation for multicenter large-scale MRI data," *Magnetic Resonance Imaging* **32**(7), 832–844 (2014).

[17] Pierson, R. K., Johnson, H. J., Harris, G., Keefe, H., Paulsen, J. S., Andreasen, N. C., and Magnotta, V. a., "Fully automated analysis using BRAINS: AutoWorkup," *NeuroImage* **54**(1), 328–336 (2011).

[18] Dubuisson, M. and Jain, A., "A modified hausdorff distance for object matching," *12th IAPR International Conference* **1**, 566–568 (1994).