

Bootstrapping Multi-atlas Segmentation Using Multiple Automatically Generated Templates for the Segmentation of the Whole Hippocampus and Subfields

Jon Pipitone¹, Min Tae M. Park¹, Julie Winterburn¹, Tristram A. Lett^{1,9}, Jason P. Lerch^{2,3}, Jens C. Pruessner⁴, Martin Lepage^{4,5}, Aristotle N. Voineskos^{1,6,9}, M. Mallar Chakravarty^{1,6,7,8} and the Alzheimer's Disease Neuroimaging Initiative*

¹*Kimel Family Translational Imaging-Genetics Lab, Centre for Addiction and Mental Health, Toronto, ON, Canada*

²*Neurosciences and Mental Health Laboratory, Hospital for Sick Children, Toronto, ON, Canada*

³*Department of Medical Biophysics, University of Toronto, Toronto, ON, Canada*

⁴*Douglas Mental Health University Institute, Verdun, QC, Canada*

⁵*Department of Psychiatry, McGill University, Montreal, QC, Canada*

⁶*Department of Psychiatry, University of Toronto, Toronto, ON, Canada*

⁷*Institute of Biomaterials and Biomedical Engineering, University of Toronto, Toronto, ON, Canada*

⁸*Rotman Research Institute, Baycrest, Toronto, ON, Canada*

⁹*Institute of Medical Science, University of Toronto, Toronto, ON, Canada*

Abstract

Introduction: Advances in image segmentation of magnetic resonance images (MRI) have demonstrated that multi-atlas approaches improve segmentation accuracy and precision over regular atlas-based approaches. These approaches often rely on a large number of such manually segmented atlases (e.g. 30-80) that take significant time and expertise to produce. We present an algorithm, MAGeT-Brain (Multiple Automatically Generated Templates), for the automatic segmentation of the hippocampus that minimizes the number of atlases needed while still achieving similar accuracy to multi-atlas approaches. Thus, our method acts as an accurate multi-atlas approach when using special, hard-to-define atlases that are laborious to construct.

Method: MAGeT-Brain works by propagating atlas segmentations to a template library, formed from a subset of target images, via transformations estimated by non-linear image registration. The resultant segmentations are then propagated to each target image and fused using a label fusion method. We conduct a 10-fold Monte Carlo cross-validation of whole hippocampal segmentation on a pool of 60 ADNI subjects over a range of parameter settings, and a leave-one-out cross-validation (LOOCV) of hippocampal subfield segmentation using five high-resolution atlases. Two final experiments assess MAGeT-Brain

*Data used in preparation of this article were obtained from the Alzheimer's Disease Neuroimaging Initiative (ADNI) database (adni.loni.ucla.edu). As such, the investigators within the ADNI contributed to the design and implementation of ADNI and/or provided data but did not participate in analysis or writing of this report. A complete listing of ADNI investigators can be found at: http://adni.loni.ucla.edu/wp-content/uploads/how_to_apply/ADNI_Acknowledgement_List.pdf

when applied to first episode psychosis and Alzheimer’s disease populations, and MAgE-T-Brain segmentations are compared with existing automated methods (FSL, FreeSurfer, MAPER) and biases are explored.

Results: Using nine atlases and 19 template images, MAgE-T-Brain achieves a mean Dice’s Similarity Coefficient (DSC) of 0.84 over 10-fold Monte Carlo cross-validation on 60 subjects relative to manual segmentations, and shows significantly lower variability in DSC than multi-atlas segmentation. In a LOOCV, MAgE-T-Brain reproduces the volumes of the cornu ammonis (CA) 1; CA4/dentate gyrus (DG); and strata radiatum (SR), strata lacunosum (SL), and strata moleculare (SM) hippocampal subfields with a percent error in volume that is at or lower than that produced by image resampling. MAgE-T-Brain produces hippocampal volumes in a first episode psychosis patient population that are highly correlated with expert manual segmentation volumes (Pearson $r = 0.877, t = 16.244, p < 0.001$). Compared to FSL and FreeSurfer, MAgE-T-Brain shows much smaller fixed volume bias (within $250mm^3$ on average) to semi-automated (SNT) segmentations available from ADNI, as well as a conservative, rather than exaggerated, proportional volume bias.

Conclusion: We demonstrate that MAgE-T-Brain produces accurate hippocampal segmentations using only 5 atlases over different hippocampal definitions, disease populations, and acquisition types, as well as showing that accurate identification of the hippocampal subfields is possible.

Contact:


Jon Pipitone and M. Mallar Chakravarty
 Kimel Family Translation Imaging-Genetics Research Laboratory
 Research Imaging Centre
 Centre for Addiction and Mental Health
 250 College St.
 Toronto, Canada M5T 1R8
 jon.pipitone@camh.ca; mallar.chakravarty@camh.ca

1 Introduction


The hippocampus is a brain structure situated in the medial temporal lobe, and has long been associated with learning and memory (den Heijer et al., 2012; Scoville and Milner, 2000). The hippocampus is of interest to clinical neuroscientists because it is implicated in many forms of brain dysfunction, including Alzheimer’s disease (Sabuncu et al., 2011) and schizophrenia (Narr et al., 2004; Karnik-Henry et al., 2012). In neuroimaging studies, structural magnetic resonance images (MRI) are often used for the volumetric assessment of the hippocampus. As such, accurate segmentation of the hippocampus and its subfields in MRI is a necessary first step to better understand the inter-individual variability of subject neuroanatomy.

The gold standard for neuroanatomical segmentation is manual delineation by an expert human rater. However, with the availability of increasingly large MRI datasets, the time and expertise required for manual segmentation becomes prohibitive (Mazziotta et al., 1995, 2001; Mazziotta et al.; Pausova et al., 2007). This effort is complicated by the fact that there is significant variation between segmentation protocols with respect to specific anatomical boundaries of the hippocampus (Geuze et al., 2004) and this has led to efforts to create an unified hippocampal segmentation protocol (Jack et al., 2011). In addition, there is controversy over the appropriate manual segmentation protocol to use in a particular imaging study (Nestor et al., 2012). Thus, a segmentation algorithm that can easily adapt to different manual segmentation definitions would be

of significant benefit to the neuroimaging community.

Automated segmentation techniques that are reliable, objective, and reproducible can be considered complementary to manual segmentation. In the case of classical model-based segmentation methods (Haller et al., 1997; Csernansky et al., 1998), an MRI atlas that was previously manually labelled by an expert rater is matched to target images using nonlinear registration methods. The resulting nonlinear transformation is applied to the manual labels (i.e. *label propagation*) to  them into the target image space. While this methodology has been used successfully in several contexts (Chakravarty et al., 2008, 2009; Collins et al., 1995; Haller et al., 1997), it is limited in accuracy due to error in the estimated nonlinear transformation itself, partial volume effects in label resampling, and irreconcilable differences between the neuroanatomy represented within the atlas and target images.

One methodology that can be used to mitigate these sources of errors involves the use of multiple manually segmented atlases and probabilistic segmentation techniques, such as those found in the FreeSurfer package (Fischl et al., 2002). FreeSurfer uses a probabilistic atlas of anatomical and tissue classes along with spatial constraints for class labels encoded using a Markov random field model to segment the entire brain.

More recently, many groups have used multiple atlases to improve overall segmentation accuracy (i.e. multi-atlas  segmentation) over model-based approaches (Heckemann et al., 2006a, 2011; Collins and Pruessner, 2010; Lötjönen et al., 2010; Aljabar et al., 2009; Leung et al., 2010; Wolz et al., 2010). Each atlas image is registered to a target image, and label propagation is performed to produce several labellings of the target image (one from each atlas). A *label fusion* technique, such as voxel-wise voting, is used to merge these labels into the definitive segmentation for the target. In addition, weighted voting procedures that use *atlas selection* techniques are often used to exclude atlases from label fusion that are dissimilar to a target image in order to reduce error from unrepresentative anatomy (Aljabar et al., 2009). This involves the selection of a subset of atlases using a similarity metric such as cross-correlation (Aljabar et al., 2009) or normalized mutual information. Such selection has the added benefit of significantly reducing the number of nonlinear registrations. For example Collins and Pruessner (2010) demonstrated that only 14 atlases, selected based on highest similarity between medial temporal lobe neuroanatomy as evaluated by normalized mutual information (Studholme et al., 1999) from a library of 80 atlases, were required to achieve accurate segmentations of the hippocampus. Additionally, several methods have been explored for label fusion, including the STAPLE algorithm (Simultaneous Truth And Performance Level Estimation; Warfield et al. (2004)) that computes a probabilistic segmentation using an expectation-maximization framework from an set of competing segmentations, or others where a subset of segmentations can be estimated using metrics such as the sum of squared differences in the regions of interest to be segmented (Coupé et al., 2012).

However, many of these methods require significant investment of time and resources for the creation of the atlas library ranging between 30 (Heckemann et al., 2006a) and 80 (Collins and Pruessner, 2010) manually segmented atlases. This strategy has the main drawback of being inflexible as it does not easily accommodate varying the definition of the hippocampal anatomy (such as the commonly used heuristic of subdividing the hippocampus into head, body, and tail (Poppenk and Moscovitch, 2011; Pruessner et al., 2000)). Furthermore, none of these methods have demonstrated sufficient flexibility to accommodate atlases that are somehow exceptional such as those derived from serial histological data (Chakravarty et al., 2006; Yelnik et al., 2007) or high-resolution MRI data that enables robust identification of hippocampal subfields (Winterburn et al., 2013; Yushkevich et al., 2009; Mueller and Weiner, 2009; Van Leemput et al., 2009; Wisse et al., 2012). Due to the recent availability of the latter, there has been increased interest in the use

of probabilistic methods for the identification of the hippocampal subfields on standard T1-weighted images. Our group recently demonstrated that through use of an intermediary automated segmentation stage, robust and accurate segmentation of the striatum, pallidum, and thalamus using a single atlas derived from serial histological data is possible (Chakravarty et al., 2012). The novelty of this manuscript is the extension of our multi-atlas methodology to the hippocampus using more than a single input atlas, while simultaneously limiting the number of inputs used during segmentation, and demonstrating that accurate identification of the hippocampal subfields is indeed possible using this methodology.

Of central relevance to the present work is the LEAP algorithm (Learning Embeddings for Atlas Propagation; Wolz et al. (2010)) because of its focus on performing multi-atlas segmentation with a limited number of input atlases. The LEAP algorithm is a clever modification to the basic multi-atlas strategy in which an atlas library is grown, beginning with a set of manually labelled atlases, by successively incorporating unlabelled target images once they themselves have been labelled using multi-atlas techniques. The sequence in which target images are labelled is chosen so that the similarity between the atlas images and the target images is minimised at each step, effectively allowing for deformations between very dissimilar images to be broken up into sequences of smaller deformations. Although Wolz et al. (2010) begin with an atlas library of 30 MR images, this method could theoretically work using a much smaller atlas library. In their validation, LEAP was used to segment the whole hippocampus in the ADNI-1 baseline dataset, achieving a mean Dice score of 0.85 against Gold Standard segmentations.

Also of interest to this manuscript are methods that attempt to define hippocampal subfields using standard T1-weighted data. To the best of our knowledge, there are only two automated segmentation algorithms that attempt this problem. The first is included with the FreeSurfer package (Van Leemput et al., 2009). This work is limited as it omits the tail of the hippocampus and the segmentation protocol has yet to be fully validated. Nonetheless, it demonstrates that the applicability of hippocampal subfield segmentation using data from 10 subjects. In the second method, Yushkevich et al. (2009) hippocampal subfields were acquired and labelled on high-resolution MRI data from post-mortem medial temporal lobe samples. Using nonlinear registration guided by manually derived hippocampus masks and specific landmarks, the authors demonstrated accurate parcellation of hippocampal subfields in unlabelled MRI volumes.

In this paper we have describe a thorough validation of the MAGeT-Brain algorithm for segmentation of the hippocampus and its subfields. First, we address the feasibility and accuracy of whole hippocampus segmentation with a limited number of input atlases (Chakravarty et al., 2012) by performing a multi-fold experiment using a subset of the Alzheimer’s Disease Neuroimaging Initiative (ADNI) dataset over a range of parameters. We then perform a leave-one-out validation to determine if hippocampal subfields can be accurately identified, using the best parameters discovered in the first experiment. To ensure that we have not overfit our parameters to the aging or neurodegenerative brain, we also apply MAGeT-Brain to a dataset of individuals suffering from first episode psychosis. Finally, we validate our algorithm using all of the data available in the ADNI1:Complete 1Yr 1.5T sample and compare our segmentations to other popular segmentation algorithms.

2 Methods

2.1 The MAGeT-Brain Algorithm

In this paper, we use the term *label* to mean any segmentation (manual or derived) of an MR image. *Label propagation* is the process by which two images are registered and the resulting transformation is applied to the labels from one image to bring them into alignment with the other image. We use the term *atlas* to mean a manually segmented image, and the term *template* to mean an automatically segmented image (i.e. via label propagation). The terms *atlas library* and *template library* describe any set of such images. Additionally, we use the term *target* to refer to an unlabelled image that is undergoing segmentation.

The simplest form of multi-atlas segmentation, which we call *basic multi-atlas segmentation*, involves three steps. First, each labelled input image (i.e. atlas or template) is registered to an unlabelled target image. Second, the labels from each image are propagated to the target image space. Third, the labels are combined into a single label by label fusion (Heckemann et al., 2006a, 2011). The basic multi-atlas segmentation method is described in detail in other publications (Collins and Pruessner, 2010; Heckemann et al., 2011; Aljabar et al., 2009). When only a single atlas is used, basic multi-atlas segmentation degenerates into model-based segmentation: labels are propagated from the atlas to a target, and no label fusion is needed.

MAGeT-Brain (Multiple Automatically Generated Templates) bootstraps the creation of a large template library given a limited input atlas library, and then uses the template library in basic multi-atlas segmentation. Images for the template library are selected from a set of input target images, either arbitrarily or so as to reflect the neuroanatomy or demographics of the target set as a whole (for instance, by sampling equally from cases and controls). The template library images are automatically labelled by each of the atlases via label propagation. Basic multi-atlas segmentation is then conducted using the template library to segment the entire set of target images (including the target images used in the construction of the template library). Since each template library image has multiple labels (one from each atlas), the final number of labels to be fused for each target may be quite large (i.e. # of atlas \times # of templates).

Figure 1 illustrates the MAGeT-Brain algorithm graphically. Source code for MAGeT-Brain can be found at <http://github.com/pipitone/MAGeTbrain>.

2.2 Experiments

The following section describes for experiments conducted to assess the segmentation quality of the MAGeT-Brain algorithm. The first two experiments assess the validity of MAGeT-Brain using cross-validation designs. Experiment 1 investigates the accuracy of MAGeT-Brain whole hippocampus segmentation over a wide range of parameter settings. The results of this experiment enable us to choose the parameter settings offering the best performance for use in subsequent experiments. Experiment 2 assesses hippocampal subfield segmentation quality in a leave-one-out cross-validation design. The last two experiments assess the validity of the MAGeT-Brain algorithm when applied to different diseases: first episode schizophrenia (Experiment 3), and Alzheimer’s disease (Experiment 4). Additionally, in Experiment 4, we compare MAGeT-Brain segmentations with those of well-known automated methods and assessed segmentation bias.

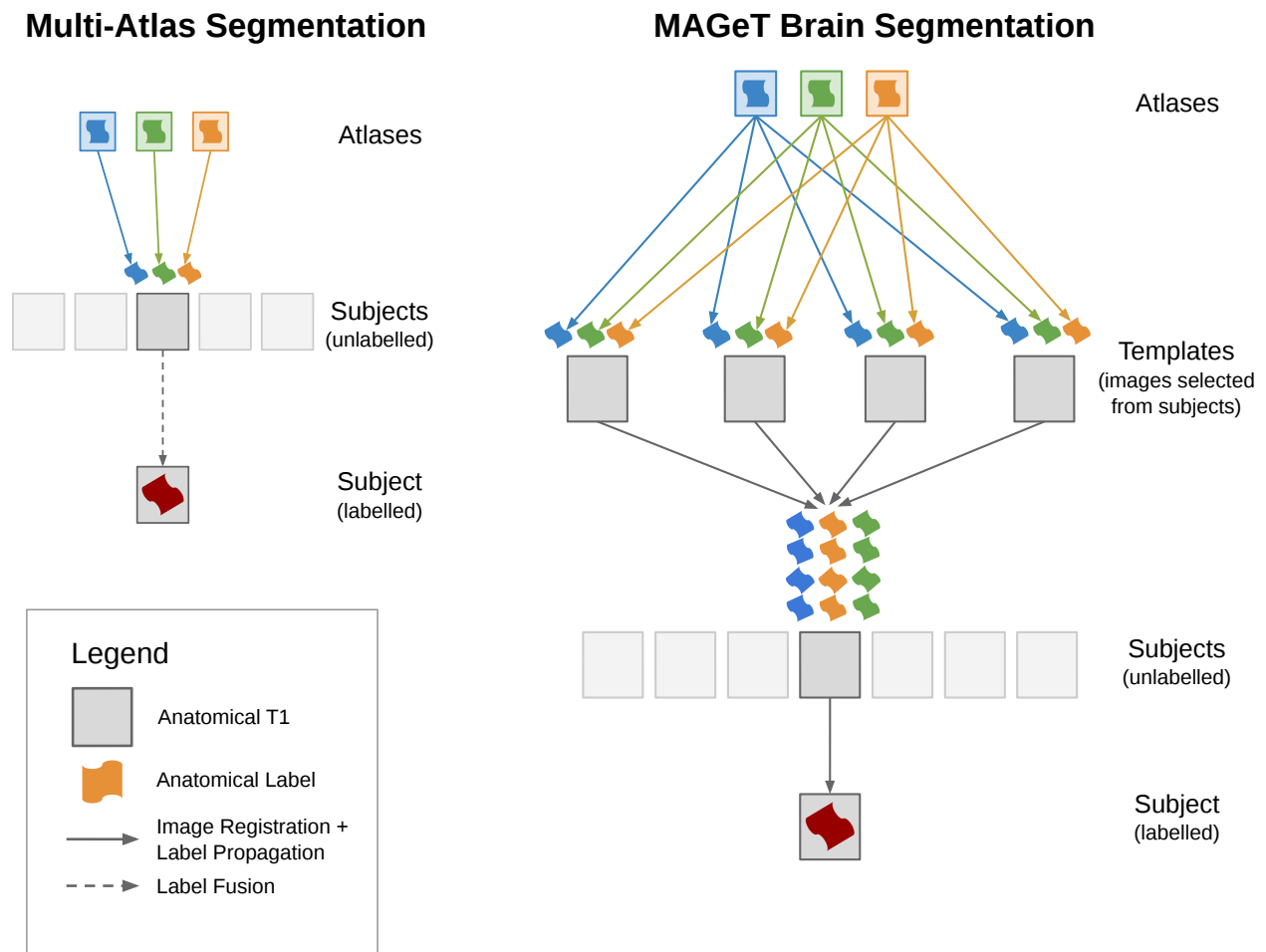


Figure 1: A schematic illustration of basic multi-atlas segmentation and MAgE-T-Brain segmentation

2.2.1 Experiment 1: Whole Hippocampus Cross-Validation

Monte Carlo Cross-Validation (MCCV) (Shao, 1993) was performed using images from the ADNI dataset (Hsu et al., 2002), and manual whole hippocampus segmentations following the protocol laid out in (Pruessner et al., 2000).

This form of cross-validation allows us to rigorously validate a large number of parameter settings of MAGeT-Brain (atlas and template library sizes, registration algorithm, and label fusion method) and select the best parameters to use in subsequent experiments.

ADNI1:Complete 1Yr 1.5T dataset Data used in the preparation of this article were obtained from the Alzheimer’s Disease Neuroimaging Initiative (ADNI) database (adni.loni.ucla.edu). The ADNI was launched in 2003 by the National Institute on Aging (NIA), the National Institute of Biomedical Imaging and Bioengineering (NIBIB), the Food and Drug Administration (FDA), private pharmaceutical companies and non-profit organizations, as a \$60 million, 5-year public-private partnership. The primary goal of ADNI has been to test whether serial magnetic resonance imaging (MRI), positron emission tomography (PET), other biological markers, and clinical and neuropsychological assessment can be combined to measure the progression of mild cognitive impairment (MCI) and early Alzheimer’s disease (AD). Determination of sensitive and specific markers of very early AD progression is intended to aid researchers and clinicians to develop new treatments and monitor their effectiveness, as well as lessen the time and cost of clinical trials.

The Principal Investigator of this initiative is Michael W. Weiner, MD, VA Medical Center and University of California San Francisco. ADNI is the result of efforts of many co- investigators from a broad range of academic institutions and private corporations, and subjects have been recruited from over 50 sites across the U.S. and Canada. The initial goal of ADNI was to recruit 800 subjects but ADNI has been followed by ADNI-GO and ADNI-2. To date these three protocols have recruited over 1500 adults, ages 55 to 90, to participate in the research, consisting of cognitively normal (CN) older individuals, people with early or late MCI, and people with early AD. The follow up duration of each group is specified in the protocols for ADNI-1, ADNI-2 and ADNI-GO. Subjects originally recruited for ADNI-1 and ADNI-GO had the option to be followed in ADNI-2. For up-to-date information, see www.adni-info.org.

60 1.5T images were arbitrarily selected from the baseline scans in the *ADNI1:Complete 1Yr 1.5T* standardized dataset. 20 subjects were chosen from each disease category: cognitively normal (CN), mild cognitive impairment (MCI) and Alzheimer’s disease (AD). Demographics for this subset are shown in Table 1. Fully manual segmentations of the left and right whole hippocampi in these images were provided by one author (JCP) according to the segmentation protocol specified in Pruessner et al. (2000).

Clinical, demographic and pre-processed T1-weighted MRI were downloaded by the authors from the ADNI database (adni.loni.ucla.edu) between March 2012 and August 2012. The image dataset used was the “ADNI1:Complete 1Yr 1.5T” standardized dataset available from ADNI ¹ (Wyman et al., 2012). This image collection contains uniformly pre-processed images which have been designated to be the “best” after quality control. All images were acquired using 1.5T scanners (General Electric Healthcare, Philips Medical Systems or Siemens Medical Solutions) at multiple sites using the protocol described in (Jack et al., 2008). Representative 1.5T imaging parameters were TR = 2400ms, TI = 1000ms, TE = 3.5ms, flip angle = 8°, field of view = 240 x 240mm, a 192 × 192 × 166 matrix (*x*, *y*, and *z* directions) yielding a voxel dimensions

¹<http://adni.loni.ucla.edu/methods/mri-analysis/adni-standardized-data/>

Table 1: **ADNI-1 cross-validation subset demographics.** CN - Cognitively Normal. LMCI - Late-onset Mild Cognitive Impairment. AD - Alzheimer’s Disease. Hisp - Hispanic. CDR-SB - Clinical Dementia Rating-Sum of Boxes. ADAS - Alzheimer’s Disease Assessment Scale. MMSE - Mini-Mental State Examination.

	CN	LMCI	AD	Combined
	$N = 20$	$N = 20$	$N = 20$	$N = 60$
Age at baseline Years	72.2 75.5 80.3	70.9 75.6 80.4	69.4 74.9 80.1	70.9 75.2 80.2
Sex : Female	50% (10)	50% (10)	50% (10)	50% (30)
Education	14.0 16.0 18.0	13.8 16.0 16.5	12.0 15.5 18.0	13.0 16.0 18.0
CDR-SB	0.00 0.00 0.00	1.00 2.00 2.50	3.50 4.00 5.00	0.00 1.75 3.62
ADAS 13	6.00 7.67 11.00	14.92 20.50 25.75	24.33 27.00 32.09	9.50 18.84 26.25
MMSE	28.8 29.5 30.0	26.0 27.5 28.2	22.8 23.0 24.0	24.0 27.0 29.0

a b c represent the lower quartile a , the median b , and the upper quartile c for continuous variables. Numbers after percents are frequencies.

of $1.25mm \times 1.25mm \times 1.2mm$.

Experiment details Monte Carlo Cross-Validation (MCCV), also known as repeated random sub-sampling cross-validation, consists of repeated rounds of validation conducted on a fixed dataset (Shao, 1993). In each round, the dataset is randomly partitioned into a training set and a validation set. The method to be validated is then given the training data, and its output is compared with the validation set.

In this experiment, our dataset consists of 60 1.5T images and corresponding Pruessner-protocol manual segmentations. In each validation round, the dataset is partitioned into a training set consisting of images and manual segmentations used as an atlas library, and a validation set consisting of the remaining images to be segmented by both MAgE-T-Brain and multi-atlas. The computed segmentations are compared to the manual segmentations (see Evaluation below).

A total of ten validation rounds were performed on each subject in the dataset, over each combination of parameter settings. The parameter settings explored are: atlas library size (1-9), template library size (1-20), registration method (ANTS or ANIMAL, described below), and label fusion method (majority vote, cross-correlation weighted majority vote, and normalized mutual information weighted majority vote, described below). In each validation round, both a MAgE-T-Brain and multi-atlas segmentation is produced. A total of $10 \times 60 \times 9 \times 20 \times 2 \times 3 = 6.48 \times 10^5$ validation rounds were conducted and resulting segmentations analysed.

Before registration, all images underwent preprocessing with the N3 algorithm (Sled et al., 1998) to minimize intensity nonuniformity. In this experiment we compared two non-linear image registration methods:

Automatic Normalization and Image Matching and Anatomical Labeling (ANIMAL) The ANIMAL algorithm carries out image registration in two phases. In the first, a 12-parameter linear transformation (3 translations, rotations, scales, shears) is estimated between images using an algorithm that maximizes the correlation between blurred MR intensities and gradient magnitude over the whole brain (Collins et al., 1994). In the second phase, nonlinear registration is completed using the ANIMAL algorithm (Collins et al., 1995): an iterative procedure that estimates a 3D deformation field between two MR images. At first, large deformations are estimated using a blurred version of the input data. These larger deformations are then input to subsequent steps where the fit is refined by estimating smaller deformations on data

Table 2: **ANIMAL registration parameters.**

Parameters	Stage 1	Stage 2	Stage 3
Model Blur (FWHM)	8	8	4
Input Blur (FWHM)	8	8	4
Iterations	30	30	10
Step	8x8x8	4x4x4	2x2x2
Sub-Lattice	6	6	6
Lattice Diameter	24x24x24	12x12x12	6x6x6

blurred with a Gaussian kernel with a smaller full width at half maximum (FWHM). The final transformation is a set of local translations defined on a bed of equally spaced nodes that were estimated through the optimization of the correlation coefficient. For the purposes of this work we used the regularization parameters optimized in Robbins et al. (2004), displayed in Table 2.

Automatic Normalization Tools (ANTS) ANTS is a diffeomorphic registration algorithm which provides great flexibility over the choice of transformation model, objective function, and the consistency of the final transformation (Avants et al., 2008). The transformation is estimated in a hierarchical fashion where the MRI data is subsampled, allowing large deformations to be estimated and successively refined at later hierarchical stages (where the data is subsampled to a finer grid). The deformation field and the objective function are regularized with a Gaussian kernel at each level of the hierarchy. The ANTS algorithm is freely available <http://www.picsl.upenn.edu/ANTS/>. We used an implementation of the ANTS algorithm compatible with the MINC data format, mincANTS <https://github.com/vfonov/mincANTS>.

We used the following command line when running ANTS:

```
mincANTS 3 -m PR[target_file.mnc,source_file.mnc,1,4]
--number-of-affine-iterations 10000x10000x10000x10000x10000
--affine-gradient-descent-option 0.5x0.95x1.e-4x1.e-4
--use-Histogram-Matching --MI-option 32x16000
-r Gauss[3,0] -t SyN[0.5] -i 100x100x100x20
-o transformation.xfm
```

These settings were adapted from the "reasonable starting point" given in the ANTS manual ².

Label fusion methods Label fusion is a term given to the process of combining the information from several candidate labels for an image into a single labelling. In this experiment we explore three fusion methods:

Voxel-wise Majority Vote Labels are propagated from all template library images to a target. Each output voxel is given the most frequent label at that voxel location amongst all candidate labels.

Cross-correlation Weighted Majority Vote An optimal combination of targets from the template library has previously been shown to improve segmentation accuracy (Aljabar et al., 2009; Collins and

²<https://sourceforge.net/projects/advants/files/Documentation/>

Pruessner, 2010). In this method, each template library image is ranked in similarity to each unlabelled image by the normalized cross-correlation (CC) of image intensities after linear registration, over a region of interest (ROI) generously encompassing the hippocampus. Only the top ranked template library image labels are used in a voxel-wise majority vote. The ROI is heuristically defined as the extent of all atlas labels after linear registration to the template, dilated by three voxels (Chakravarty et al., 2012). The number of top ranked template library image labels is a configurable parameter and displayed as the size of the template library in the rest of the paper.

The `xcorr_vol` utility from the ANIMAL toolkit is used to calculate the cross-correlation similarity measure.

Normalised Mutual Information Weighted Majority Vote This method is similar to cross-correlation weighted voting except that image similarity is calculated by the normalised mutual information score over the region of interest (Studholme et al., 2001). The `itk.similarity` utility from the EZMinc toolkit³ is used to calculate the normalised mutual information measure between two images.

Evaluation method The Dice similarity coefficient (DSC), also known as Dice’s Kappa, assesses the agreement between two segmentations. It is one of the most widely used measures of segmentation accuracy, and we use it as the basis of comparison in this experiment.

$$\text{Dice's coefficient (DSC)} = \frac{2|A \cap B|}{|A| + |B|}$$

where A and B are the regions being compared, and the cardinality is the volume measured in voxels.

The labels produced by MAgE-T-Brain and multi-atlas segmentation are compared to the manual labels using the Dice similarity coefficient, and the recorded value for each subject at each parameter setting explored in this experiment is the average over ten validation rounds.

Additionally, the sensitivity of MAgE-T-Brain multi-atlas and template library composition is evaluated by comparing the variability in Dice scores over all validation rounds at fixed parameter settings. This is achieved by first computing the variance of DSC scores in each block of ten validation rounds per subject. The distribution of these statistics across all subjects is then compared between MAgE-T-Brain and multi-atlas using a Student’s t-test. A significant difference between distributions is taken to show either a larger or smaller level of variability between methods.

2.2.2 Experiment 2: Hippocampal Subfield Cross-Validation

The previous experiment assesses MAgE-T-Brain performance on whole hippocampus segmentation. In this experiment, we assess MAgE-T-Brain hippocampal subfield segmentation of standard 3T T1-weighted images using a leave-one-out cross-validation (LOOCV) design.

Winterburn Atlases The Winterburn atlases (Winterburn et al., 2013) are digital hippocampal segmentations of five in-vivo 0.3mm-isotropic T1-weighted MR images. The segmentations include subfield segmentations for the cornu ammonis (CA) 1; CA2 and CA3; CA4 and dentate gyrus; subiculum; and strata radiatum (SR), strata lacunosum (SL), and strata moleculare (SM). Subjects in the Winterburn atlases range in age from 29-57 years (mean age of 37), and include two males and three females.

³<https://github.com/vfonov/EZminc>

Table 3: **Demographics for the hippocampal subfield cross-validation healthy control subject sample.** Education is shown in years.

	N	Control N = 16
Age	16	31.0 53.0 63.8
Sex : Male	16	38% (6)
Education : 0.01	15	7% (1)
12		13% (2)
13		13% (2)
14		20% (3)
16		13% (2)
18		33% (5)
Handedness : R	16	94% (15)

a b c represent the lower quartile a , the median b , and the upper quartile c for continuous variables.
 N is the number of non-missing values.
Numbers after percents are frequencies.

Experiment details Leave-one-out cross-validation (LOOCV) is a validation approach in which an algorithm is given all but one item in a dataset as training data (in our case, atlas images and labels) and then is applied to the left-out item. This is done, in turn, for each item in the dataset and the output across all items is evaluated.

In this experiment, the high-resolution $0.3mm^3$ voxel Winterburn atlases are used as the MAgE-T-Brain atlas library, but in each round of LOOCV a $0.9mm^3$ voxel target image corresponding to the left-out atlas is used so as to evaluate MAgE-T-Brain on standard 3T T1-weighted resolution images.

Specifically, two sets of $0.9mm^3$ voxel target images are used in this experiment: separately acquired T1 acquisitions of four of the five Winterburn atlas subjects (referred to as the *BRAVO* set because of the scan parameters used, described below), and computationally subsampled versions of the Winterburn atlas images (the *Subsampled* set).

The 3T T1 BRAVO acquisitions of four of the five Winterburn atlas subjects were separately obtained within a short time of the original atlas image acquisitions. Images were acquired on a 3T GE Discovery MR 750 system (General Electric, Milwaukee, WI) using an 8-channel head coil with the enhanced fast gradient recalled echo 3-dimensional acquisition protocol, EFGRE-BRAVO, with the following parameters: $TE/TR/TI = 3.0ms/6.7ms/650ms$, flip angle = 8° , $FOV = 15.8cm$, slice thickness = $0.9mm$, 176 in-plane steps for an approximate isotropic resolution of $0.9mm$ dimension voxels.

Image subsampling of the Winterburn atlases was performed using trilinear interpolation techniques.

The template library is composed of all $0.9mm^3$ voxel target images (either *BRAVO* or *Subsampled*), plus an additional set of 16 of healthy subjects 3T T1 images acquired separately (Table 3). These images were acquired on a 3T GE Discovery MR 750 system (General Electric, Milwaukee, WI) using an 8-channel head coil with the enhanced fast gradient recalled echo 3-dimensional acquisition protocol, EFGRE-BRAVO, with the following parameters: $TE/TR/TI = 3.0ms/6.7ms/650ms$, flip angle = 8° , $FOV = 15.3cm$, slice thickness = $0.9mm$, 170 in-plane steps for an approximate isotropic resolution of $0.9mm$ dimension voxels.

The optimal size of template library, registration method, and label fusion method found in Experiment 1 are used. Figure 2 illustrates schematically the experimental set up.

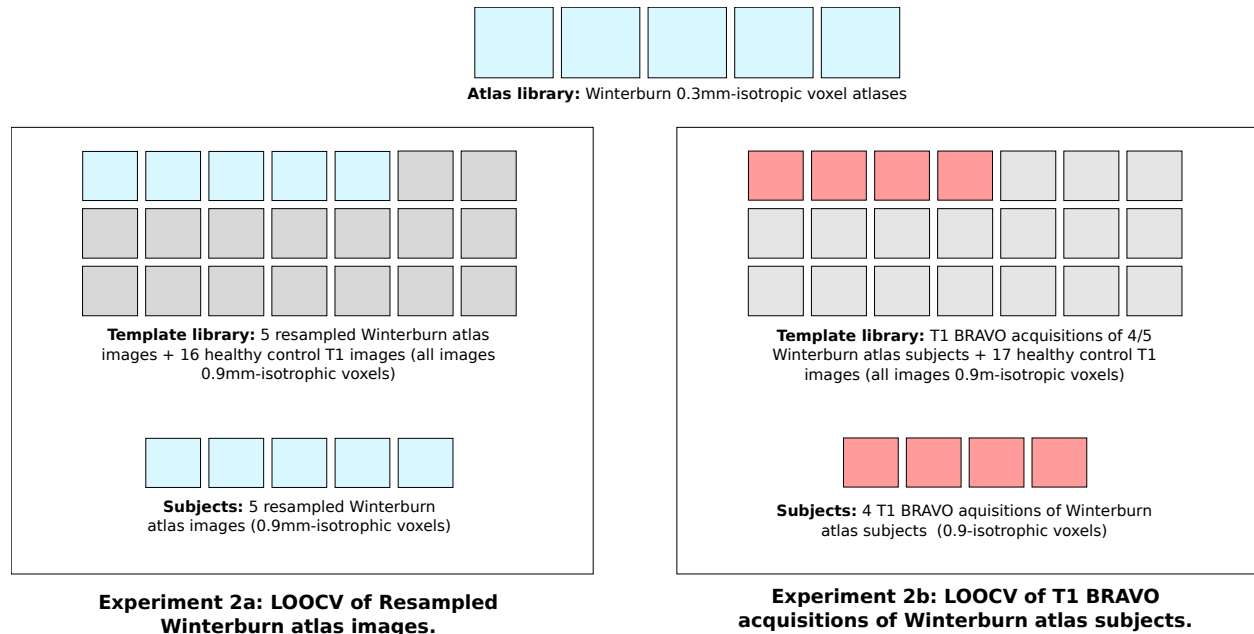


Figure 2: A schematic illustration of the experimental set up of the leave-one-out cross-validation (LOOCV) subfield experiment.

Evaluation method In this experiment we must evaluate MAgE-T-Brain on a task which is ill-defined since there are no segmentation protocols of the hippocampal subfields for T1 images at $0.9mm^3$ voxels (in fact, whether subfields can even be differentiated at this resolution is debated). Therefore, we evaluate the *precision* with MAgE-T-Brain produces hippocampal subfields that correspond to the segmentation protocol used by the given atlas library images.

Specifically, we compute the relative percent error in volume of the MAgE-T-Brain-produced hippocampal subfields with respect to the full-resolution Winterburn atlas segmentations. In addition, by directly sampling the Winterburn atlas segmentations to $0.9mm^3$ voxel images (using nearest-neighbour interpolation) we can obtain a baseline relative percent error to compare MAgE-T-Brain against.

2.2.3 Experiment 3: Application to the segmentation of first episode schizophrenia patients

To validate that the MAgE-T-Brain works effectively in the context of other neurological disorders, in this experiment we use the Winterburn atlases to derive whole hippocampal segmentations of a dataset of patients with Schizophrenia. The resulting segmentations are assessed for quality by comparison with expert manual segmentations.

SZ-FEP dataset All patients were recruited and treated through the Prevention and Early Intervention Program for Psychoses (PEPP-Montreal), a specialized early intervention service at the Douglas Mental Health University Institute in Montreal, Canada. People aged 14 to 35 years from the local catchment area suffering from either affective or non-affective psychosis who had not taken antipsychotic medication for more than one month with an IQ above 70 were consecutively admitted as either in- or out-patients. Of

Table 4: **Schizophrenia First Episode Patient Demographics.** ambi - ambidextrous. SES - Socioeconomic Status score. FSIQ - Full Scale IQ.

	N	FEP
		<i>N</i> = 81
Age	80	21 23 26
Gender : M	81	63% (51)
Handedness : ambi	81	6% (5)
left		5% (4)
right		89% (72)
Education	81	11 13 15
SES : lower	81	31% (25)
middle		54% (44)
upper		15% (12)
FSIQ	79	88 102 109

a b c represent the lower quartile *a*, the median *b*, and the upper quartile *c* for continuous variables.

N is the number of non-missing values.

Numbers after percents are frequencies.

those treated at PEPP, only patients aged 18 years with no previous history of neurological disease or head trauma causing loss of consciousness were eligible for the neuroimaging study; only those suffering from schizophrenia spectrum disorders were considered for this analysis. For complete program details see Malla et al. (2003).

Scanning of 81 subjects was carried out at the Montreal Neurological Institute on a 1.5-T Siemens whole body MRI system. Structural T1 volumes were acquired for each participant using a three-dimensional (3D) gradient echo pulse sequence with sagittal volume excitation (repetition time=22ms, echo time=9.2ms, flip angle=30°, 180 1mm contiguous sagittal slices). The rectangular field-of-view for the images was 256mm (SI)×204mm (AP). Subject demographics are shown in Table 4.

Manual segmentation of each subject whole hippocampus is produced following a validated segmentation protocol (Pruessner et al., 2000).

Experiment details MAGeT-Brain is configured with an atlas library composed of the Winterburn T1 atlases (see Experiment 2) ignoring subfield delineations. All images from the SZ-FEP dataset are segmented by MAGeT-Brain. The optimal size of template library, registration method, and label fusion method found in Experiment 1 are used.

Evaluation method The Pruessner and Winterburn hippocampal segmentation protocols differ slightly in the neuroanatomical features that are delineated (Winterburn et al., 2013). This difference poses a problem for evaluation by measuring overlap. That is, since different protocols will necessarily produce segmentations that do not perfectly overlap, the degree of overlap cannot be solely used to compare segmentation methods using different protocols. In place of an overlap metric, we can assess the degree of (Pearson) correlation in average bilateral hippocampal volume of the subjects produced by each method.

2.2.4 Experiment 4: Application to the segmentation of Alzheimer’s disease patients

To validate MAGeT-Brain segmentation quality with respect to other established automated hippocampal segmentation methods, MAGeT-Brain was applied to large dataset from the ADNI project and the resulting segmentations were compared to those produced by FreeSurfer, FSL, MAPER, as well as semi-automated hippocampal segmentations (SNT) provided by ADNI.

ADNI1 dataset revisited The *ADNI1:Complete 1Yr 1.5T* standardized dataset contains 1919 images in total. SNT, MAPER, FreeSurfer hippocampal volumes for a subset of images were provided by ADNI, along with quality control data for each FreeSurfer segmentation (guidelines described in (Hartig et al., 2010)). Clinical and demographic data for the entire ADNI1:Complete 1Yr 1.5T dataset are shown in Table 5.

Semi-automated segmentations of the left and right whole hippocampi are made available with a subset of ADNI images (Hsu et al., 2002). These labels have been generated using the SNT tool from Medtronic Surgical Navigation Technologies, Louisville, CO (see Supplementary Materials for detailed discussion of the segmentation process used).

Experiment details MAGeT-Brain was configured with an atlas library composed of the five Winterburn T1 atlases (described in Experiment 2), and size of template library, registration method, and label fusion method were determined by the optimal settings found in Experiment 1. The FSL segmentation method was used via the `run_first_all` script according to the FIRST user guide⁴. All images in the ADNI1:Complete 1Yr 1.5T dataset were segmented by both methods.

One author (MP) performed visual quality inspection for MAGeT-Brain and FreeSurfer segmentations using similar quality control guidelines (if either hippocampus was under or over segmented by 10mm or greater in three or more slices then the segmentation did not pass). Only images meeting the conditions of having segmentations from all methods (SNT, MAPER, FreeSurfer, FSL, and MAGeT-Brain) and also passing quality control inspection were included in the analysis.

Evaluation method As in Experiment 3, the SNT and Winterburn hippocampal segmentation protocols differ in the neuroanatomical features delineated, and so we assessed MAGeT-Brain by the degree of (Pearson) correlation of average hippocampal volume across subjects. We also computed the correlation in hippocampal volume between existing, established automated segmentation methods – FSL, FreeSurfer, and MAPER, and SNT semi-automated segmentations. Additionally, we evaluate the volume-related fixed and proportional biases in all segmentation methods using Bland-Altman plots (Bland and Altman, 1986).

⁴<http://fsl.fmrib.ox.ac.uk/fsl/fslwiki/FIRST/UserGuide>

Table 5: **ADNI1 1.5T Complete 1Yr dataset demographics.** CN - Cognitively Normal. LMCI - Late-onset Mild Cognitive Impairment. AD - Alzheimer’s Disease. Hisp - Hispanic. CDR-SB - Clinical Dementia Rating-Sum of Boxes. ADAS - Alzheimer’s Disease Assessment Scale. MMSE - Mini-Mental State Examination.

	N	CN <i>N</i> = 584			LMCI <i>N</i> = 931			AD <i>N</i> = 404			Combined <i>N</i> = 1919		
Age at baseline Years	1919	72.4	75.8	78.5	70.5	75.1	80.4	70.1	75.3	80.2	71.1	75.3	79.8
Sex : Female	1919	48% (278)			35% (327)			49% (198)			42% (803)		
Education	1919	14 16 18			14 16 18			12 15 17			13 16 18		
CDR-SB	1911	0.0 0.0 0.0			1.0 1.5 2.5			3.5 4.5 6.0			0.0 1.5 3.0		
ADAS 13	1895	5.67	8.67	12.33	14.67	19.33	24.33	24.67	30.00	35.33	10.67	18.00	25.33
MMSE	1917	29 29 30			25 27 29			20 23 25			25 27 29		

a b c represent the lower quartile *a*, the median *b*, and the upper quartile *c* for continuous variables.

N is the number of non-missing values.

Numbers after percents are frequencies.

3 Results

3.1 Experiment 1 Results: Whole Hippocampus Cross-Validation

In this experiment we conduct 10 rounds of MAgE-T-Brain and multi-atlas segmentation of 60 subjects using atlas and template library sizes varying from 1-9 and 1-20 images respectively, two registration algorithms (ANTS or ANIMAL), and three label fusion techniques (unweighted, cross-correlation, and normalised mutual information weighted majority vote). Hippocampal MAgE-T-Brain-based segmentations using both ANIMAL and ANTS registration algorithms demonstrate good overlap with manual segmentations (maximum mean DSC of 0.87 when using 9 atlases, 19 templates, ANTS registration, and majority vote label fusion); Figure 3). Qualitatively, both ANIMAL and ANTS-based segmentations demonstrate trend overlap accuracy that increases with the size of atlas library and template library. Improvement in accuracy plateaus with template libraries larger than ten images.

The use of MAgE-T-Brain with ANTS registration shows a pronounced increase in segmentation accuracy over MAgE-T-Brain with ANIMAL registration, across all other variable settings we tested. Additionally, by itself, using a weighted voting strategy did not significantly improve segmentation accuracy, contrary to the findings of Aljabar et al. (2009) in basic multi-atlas segmentation. Given these findings, in the remainder of this section only results using the ANTS registration algorithm and majority vote fusion will be shown.

With an increasing number of templates, MAgE-T-Brain shows improvement over multi-atlas-based segmentation in overlap accuracy when using the same number of atlases and voting method (Figure 4). The two methods converge in accuracy when using several atlases. Peak improvement in MAgE-T-Brain accuracy (0.87 DSC) is found when one atlas is used with a template library of 19 images.

In addition to an improvement in accuracy over multi-atlas-based segmentation, MAgE-T-Brain also shows a decrease in the variability of segmentation accuracy (Figure 5). The size of template library necessary to reach a decrease ($p < 0.05$) in variance and standard deviation grows with the size of atlas library used. A template library of 19 images is sufficient to show significant decrease in variance and standard deviation for 3-7 atlases.

413 We have omitted results obtained when using an even number of atlases or templates since with this
414 configuration we found significantly decreased performance. We believe this is as a result of an inherent bias
415 in the majority vote fusion method used (see Discussion).

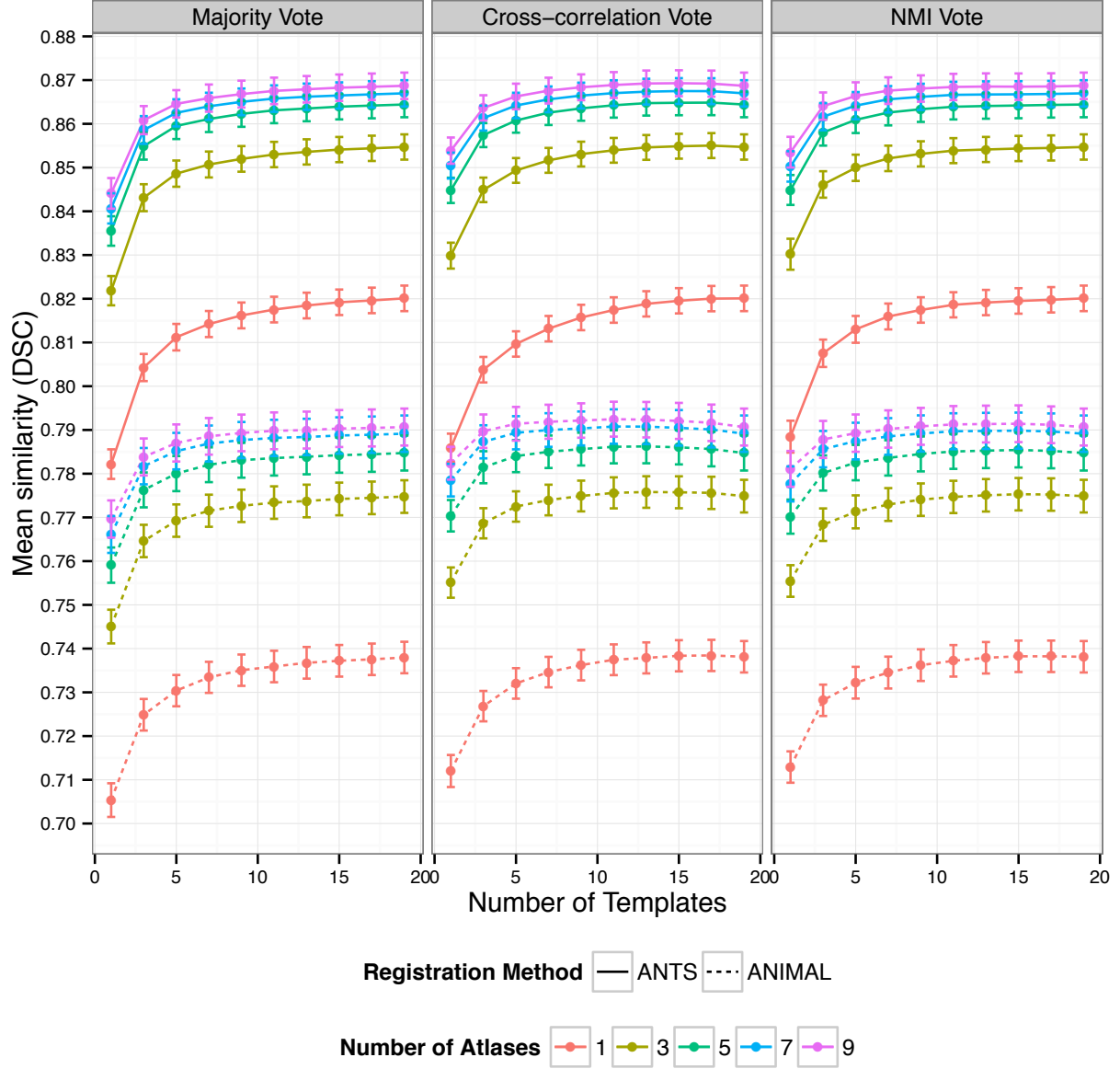


Figure 3: Mean Dice's Similarity Coefficient of MAGeT-Brain segmentations relative to Pruessner-protocol manual segmentations for 60 ADNI subjects vs. atlas and template library size, registration algorithm, and label fusion method. Error bars indicate standard error.

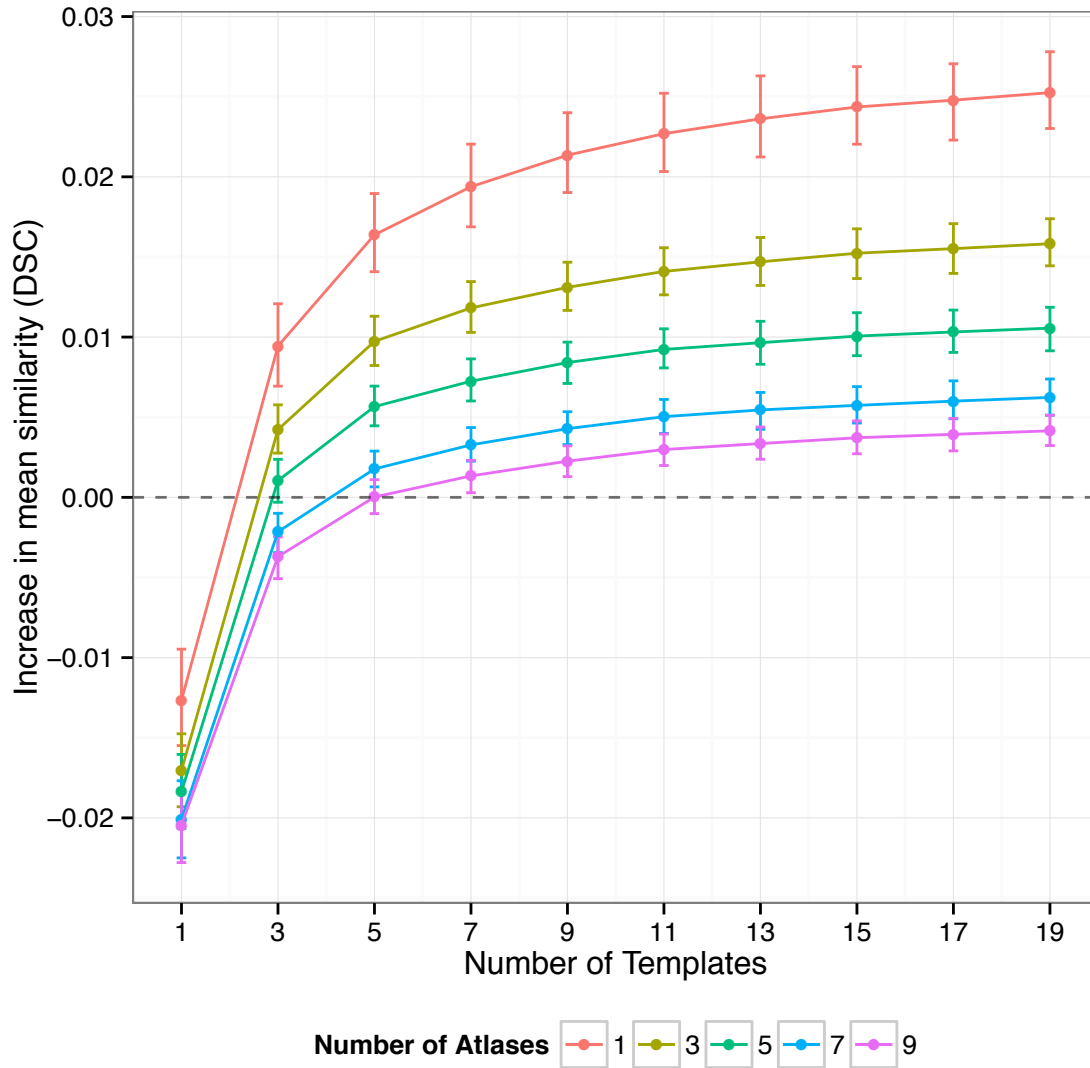


Figure 4: Increase in mean Dice's Similarity Coefficient of MAgE-T-Brain over multi-atlas segmentations vs. atlas and template library size when using the ANTS registration method, and majority-vote label fusion. Segmentation similarity is computed against Pruessner-protocol manual segmentations. Error bars indicate standard error.

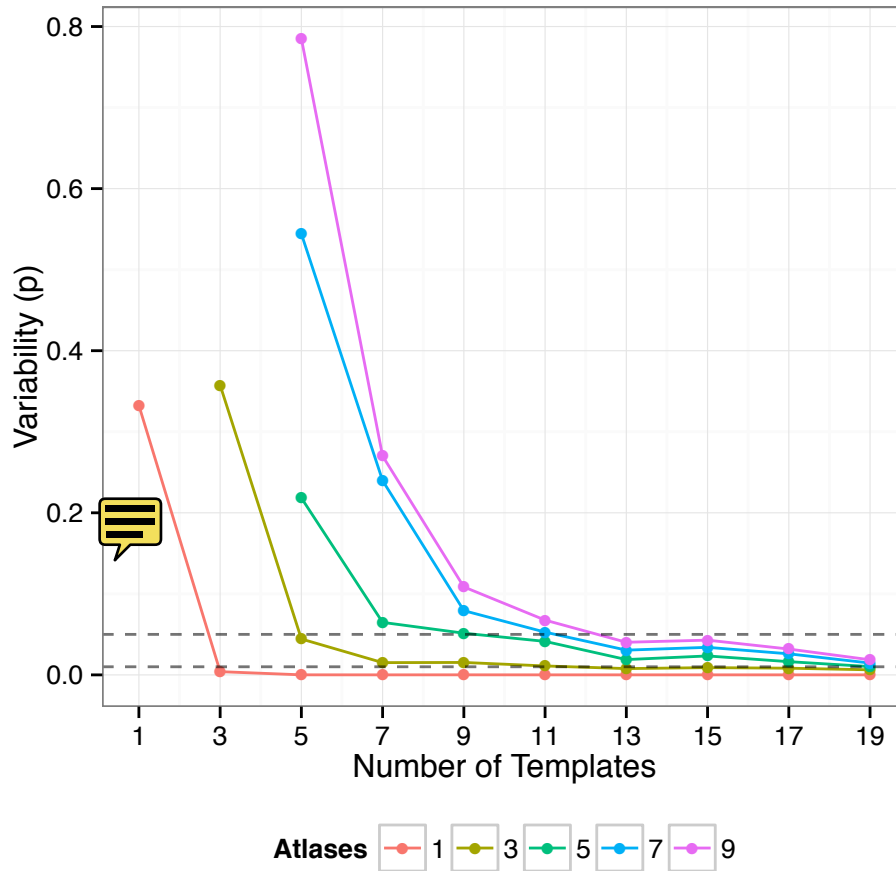


Figure 5: **Difference in variability of MAGEt-Brain vs. multi-atlassegmentation with respect to manual segmentation.** Variance of segmentation accuracy between MAGEt-Brain and multi-atlassegmentation is computed for each subject across all ten rounds of validation. Shown on the y-axis is the p-value resulting from a t-test comparing the distribution of variances at each parameter setting (atlas/template library size). Only points where MAGEt-Brain mean variability is lower than multi-atlas are shown. Dashed lines indicate a p-value of 0.05 and 0.01.

3.2 Experiment 2 Results: Winterburn Atlases Cross-Validation

This experiment explores MAGEt-Brain segmentation of hippocampal subfields. To achieve this, a leave-one-out validation was conducted in which lower-resolution images ($0.9mm^3$ voxels) of each Winterburn atlas subject were segmented using the remaining high-resolution Winterburn atlas subjects' images.

In general, across hippocampal subfields the percent error in volume of MAGEt-Brain segmentations (relative to the full-resolution Winterburn atlas segmentation) compares favourably to the error resulting from image resampling (Figure 6). In particular, there is no significant difference between percent error of MAGEt-Brain volumes compared to image resampling error for the CA1, Subiculum, and CA2/CA3 subregions. The MAGEt-Brain volumes of CA4/DG and SR/SL/SM subregions show both significantly different (lower) percent error ($p < 0.001$).

Figure 7 shows a qualitative comparison of MAGEt-Brain subfield segmentations for a single subject.

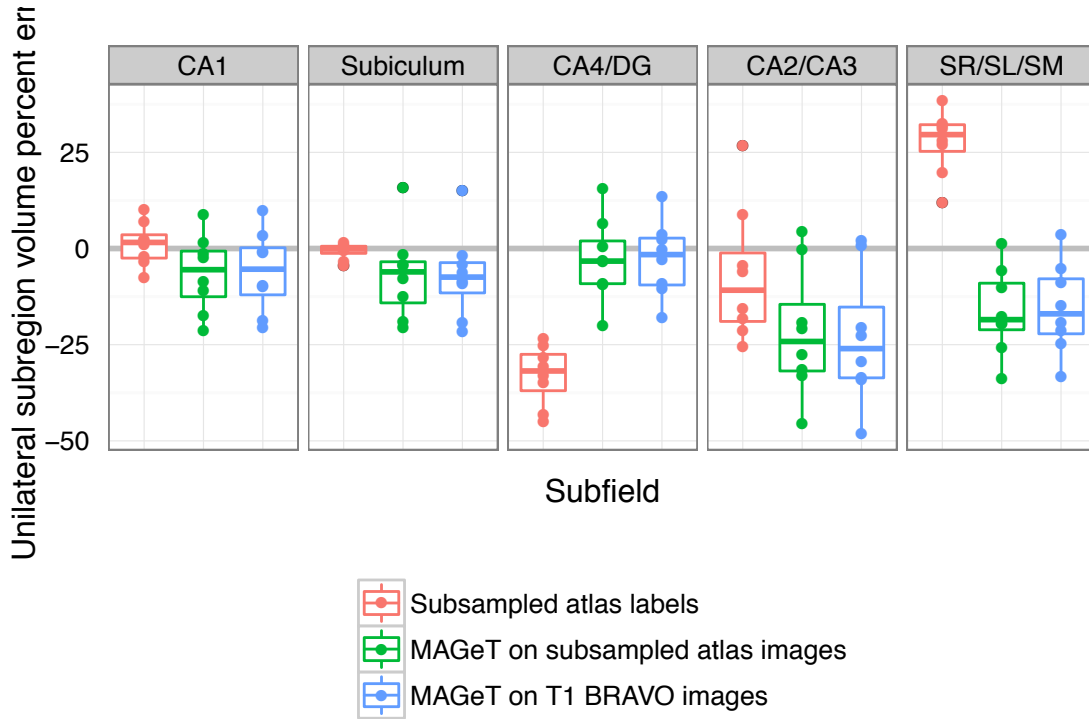
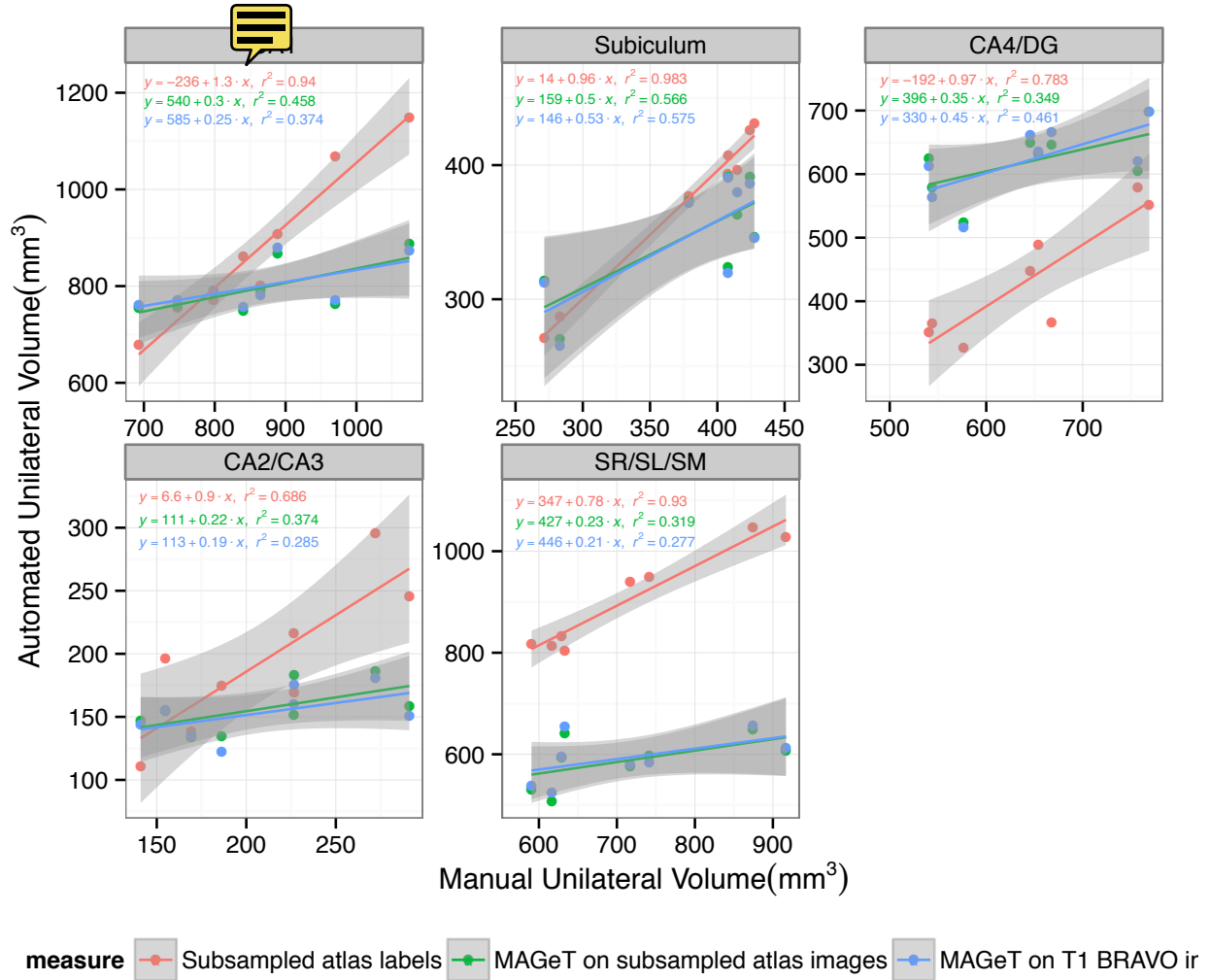


Figure 6: **Percent error in unilateral subfield volume.** Percent error is measured against the volumes of the unmodified Winterburn atlas subfield segmentations. **Subsampled atlas labels** volumes of the manual segmentations of the Winterburn atlases after resampling to $0.9mm^3$ voxels. **MAGEt on subsampled atlas images** volumes are MAGEt-Brain segmentations of the Winterburn atlas images after resampling to $0.9mm^3$ voxels. **MAGEt on T1 BRAVO** volumes are MAGEt-Brain segmentations of T1 BRAVO images ($0.9mm^3$ voxels) acquired separately of four of the five Winterburn atlas subjects.



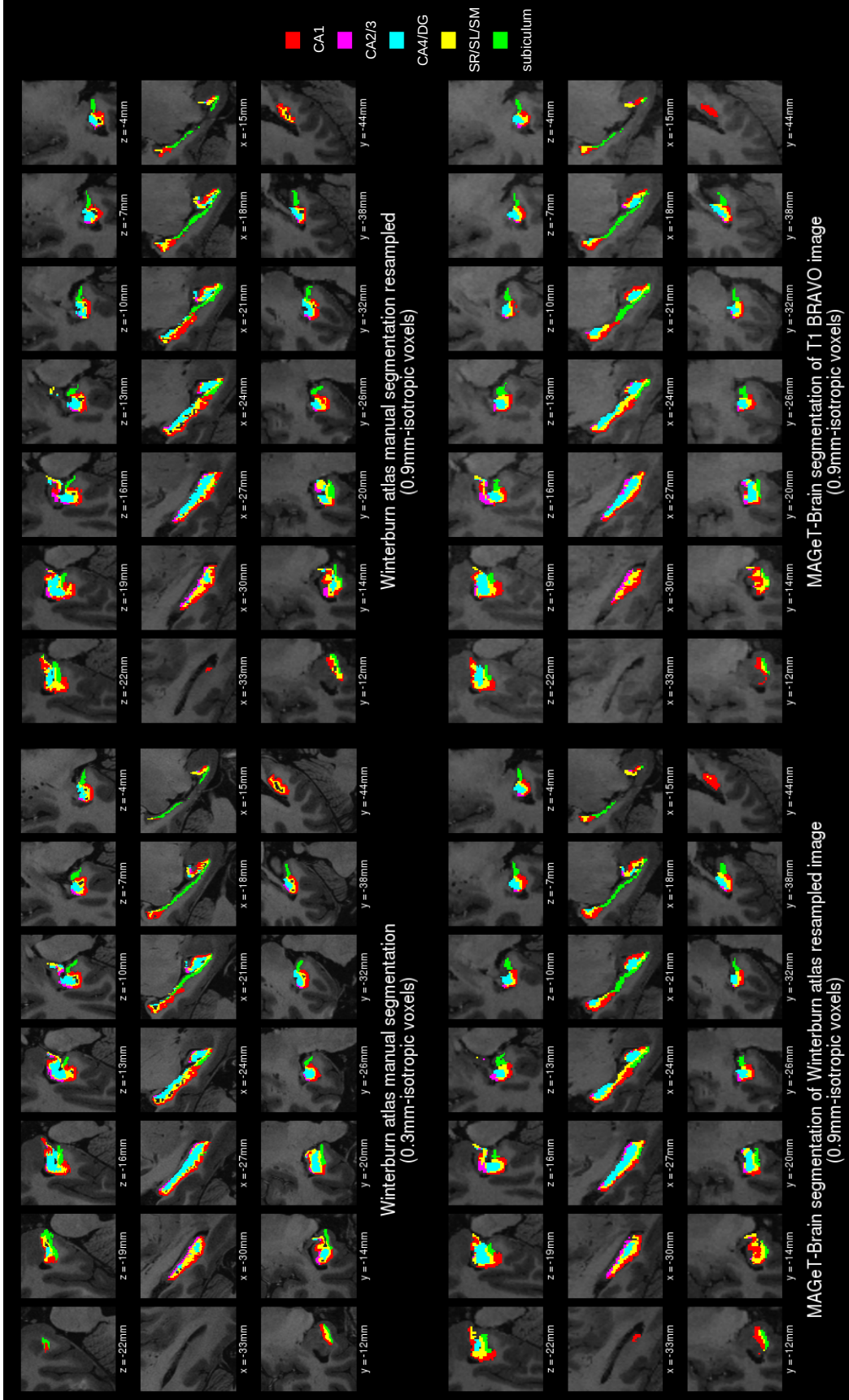


Figure 7: Detailed subfield segmentation results for a single subject. In the upper left corner is the original high-resolution Winterburn atlas manual subfield segmentation; in the upper right corner is the Winterburn atlas segmentation resampled from 0.3mm- to 0.9mm-isotropic voxels; in the lower left corner is the MAGE-T-Brain segmentation of the resampled Winterburn atlas images; in the lower right corner is the MAGE-T-Brain segmentation of a separately acquired T1 BRAVO image of the same subject. In each segmentation, slices from the left hemisphere are shown in Talairach-like ICBM152 space: the first row shows axial slices from inferior to superior; the second row shows sagittal slices from lateral to medial; the third row shows coronal slices from anterior to posterior.

3.3 Experiment 3 Results: Application to the segmentation of first episode schizophrenia patients

In this experiment MAgE-T-Brain is applied to a dataset of images of first episode schizophrenia patients, using the Winterburn atlases and a template library of 21 subject images selected at random. Expert manual whole hippocampal segmentations are used as a gold standard.

MAgE-T-Brain produces hippocampal volumes that are highly correlated with manual segmentation volumes (Pearson $r = 0.877$, $t = 16.244$, $p < 0.001$; Figure 8).

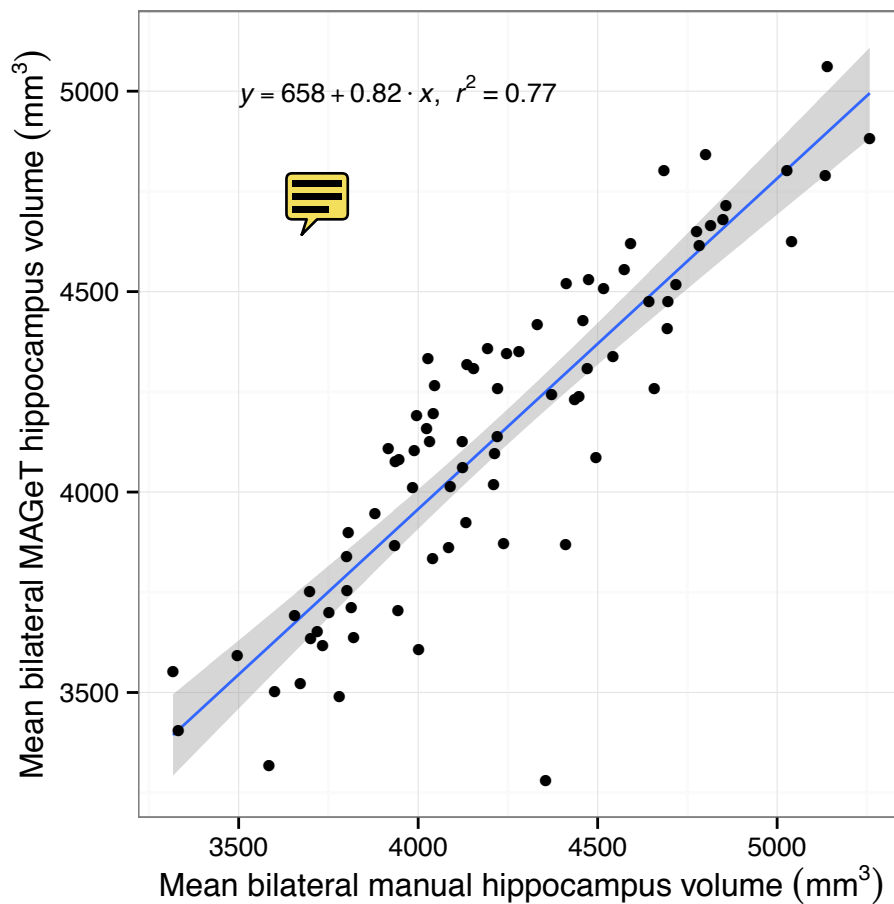


Figure 8: Mean bilateral hippocampus volume as measured by MAGeT-Brain vs. manually segmented volumes from the First Episode Patients with Schizophrenia dataset. A linear fit line is shown, with shaded region showing standard error.

Table 6: Number of segmented images and quality control failures of ADNI1:Complete 1Yr 1.5T dataset by method.

X	SNT	MAGeT	MAPER	FSL	FS
Images	368	368	368	368	368
Failures	n/a	30	n/a	20	88

3.4 Experiment 4 Results: Application to the segmentation of Alzheimer’s disease patients

Based on the results from Experiment 1, in this experiment MAGeT-Brain was configured with a template library of 21 randomly chosen subject images (7 from each disease class) and majority vote label fusion. The entire ADNI1:Complete 1Yr 1.5T dataset was segmented by MAGeT-Brain, and the resulting volumes compared with those obtained by expert semi-automated segmentation (SNT), and three other automated segmentation techniques: MAPER, FreeSurfer, and FSL. Table 6 shows the total count of segmentations available, including a count of those which have failed a quality control inspection. A total of 246 images are included in the following analysis, having met quality criteria and having segmentations from every method.

We found a close relationship in total bilateral hippocampal volume between all methods and the SNT semi-automated label volumes (Figure 10). Volumes are correlated with Pearson $r > 0.78$ for all methods across disease categories. Within disease categories (Figure 9), MAGeT-Brain is consistently well correlated to manual volumes (Pearson $r > 0.85$), but appears to slightly over-estimate the volume of the AD hippocampus.

Bland-Altman plots illustrate the level of agreement of each method with SNT segmentation hippocampal volumes (Figure 11). As Bland and Altman (1986) noted, high correlation amongst measures of the same quantity does not necessarily imply agreement (as correlation can be driven by a large range in true values, for instance). All methods show an obvious proportional bias: FreeSurfer and FSL markedly under-estimate smaller hippocampi and over-estimate large hippocampi, whereas MAPER and MAGeT-Brain show a reverse, conservative bias (Figure 11). Additionally, all methods show a fixed volume bias, with FreeSurfer and FSL most dramatically over-estimating hippocampal volume by $2600mm^3$ and $2800mm^3$ on average, respectively, and MAPER and MAGeT-Brain within $250mm^3$ on average.

Figure 12 shows a qualitative comparison of MAGeT-Brain and SNT hippocampal segmentations for 10 randomly selected subjects in each disease category, and illustrates some of the common errors found during visual inspection. Mostly frequently, we found MAGeT-Brain improperly includes the vestigial hippocampal sulcus and, although not anatomically incorrect, MAGeT-Brain under-estimates the hippocampal body in comparison to the SNT segmentation.

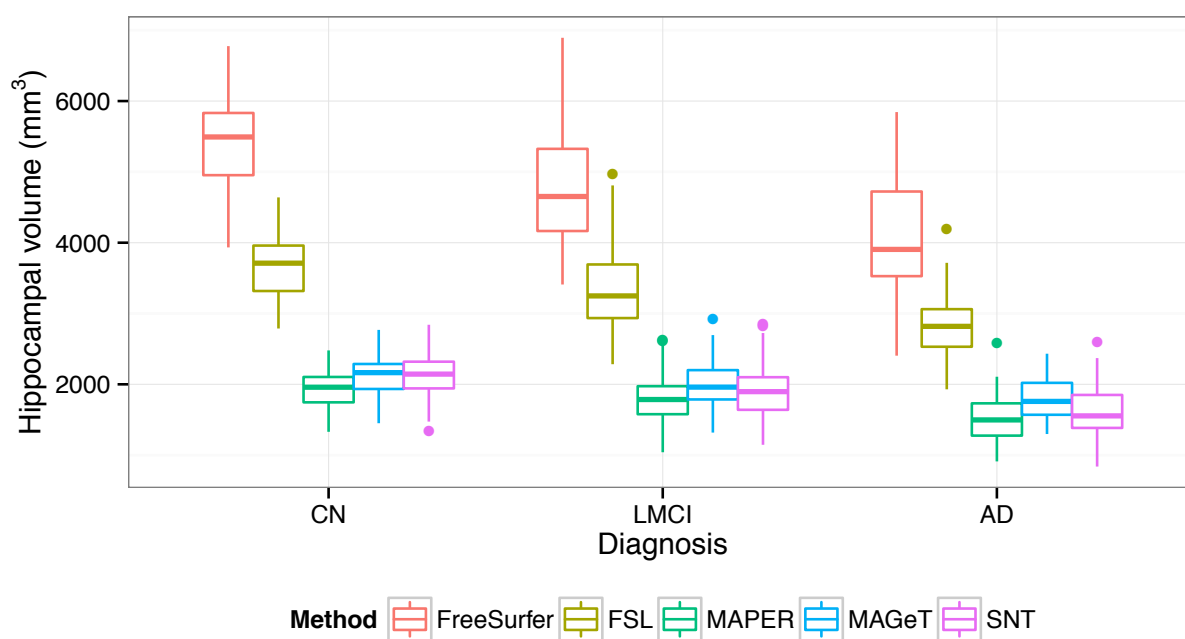


Figure 9: Subject mean hippocampal volume as measured in the ADNI1:Complete 1Yr 1.5T dataset by FreeSurfer, FSL, MAPER, MAGeT-Brain, and SNT vs. disease category.

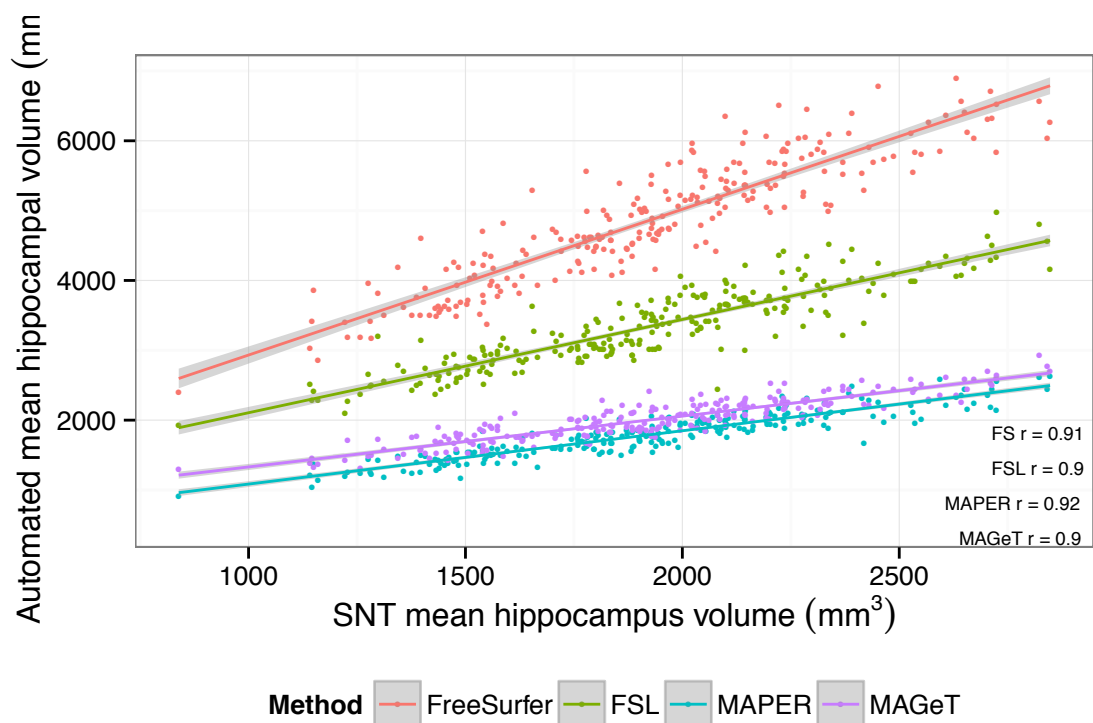


Figure 10: Subject mean hippocampal volume as measured in the ADNI1:Complete 1Yr 1.5T dataset by each of the four automated methods investigated (FreeSurfer (FS), FSL, MAPER, MAGeT-Brain) vs. SNT. Linear fit lines and Pearson correlations are shown for each method.

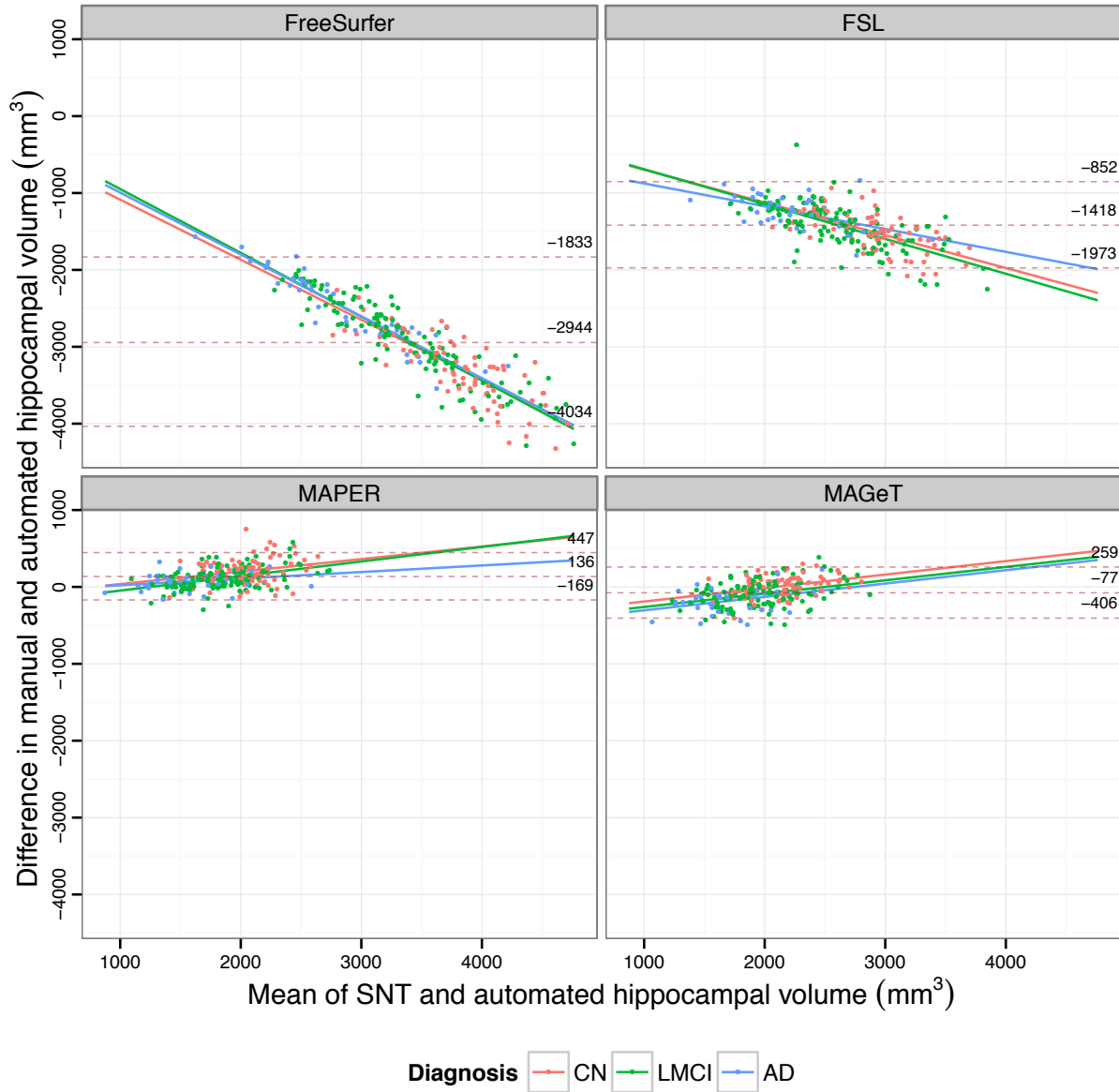


Figure 11: Bland-Altman plots comparing subject mean hippocampal volume as measured in the ADNI1:Complete 1Yr 1.5T dataset by SNT segmentation and each of the four automated methods investigated (FreeSurfer, FSL, MAPER, MAGeT-Brain). The overall mean difference in volume, and limits of agreement ($\pm 1.96SD$) are shown by dashed horizontal lines. Linear fit lines are shown for each diagnosis group. Note, points below the mean difference indicate overestimation of the volume with respect to the SNT volume, and vice versa.

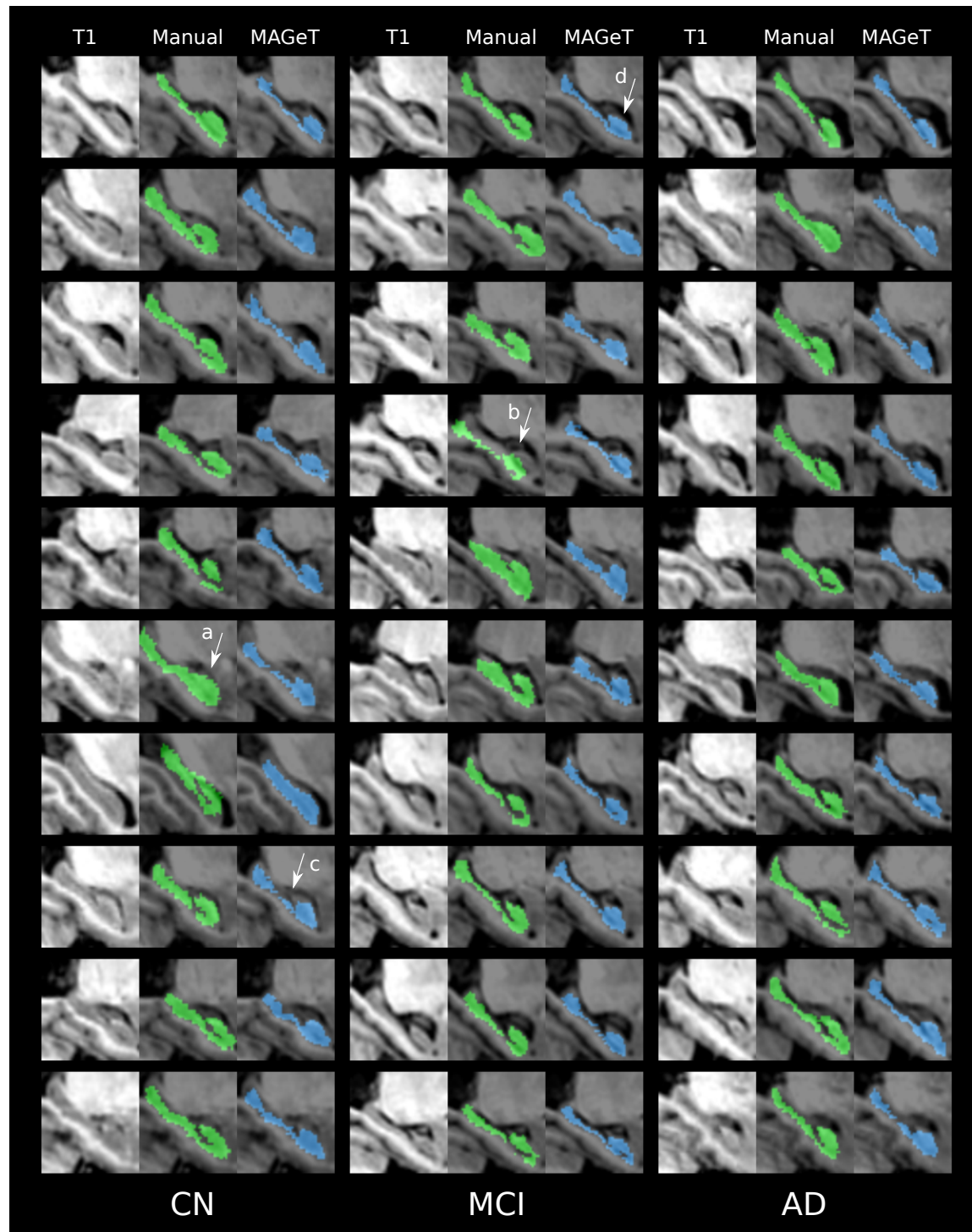


Figure 12: SNT and MAGeT-Brain segmentations for 30 ADNI subjects (10 subjects randomly selected from each disease category in the subject pool used in Experiment 1). Sagittal slices are shown for each unlabelled T1-weighted anatomical image. SNT labels appear in green, and MAGeT-Brain labels appear in blue. Noted are examples of common segmentation idiosyncrasies: (a) over-estimation of hippocampal head and (b) translated segmentation (seen in SNT segmentations only); (c) under-estimation of hippocampal body and (d) improper inclusion of the vestigial hippocampal sulcus by MAGeT-Brain.

Table 7: **Automated segmentation accuracy (overlap with SNT labels) of the ADNI dataset.** For each method, the number of labelled atlases used for training, the best Dice’s overlap measure, the disease classes measured, and the validation procedure are shown. Unless specified, validation datasets are composed equally of subjects diagnosed with Alzheimer’s Disease (AD), Mild Cognitive Impairment (MCI), and Cognitively Normal (CN). LOOCV = Leave-one-out cross-validation. Some studies of automated segmentation of ADNI images are excluded because they do not provide overlap measures for the hippocampus (Heckemann et al., 2011; Chupin et al., 2009).

Method	Atlases	DSC	Reference	Validation
MAGeT-Brain	9	0.869		10 rounds of Monte Carlo CV on a pool of 69 subjects
LEAP	30	0.848	Wolz et al. (2010)	Segmentation of 60 subjects
ACM (AdaBoost-based)	21	0.862	Morra et al. (2008)	LOOCV on atlases
Patch-based label fusion	16	0.883 (CN) 0.838 (AD)	Coupe et al. (2011)	LOOCV on atlases
Multi-atlas	30	0.885	Lötjönen et al. (2010)	Segmentation of 60 subjects
Multi-atlas + weighted fusion	20	0.898 (CN) 0.798 (left HC, MCI)	Wang et al. (2011)	10 rounds of Monte Carlo CV on 20 subjects, pool of 139 (CN/MCI)
Multi-atlas (MAPS)	55	0.890	Leung et al. (2010)	Segmentation of 30 subjects (10 AD, MCI, and CN)

4 Discussion

In this manuscript we have presented the implementation and validation of the MAGeT-Brain framework – a methodology that requires very few input atlases in order to provide accurate and reliable segmentations. Experiment 1 establishes MAGeT-Brain as effective and as comparable to multi-atlas segmentation. This experiment also rigorously characterizes the behaviour of MAGeT-Brain whole hippocampal segmentation under various parameter settings, allowing us to choose an optimal setting for subsequent experiments. Experiment 2 demonstrates the accuracy of MAGeT-Brain hippocampal subfield identification despite contrast and resolution limitations in standard T1-weighted image volumes. Experiments 2 and 3 validate MAGeT-Brain whole hippocampal segmentation accuracy and consistency on populations with different ageing and neuropsychiatric characteristics. Furthermore, taken together, these experiments demonstrate that algorithmic performance is not dependent on a single definition of the hippocampus but is effective with differing hippocampal definitions (Winterburn et al., 2013; Pruessner et al., 2000; Hsu et al., 2002).

Throughout the cross validation in Experiment 1 (10-fold Monte Carlo cross validation in the ADNI1:Complete 1Yr 1.5T dataset subsample) we find that two parameter choices improve segmentation accuracy: increasing the number of atlases, and the number of templates. However, after setting the parameters to 5 atlases and 15 templates there are diminishing returns with respect to this improvement.