

Manuscript Number: NIMG-13-2171R2

Title: Multi-atlas Segmentation of the Whole Hippocampus and Subfields Using Multiple Automatically Generated Templates

Article Type: Regular Article

Corresponding Author: Mr. Jonathan Pipitone, MSc

Corresponding Author's Institution: Centre for Addiction and Mental Health

First Author: Jonathan Pipitone, MSc

Order of Authors: Jonathan Pipitone, MSc; Min Tae M Park, BSc; Julie Winterburn, BSc; Tristram A Lett, BSc; Jason P Lerch, PhD; Jens C Pruessner, PhD; Martin Lepage, PhD; Aristotle N Voineskos, PhD, MD; Mallar M Chakravarty, PhD

**Abstract:** Introduction: Advances in image segmentation of magnetic resonance images (MRI) have demonstrated that multi-atlas approaches improve segmentation 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 agreement to multi-atlas approaches. Thus, our method acts as a reliable multi-atlas approach when using special or 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 nonlinear image registration. The resulting segmentations are then propagated to each target image and fused using a label fusion method. We conduct two separate Monte Carlo cross-validation experiments comparing MAGeT-Brain and multi-atlas whole hippocampal segmentation using differing atlas and template library sizes, and registration and label fusion methods. The first experiment is a 10-fold validation (per parameter setting) over 60 subjects taken from the Alzheimer's Disease Neuroimaging Database (ADNI), and the second is a five-fold validation over 81 subjects having had a first episode of psychosis. In both cases, automated segmentations are compared with manual segmentations following the Pruessner-protocol. Using the best settings found from these experiments, we segment 246 images of the ADNI1:Complete 1Yr 1.5T dataset and compare these with segmentations from existing automated methods: FSL FIRST, FreeSurfer, MAPER, and SNT. Finally, we conduct a leave-one-out cross-validation (LOOCV) of hippocampal subfield segmentation in standard 3T T1-weighted images, using five high-resolution manually segmented atlases (Winterburn et al., 2013).

**Results:** In the ADNI cross-validation, using 9 atlases MAGeT-Brain achieves a mean Dice's Similarity Coefficient (DSC) score of 0.869 with respect to manual whole hippocampus segmentations, and also exhibits significantly lower variability in DSC scores than multi-atlas segmentation. In the younger, psychosis dataset, MAGeT-Brain achieves a mean DSC score of 0.892 and produces volumes which agree with manual segmentation volumes better than those produced by the FreeSurfer and FSL FIRST methods (mean difference in volume: 80mm<sup>3</sup>, 1600mm<sup>3</sup>, and 800mm<sup>3</sup>, respectively). Similarly, in the ADNI1:Complete 1Yr 1.5T dataset, MAGeT-Brain produces hippocampal segmentations well correlated

( $r > 0.85$ ) with SNT semi-automated reference volumes within disease categories, and shows a conservative bias and a mean difference in volume of 250mm<sup>3</sup> across the entire dataset, compared with FreeSurfer and FSL FIRST which both overestimate volume differences by 2600mm<sup>3</sup> and 2800mm<sup>3</sup> on average, respectively. Finally, MAGeT-Brain segments the CA1, CA4/DG and subiculum subfields on standard 3T T1-weighted resolution images with DSC overlap scores of 0.56, 0.65, and 0.58, respectively, relative to manual segmentations.

Conclusion: We demonstrate that MAGeT-Brain produces consistent whole hippocampal segmentations using only 9 atlases, or fewer, with various hippocampal definitions, disease populations, and image acquisition types. Additionally, we show that MAGeT-Brain identifies hippocampal subfields in standard 3T T1-weighted images with overlap scores comparable to competing methods.



---

February 18, 2014

Editorial Board, Neuroimage

Dear Sir or Madam:

Please find enclosed our manuscript entitled, *Bootstrapping Multi-atlas Segmentation Using Multiple Automatically Generated Templates for the Segmentation of the Whole Hippocampus and Subfields* by Jon Pipitone and colleagues.

In this manuscript we describe a novel automated MRI hippocampal segmentation algorithm optimised to perform well with only a small number of manually segmented training images. We believe this is an important contribution because the expertise and effort needed to perform manual segmentation can be prohibitive for many clinicians and researchers, and yet existing automated methods generally require between 30 to 80 training segmentations. This situation makes it infeasible to use these methods for segmentations based on histological-based digital segmentations (because of their rarity), high-resolution digital atlases (because of the time needed to segment these images), or when exploring new protocols or the effect of variations on a segmentation protocol.

It is for these reasons that we developed the automated segmentation algorithm, mischievously named MAGeT-Brain, which takes advantage of the neuroanatomical variability that exists in the target population being studied to bootstrap a large template library from a small set of manually segmented images. In this manuscript we rigorously validate this approach on multiple disease populations, and compare our segmentations with existing popular methods (e.g. FSL and FreeSurfer). We find that MAGeT-Brain produces very reliable and consistent segmentations of the whole hippocampus when compared with manual segmentations, and is competitive with exiting methods. Finally, we have made our algorithm available publically online for use by other groups, and are pursuing the contribution of our segmentations to the Alzheimer's Disease Neuroimaging Initiative image database.

We believe that the technique we have developed and our findings are a significant contribution to the neuroimaging community, specifically for those



---

researchers interested in large scale studies of the hippocampus in the context of normal brain function and different forms of brain dysfunction.

We hope you find the enclosed manuscript meets the high standards of NeuroImage.

Sincerely,

Jon Pipitone  
Kimel Family Translational  
Imaging-Genetics Laboratory  
Research Imaging Centre  
Centre for Addiction and Mental  
Health  
Toronto, Ontario

encl: Manuscript

## Highlights

- We propose an automated MR image hippocampus (and subfield) segmentation method
- Our method is optimised for use with a small number (< 10) of training images
- Consistent, accurate identification of the whole hippocampus and subfields
- Validated on healthy, Alzheimer's disease, and first episode schizophrenia subjects
- Source code and high-resolution training subfield atlases available online

We would like to thank both reviewers for their careful review of our manuscript, and reviewer #1 for the detailed and insightful comments. What follows is a point-by-point response to the reviewer's comments. The reviewers' comments are in bolded text, our responses to the comments in regular weight text. To aid in assessing our modifications we have attached a marked up version of our manuscript which highlights every change made.

### **Reviewer #1**

**The manuscript by Pipitone describes a technique for automated segmentation of the hippocampus and hippocampal subfields. I reviewed this manuscript previously and this version is improved. However, problems persist, particularly in the abstract, that require attention. The authors insist on using the term accuracy (p. 1) but then within the paper, state that there is not agreement for even manual segmentation of the entire hippocampus (let alone the subfields). Thus, it seems incorrect to state that their method is accurate when the "gold-stand" is not agreed on. While the authors touch on this in some form in the manuscript, I think they need to take out instances of accuracy in the abstract (and elsewhere) and instead state that the automated segmentation generally matches quite well what would be obtained in the same subject using manual tracing. Closely corresponding, consistent, even reliable, seem like much better choices here.**

We have updated all uses of the word “accuracy” or related terms. We have used the terms that the reviewer deems to be more appropriate throughout the manuscript.

**Along these lines, the manuscript continues to suffer from issues of over use of jargon, making it hard to follow. When the authors employ precision here, it wasn't clear to me that they mean the same numbers came out during repeated testing of the same algorithm? This should probably be made clearer.**

In the revised manuscript, line 512-514 now reads:

*Instead, we must evaluate how well the lower-resolution MAGE-T-Brain hippocampal subfield segmentations correspond in form to the segmentation protocol used in the high-resolution images.*

In the discussion, on lines 567-569, we have replaced the use of the word “precision” with “agreement” and specified that we see an improvement “over repeated randomized trials”:

*... we have found that generating a template library reduces the variability in segmentation agreement (i.e. MAGE-T-brain more consistently produces segmentations in greater agreement with manual segmentations than does basic-multi-atlas method, over repeated randomized trials).*

**Another example of jargon is the sentence: "meaningfully bootstrap a template library..." This sentence is opaque to me.**

To avoid confusion, where applicable we have reworded sentences that discuss “bootstrapping” a template library and instead describe the process as “generating” a template library, as was done in our earlier paper on MAGeT-brain (Chakravarty, 2013).

Specifically with respect to the quoted sentence, we have reworded this as (lines 559-561):

*The core claim the MAGeT-Brain method is based on – that a useful template library can be generated from a small set of labelled atlas images – is validated in...*

**Even the use of bootstrap in the title is a little confusing, I suggest a title that can more clearly convey the content to a wider audience.**

We have removed the use of “bootstrap” from the title. It now reads:

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

**Finally, it wasn't clear to me that FIRST provided "radically" different definitions of the hippocampus from the plots showed here. I think there is tendency in the manuscript to overstate MAGeTs accomplishments relative to other methods. The plots support that MAGeT brain is doing better overall, but I don't see the whopping advantage to justify these kinds of statements.**

In the revised manuscript, lines 654-657, we have been more specific (and less editorial) in comparing the labels from each method:

*With the exception of FSL FIRST all methods correlate well with the semi-automated SNT volumes provided in the ADNI database. However, the FreeSurfer and FSL FIRST hippocampal segmentations are on average about twice the volume of those from all other methods.*

**Regarding the inclusion of hippocampal subfields, given the resampling necessary to do this, I wasn't really convinced that much was gained here other than that MAGeT could do this, given fairly noisy input. I suggest reframing this part, if the authors still feel strongly about keeping this section in, to talk about the subfield segmentations as a proof of concept.**

We do feel the subfield experiment is necessary as it is both novel and highlights the relevance of our method in situations where manual labels are extremely expensive (time/effort-wise). Therefore, throughout the paper we have qualified this experiment as a “proof-of-concept”.

**Overall, I think the manuscript needs to be more evenly toned with its conclusions and how it compares MAGeT to other work, and greater consideration is still required with how they treat subfield segmentation.**

We hope this concern is addressed by the modifications made above, e.g. by reframing the subfield results as proof-of-concept, simplifying language and editorializing with more concrete description.

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

Jon Pipitone<sup>1</sup>, Min Tae M. Park<sup>1</sup>, Julie Winterburn<sup>1</sup>, Tristram A. Lett<sup>1,9</sup>, Jason P. Lerch<sup>2,3</sup>, Jens C. Pruessner<sup>4</sup>, Martin Lepage<sup>4,5</sup>, Aristotle N. Voineskos<sup>1,6,9</sup>, M. Mallar Chakravarty<sup>1,6,7,8</sup> and the Alzheimer's Disease Neuroimaging Initiative\*

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

<sup>2</sup>*Neurosciences and Mental Health Laboratory, Hospital for Sick Children, Toronto, ON, Canada*

<sup>3</sup>*Department of Medical Biophysics, University of Toronto, Toronto, ON, Canada*

<sup>4</sup>*Douglas Mental Health University Institute, Verdun, QC, Canada*

<sup>5</sup>*Department of Psychiatry, McGill University, Montreal, QC, Canada*

<sup>6</sup>*Department of Psychiatry, University of Toronto, Toronto, ON, Canada*

<sup>7</sup>*Institute of Biomaterials and Biomedical Engineering, University of Toronto, Toronto, ON, Canada*

<sup>8</sup>*Rotman Research Institute, Baycrest, Toronto, ON, Canada*

<sup>9</sup>*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 (**M**ultiple **A**utomatically **G**enerated **T**emplates), for the automatic segmentation of the hippocampus that minimizes the number of atlases needed while still achieving similar agreement to multi-atlas approaches. Thus, our method acts as an ~~accurate-reliable~~ multi-atlas approach when using special ~~-or~~ 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 nonlinear image registration. The resulting segmentations are then propagated to each target image and fused using a label fusion method.

We conduct two separate Monte Carlo cross-validation experiments comparing MAGeT-Brain and multi-atlas whole hippocampal segmentation using differing atlas and template library sizes, and registration and label fusion methods. The first experiment is a 10-fold validation (per parameter setting) over 60 subjects taken from the Alzheimer's Disease Neuroimaging Database (ADNI), and the second is a five-fold validation over 81 subjects having had a first episode of psychosis. In both cases, automated segmentations are compared with manual segmentations following the Pruessner-protocol. Using

\*Data used in preparation of this article were obtained from the Alzheimer's Disease Neuroimaging Initiative (ADNI) database (adni.loni.usc.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.usc.edu/wp-content/uploads/how\\_to\\_apply/ADNI\\_Acknowledgement\\_List.pdf](http://adni.loni.usc.edu/wp-content/uploads/how_to_apply/ADNI_Acknowledgement_List.pdf)

20 the best settings found from these experiments, we segment 246 images of the ADNI1:Complete 1Yr  
21 1.5T dataset and compare these with segmentations from existing automated methods: FSL FIRST,  
22 FreeSurfer, MAPER, and SNT. Finally, we conduct a leave-one-out cross-validation (LOOCV) of hip-  
23 pocampal subfield segmentation in standard 3T T1-weighted images, using five high-resolution manually  
24 segmented atlases (Winterburn et al., 2013).

25 **Results:** In the ADNI cross-validation, using 9 atlases MAGeT-Brain achieves a mean Dice's Sim-  
26 ilarity Coefficient (DSC) score of 0.869 with respect to manual whole hippocampus segmentations, and  
27 also exhibits significantly lower variability in DSC scores than multi-atlas segmentation. In the younger,  
28 psychosis dataset, MAGeT-Brain achieves a mean DSC score of 0.892 and produces volumes which  
29 agree with manual segmentation volumes better than those produced by the FreeSurfer and FSL FIRST  
30 methods (mean difference in volume:  $80mm^3$ ,  $1600mm^3$ , and  $800mm^3$ , respectively). Similarly, in the  
31 ADNI1:Complete 1Yr 1.5T dataset, MAGeT-Brain produces hippocampal segmentations well correlated  
32 ( $r > 0.85$ ) with SNT semi-automated reference volumes within disease categories, and shows a conserva-  
33 tive bias and a mean difference in volume of  $250mm^3$  across the entire dataset, compared with FreeSurfer  
34 and FSL FIRST which both overestimate volume differences by  $2600mm^3$  and  $2800mm^3$  on average, re-  
35 spectively. Finally, MAGeT-Brain segments the CA1, CA4/DG and subiculum subfields on standard 3T  
36 T1-weighted resolution images with DSC overlap scores of 0.56, 0.65, and 0.58, respectively, relative to  
37 manual segmentations.

38 **Conclusion:** We demonstrate that MAGeT-Brain produces ~~accurate~~consistent whole hippocampal  
39 segmentations using only 9 atlases, or fewer, with various hippocampal definitions, disease populations,  
40 and image acquisition types. Additionally, we show that MAGeT-Brain identifies hippocampal subfields  
41 in standard 3T T1-weighted images with overlap scores comparable to competing methods.

#### 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

## 43 1 Introduction

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

52 The gold standard for neuroanatomical image segmentation is manual delineation by an expert human  
53 rater. However, with the availability of increasingly large MRI datasets the time and expertise required  
54 for manual segmentation becomes prohibitive (Mazziotta et al., 1995, 2001; Mazziotta et al.; Pausova et al.,  
55 2007). This effort is complicated by the fact that there is significant variation between segmentation protocols

56 with respect to specific anatomical boundaries of the hippocampus (Geuze et al., 2004) and this has led to  
57 efforts to create an unified hippocampal segmentation protocol (Jack et al., 2011; Boccardi et al., 2013b,a).  
58 In addition, there is controversy over the appropriate manual segmentation protocol to use in a particular  
59 imaging study (Nestor et al., 2012). Thus, a segmentation algorithm that can easily adapt to different  
60 manual segmentation definitions would be of significant benefit to the neuroimaging community.

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

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

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

92 However, many of these methods require significant investment of time and resources for the creation  
93 of the atlas library ranging between 30 (Heckemann et al., 2006a) and 80 (Collins and Pruessner, 2010)  
94 manually segmented atlases. This strategy has the main drawback of being inflexible as it does not easily  
95 accommodate varying the definition of the hippocampal anatomy (such as the commonly used heuristic of  
96 subdividing the hippocampus into head, body, and tail (Poppenk and Moscovitch, 2011; Pruessner et al.,  
97 2000)). Furthermore, none of these methods have demonstrated sufficient flexibility to accommodate atlases  
98 that are somehow exceptional such as those derived from serial histological data (Chakravarty et al., 2006;

---

99 Yelnik et al., 2007) or high-resolution MRI data that enables robust identification of hippocampal subfields  
100 (Winterburn et al., 2013; Yushkevich et al., 2009; Mueller and Weiner, 2009; Van Leemput et al., 2009;  
101 Wisse et al., 2012). Due to the recent availability of the latter, there has been increased interest in the  
102 use of probabilistic methods for the identification of the hippocampal subfields on standard T1-weighted  
103 images. Our group recently demonstrated that through use of an intermediary automated segmentation  
104 stage, robust and ~~accurate~~reliable segmentation of the striatum, pallidum, and thalamus using a single atlas  
105 derived from serial histological data is possible (Chakravarty et al., 2013). The novelty of this manuscript  
106 is the extension of our multi-atlas methodology to the segmentation of hippocampus. Additionally, in this  
107 paper we rigorously explore the effects of using multiple input atlases, of varying the size of the template  
108 library constructed, and registration and label fusion methods. As a result, we aim to demonstrate that it is  
109 indeed possible to reliably apply the segmentation represented in a very small set of segmented input atlases  
110 to an unlabelled target image set.

111 Of particular relevance to the present work is the LEAP algorithm (Learning Embeddings for Atlas  
112 Propagation; Wolz et al. (2010)) because of its focus on performing multi-atlas segmentation with a limited  
113 number of input atlases. The LEAP algorithm is a clever modification to the basic multi-atlas strategy in  
114 which an atlas library is grown, beginning with a set of manually labelled atlases, by successively incorpo-  
115 rating unlabelled target images once they themselves have been labelled using multi-atlas techniques. The  
116 sequence in which target images are labelled is chosen so that the similarity between the atlas images and the  
117 target images is minimised at each step, effectively allowing for deformations between very dissimilar images  
118 to be broken up into sequences of smaller deformations. Although Wolz et al. (2010) begin with an atlas  
119 library of 30 MR images, this method could theoretically work using a much smaller atlas library. In their  
120 validation, LEAP was used to segment the whole hippocampus in the ADNI1 baseline dataset, achieving a  
121 mean Dice score of 0.85 against semi-automated segmentations.

122 Also of interest to this manuscript are the methods that attempt to define hippocampal subfields using  
123 standard T1- or T2-weighted data, of which there are few. Van Leemput et al. (2009) demonstrate that  
124 the applicability of hippocampal subfield segmentation in T1-weighted images by Bayesian techniques using  
125 Markov random field shape priors learned from 10 manual segmentations. This work, available as part of  
126 the FreeSurfer package, is limited as the segmentation omits the tail of the hippocampus and the protocol  
127 has yet to be fully validated. Yushkevich et al. (2009) manually segment hippocampal subfields on high-  
128 resolution (either 0.2mm-isotropic or 0.2mm × 0.3mm × 0.2mm resolution voxels) T2-weighted MR images  
129 acquired from five post-mortem medial temporal lobe samples. Then, using nonlinear registration guided by  
130 shape-based models of the subfield segmentations ~~and~~ and manually derived hippocampus masks of the target  
131 images, the authors demonstrate accurate parcellation of hippocampal subfields, with respect to manual  
132 segmentations, in clinical 3T T1-weighted MRI volumes. Using multi-atlas with bias correction techniques,  
133 Yushkevich et al. (2010) demonstrate a semi-automated method of subfield segmentation on in vivo focal  
134 T2-weighted MR acquisitions of the temporal lobe. Manual input is only needed to mark divisions between  
135 the head, body and tail of the hippocampus on target images.

136 In this paper we describe a thorough validation of the MAGeT-Brain algorithm for the fully automatic  
137 segmentation of the hippocampus and its subfields a proof-of-concept validation of its application to the  
138 segmentation of hippocampal subfields in standard T1-weighted images. First, we address the very idea  
139 of bootstrapping generating a template library from a limited number of input atlases (Chakravarty et al.,  
140 2013) for whole hippocampus segmentation by conducting a multi-fold validation experiment over a range  
141 of atlas and template library sizes, registration and label fusion methods. This type of validation is done

---

142 first on a subset of the Alzheimer’s Disease Neuroimaging Initiative (ADNI) dataset with manual segmentations  
143 Pruessner-protocol, and then replicated on a first episode psychosis patient dataset to determine the  
144 behaviour of MAGeT-Brain when segmenting younger and differently diseased subjects. Next, we compare  
145 MAGeT-Brain with other popular segmentation algorithms (FreeSurfer, FSL FIRST, MAPER, and SNT)  
146 on all the images available in the ADNI1:Complete 1Yr 1.5T sample. Lastly, using the optimal parameter  
147 settings for MAGeT-Brain found from the previous experiments, we investigate hippocampal subfield seg-  
148 mentation by conducting a leave-one-out validation using the Winterburn et al. (2013) manually segmented  
149 high-resolution MR atlases.

## 150 2 The MAGeT-Brain Algorithm

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

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

164 The MAGeT-Brain (Multiple Automatically Generated Templates) bootstraps the creation of algorithm  
165 creates a large template library given a limited much smaller sized input atlas library , and then uses the  
166 this template library in basic multi-atlas segmentation . Images for to segment a set of input target images.  
167 The images used in the template library are selected from a set of input target the input images, either  
168 arbitrarily or so as to reflect the neuroanatomy or demographics of the target set as a whole (for instance,  
169 by sampling equally from cases and controls). The template library images are automatically labelled by  
170 each of the atlases via label propagation. Effectively, basic multi-atlas segmentation is then conducted using  
171 the template library to segment the entire set of target images (including the target images used in the  
172 construction of the template library). Since each template library image has multiple labels (one from each  
173 atlas), the final number of labels to be fused for each target may be quite large (i.e. # of atlas × # of  
174 templates).

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

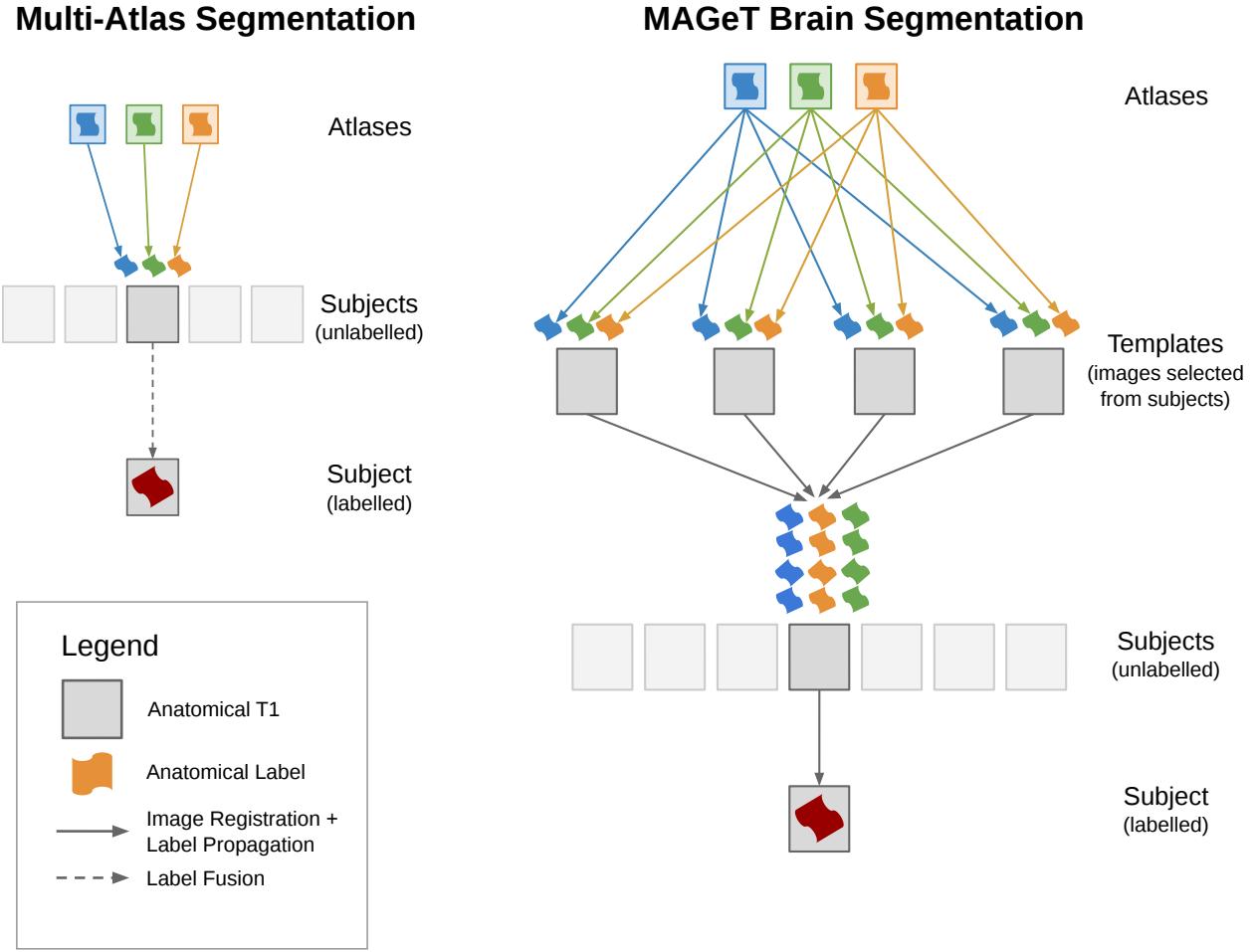


Figure 1: A schematic illustration of basic multi-atlas segmentation and MAGeT-Brain segmentation. In multi-atlas segmentation, manual labels from atlas images are warped (propagated) into subject space by applying the transformations estimated from nonlinear image registration. The resulting candidate labels from all atlas images are then fused to create a final segmentation. In MAGeT-Brain segmentation, a template library is created by sampling (either randomly or representatively) from the subject images. Atlas labels are propagated to all template images and then to each subject image (including those used in the template library). The candidate labels for a subject are then fused into a final segmentation.

---

### 177 3 Experiments

178 The following section describes experiments conducted to assess the segmentation quality of the MAGeT-  
179 Brain algorithm:

- 180 • Experiment 1 investigates MAGeT-Brain whole hippocampus segmentation of aging and Alzheimer’s  
181 diseased subjects over a wide range of parameter settings using a Monte Carlo cross-validation design.  
182 The results of this experiment enable us to choose the parameter settings offering the best performance  
183 for use in subsequent experiments.
- 184 • Experiment 2 is a similar cross-validation to explore MAGeT-Brain segmentations on the brain im-  
185 ages of young, first episode psychosis patients. In addition, MAGeT-Brain segmentations with two  
186 different atlas segmentation protocols are compared to automated segmentations by the FSL FIRST  
187 and FreeSurfer algorithms. The results of this experiment combined with the previous experiment  
188 establishes parameter settings that do not overfit to the neuroanatomical features of a specific patient  
189 cohort.
- 190 • Experiment 3 bridges MAGeT-Brain with the existing segmentation literature by comparing MAGeT-  
191 Brain whole hippocampus segmentations with those of several well-known automated methods (FreeSurfer,  
192 FSL FIRST, MAPER, SNT) on the entire ADNI1:Complete 1Yr 1.5T image dataset consisting of  
193 246 brain images of subjects diagnosed as cognitively normal, having mild cognitive impairment, or  
194 Alzheimer’s disease.
- 195 • Experiment 4 assesses hippocampal subfield segmentation quality in a leave-one-out cross-validation  
196 on the five high-resolution manually segmented Winterburn MR atlases (Winterburn et al., 2013).

#### 197 3.1 Experiment 1: Whole Hippocampus Segmentation Cross-Validation — Alzheimer’s 198 Disease

199 In this experiment we explore the very idea of ~~bootstrapping~~ generating a template library for multi-atlas-  
200 based segmentation from a small number of input atlases. To do so, we conduct repeated cross-validations  
201 of MAGeT-Brain whilst varying the composition and sizes of the atlas and template libraries used, as well  
202 as varying the registration algorithm and label fusion method. The dataset used in this experiment is  
203 images from the ADNI dataset (Jack et al., 2008) along with whole hippocampus labels manually segmented  
204 following the Pruessner-protocol (Pruessner et al., 2000).

205 Note, in the Supplementary Materials we have replicated this experiment using the SNT semi-automated  
206 segmentations included as part of the ADNI dataset.

##### 207 3.1.1 Experiment 1: Materials and Methods

208 **ADNI1:Complete 1Yr 1.5T dataset** Data used in the preparation of this article were obtained from the  
209 Alzheimer’s Disease Neuroimaging Initiative (ADNI) database ([adni.loni.usc.edu](http://adni.loni.usc.edu)). The ADNI was launched  
210 in 2003 by the National Institute on Aging (NIA), the National Institute of Biomedical Imaging and Bio-  
211 engineering (NIBIB), the Food and Drug Administration (FDA), private pharmaceutical companies and  
212 non-profit organizations, as a \$60 million, 5-year public-private partnership. The primary goal of ADNI has  
213 been to test whether serial magnetic resonance imaging (MRI), positron emission tomography (PET), other

Table 1: **ADNI1 cross-validation subset demographics.** CN - Cognitively Normal. LMCI - Late-onset Mild Cognitive Impairment. AD - Alzheimer’s Disease. CDR-SB - Clinical Dementia Rating-Sum of Boxes. ADAS - Alzheimer’s Disease Assessment Scale. MMSE - Mini-Mental State Examination.

	CN N = 20	LMCI N = 20	AD N = 20	Combined 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.

214 biological markers, and clinical and neuropsychological assessment can be combined to measure the progression  
 215 of mild cognitive impairment (MCI) and early Alzheimer’s disease (AD). Determination of sensitive and  
 216 specific markers of very early AD progression is intended to aid researchers and clinicians to develop new  
 217 treatments and monitor their effectiveness, as well as lessen the time and cost of clinical trials.

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

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

232 Clinical, demographic and pre-processed T1-weighted MRI were downloaded by the authors from the  
 233 ADNI database ([adni.loni.usc.edu](http://adni.loni.usc.edu)) between March 2012 and August 2012. The image dataset used was  
 234 the *ADNI1:Complete 1Yr 1.5T* standardized dataset available from ADNI <sup>1</sup> (Wyman et al., 2012). This  
 235 image collection contains uniformly pre-processed images which have been designated to be the “best” after  
 236 quality control. All images were acquired using 1.5T scanners (General Electric Healthcare, Philips Medical  
 237 Systems or Siemens Medical Solutions) at multiple sites using the protocol described in Jack et al. (2008).  
 238 Representative 1.5T imaging parameters were TR = 2400ms, TI = 1000ms, TE = 3.5ms, flip angle = 8°,  
 239 field of view = 240 x 240mm, a 192 × 192 × 166 matrix ( $x$ ,  $y$ , and  $z$  directions) yielding voxel dimensions of  
 240 1.25mm × 1.25mm × 1.2mm.

241 **Experiment details** Monte Carlo Cross-Validation (MCCV), also known as repeated random sub-sampling  
 242 cross-validation, consists of repeated rounds of validation conducted on a fixed dataset (Shao, 1993). In each

<sup>1</sup><http://adni.loni.usc.edu/methods/mri-analysis/adni-standardized-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

243 round, the dataset is randomly partitioned into a training set and a validation set. The method to be  
 244 validated is then given the training data, and its output is compared with the validation set.

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

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

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

259 **Automatic Normalization and Image Matching and Anatomical Labeling (ANIMAL)** The  
 260 ANIMAL algorithm carries out image registration in two phases. In the first, a 12-parameter linear trans-  
 261 formation (3 translations, rotations, scales, shears) is estimated between images using an algorithm that  
 262 maximizes the correlation between blurred MR intensities and gradient magnitude over the whole brain  
 263 (Collins et al., 1994). In the second phase, nonlinear registration is completed using the ANIMAL algorithm  
 264 (Collins et al., 1995): an iterative procedure that estimates a 3D deformation field between two MR images.  
 265 At first, large deformations are estimated using a blurred version of the input data. These larger deforma-  
 266 tions are then input to subsequent steps where the fit is refined by estimating smaller deformations on data  
 267 blurred with a Gaussian kernel with a smaller full width at half maximum (FWHM). The final transfor-  
 268 mation is a set of local translations defined on a bed of equally spaced nodes that were estimated through  
 269 the optimization of the correlation coefficient. For the purposes of this work we used the regularization  
 270 parameters optimized in Robbins et al. (2004), displayed in Table 2.

271 **Automatic Normalization Tools (ANTs)** ANTS is a diffeomorphic registration algorithm which  
 272 provides great flexibility over the choice of transformation model, objective function, and the consistency of  
 273 the final transformation (Avants et al., 2008). The transformation is estimated in a hierarchical fashion where  
 274 the MRI data is subsampled, allowing large deformations to be estimated and successively refined at later  
 275 hierarchical stages (where the data is subsampled to a finer grid). The deformation field and the objective

276 function are regularized with a Gaussian kernel at each level of the hierarchy. The ANTS algorithm is  
 277 freely available <http://www.picsl.upenn.edu/ANTS/>. We used an implementation of the ANTS algorithm  
 278 compatible with the MINC data format, mincANTS <https://github.com/vfonov/mincANTS>.

279 We used the following command line when running ANTS:

```
280 mincANTS 3 -m PR[target_file.mnc,source_file.mnc,1,4]
281 --number-of-affine-iterations 10000x10000x10000x10000x10000
282 --affine-gradient-descent-option 0.5x0.95x1.e-4x1.e-4
283 --use-Histogram-Matching --MI-option 32x16000
284 -r Gauss[3,0] -t SyN[0.5] -i 100x100x100x20
285 -o transformation.xfm
286
```

287 These settings were adapted from the "reasonable starting point" given in the ANTS manual <sup>2</sup>.

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

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

293 **Cross-correlation Weighted Majority Vote** An optimal combination of targets from the template li-  
 294 brary has previously been shown to improve segmentation accuracy [with respect to manual segmentations](#)  
 295 (Aljabar et al., 2009; Collins and Pruessner, 2010). In this method, each template library image is  
 296 ranked in similarity to each unlabelled image by the normalized cross-correlation (CC) of image intensi-  
 297 ties after linear registration, over a region of interest (ROI) generously encompassing the hippocampus.  
 298 Only the top ranked template library image labels are used in a voxel-wise majority vote. The ROI  
 299 is heuristically defined as the extent of all atlas labels after linear registration to the template, dilated  
 300 by three voxels (Chakravarty et al., 2013). The number of top ranked template library image labels is  
 301 a configurable parameter and displayed as the size of the template library in the rest of the paper.

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

304 **Normalised Mutual Information Weighted Majority Vote** This method is similar to cross-correlation  
 305 weighted voting except that image similarity is calculated by the normalised mutual information score  
 306 over the region of interest (Studholme et al., 2001). The `itk_similarity` utility from the EZMinc  
 307 toolkit<sup>3</sup> is used to calculate the normalised mutual information measure between two images.

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

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

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

<sup>3</sup><https://github.com/vfonov/EZminc>

311 where  $A$  and  $B$  are the regions being compared, and the cardinality is the volume measured in voxels. The  
312 labels produced by MAGeT-Brain and multi-atlas segmentation are compared to the manual labels using  
313 the Dice similarity coefficient, and the recorded value for each subject at each parameter setting explored in  
314 this experiment is the average over ten validation rounds.

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

321 **3.1.2 Experiment 1: Results**

322 We find that for MAGeT-Brain segmentations, similarity score increases as atlas and template library size  
323 is increased, although with diminishing returns and an eventual trend towards a plateau (Figure 2a). For  
324 instance, with 9 atlases and using ANTS for registration and majority vote fusion, the mean DSC scores for  
325 1, 5, 9 and 17 templates are 0.845, 0.865, 0.867, 0.869, respectively. A maximum similarity score of 0.869 is  
326 found when using 9 atlases, 19 templates, ANTS registration, and majority vote label fusion.

327 The ANTS registration method consistently outperforms ANIMAL registration over all variable settings  
328 we tested (mean increase in DSC is 0.079). Pearson correlations of MAGeT-Brain DSC scores when using  
329 weighted voting and when using non-weighted majority vote label fusion (with ANTS registration) for all  
330 combinations of atlases and templates are  $r > 0.899$ ,  $p < 0.001$ , with a mean difference in DSC score of 0.002.  
331 This result suggests that using a weighted voting strategy does not significantly improve MAGeT-Brain  
332 segmentation agreement, contrary to the findings of Aljabar et al. (2009) for basic multi-atlas segmentation.  
333 Thus, in the remainder of our experiments only results using the ANTS registration algorithm and majority  
334 vote fusion will be shown.

335 With at least five templates, MAGeT-Brain consistently shows a higher DSC score than multi-atlas  
336 segmentation wth the same number of atlases:  $r = 0.94$ ,  $p < 0.001$ , mean DSC increase = 0.008 (Figure 2b).  
337 The magnitude of DSC increase grows with template library size but shows diminishing returns with larger  
338 atlas libraries. Peak increase (+0.025 DSC) is found with a single atlas and template library of 19 images.

339 In addition to a mean increase in similarity score over multi-atlas-based segmentation, MAGeT-Brain also  
340 shows more consistency in similarity scores across all subjects and validation folds (Figure 2c). A template  
341 library of at least 13 images is sufficient to show significant ( $p < 0.05$ ) decrease in variance for all sizes of  
342 atlas library tested (1-9 images).

343 We find similar behaviour with respect to optimal parameter settings and increased consistency of  
344 MAGeT-Brain segmentations in the replication of this experiment (Experiment 5, Supplementary Mate-  
345 rials) where a different hippocampal definition is used (SNT labels available with the ADNI datasets). This  
346 strongly suggests that these results are independent of the segmentation protocol used and are, instead,  
347 features of the MAGeT-Brain algorithm.

348 We have omitted results obtained when using an even number of atlases or templates since with these  
349 configurations we found significantly decreased performance. We believe this results from an inherent bias  
350 in the majority vote fusion method used (see Discussion).

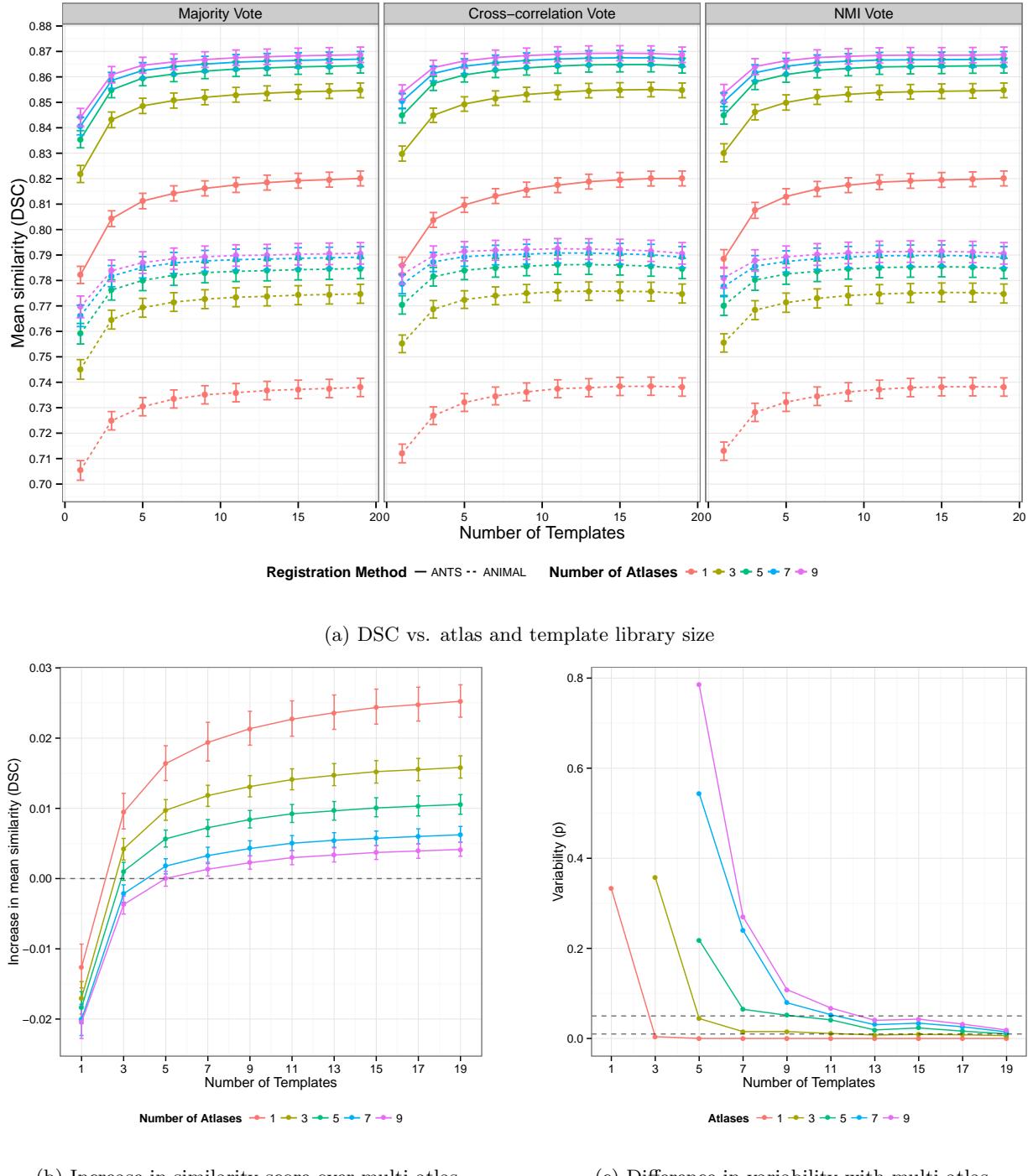


Figure 2: **Whole hippocampus segmentation cross-validation on ADNI subjects with Pruessner-protocol manual segmentations.** (2a) Average DSC score of MAGeT-Brain with manual segmentations for 60 ADNI subjects taken over 10 folds of cross-validation at each parameter setting. Error bars indicate standard error. (2b) Increase in DSC of MAGeT-Brain over multi-atlas segmentations. (2c) shows the significance of t-tests comparing the variability in DSC scores of MAGeT-Brain and multi-atlas across validation folds. Only points where MAGeT-Brain mean variability is lower than multi-atlas are shown. Dashed lines indicate p-values of 0.05 and 0.01.

351 **3.2 Experiment 2: Whole Hippocampus Segmentation Cross-Validation — First**  
352 **Episode of Psychosis**

353 To validate that the MAGeT-Brain works effectively in the context of other neurological disorders, in this  
354 experiment we replicate the cross-validation done in Experiment 1 with a dataset of patients having had a  
355 single episode of psychosis. We also compare MAGeT-Brain segmentations with those of two well-known  
356 automated segmentation methods, FSL FIRST and FreeSurfer.

357 **3.2.1 Experiment 2: Materials and Methods**

358 **First Episode Psychosis (FEP) Dataset** All patients were recruited and treated through the Prevention  
359 and Early Intervention Program for Psychoses (PEPP-Montreal), a specialized early intervention service at  
360 the Douglas Mental Health University Institute in Montreal, Canada. People aged 14 to 35 years from the  
361 local catchment area suffering from either affective or non-affective psychosis who had not taken antipsychotic  
362 medication for more than one month with an IQ above 70 were consecutively admitted as either in- or out-  
363 patients. Of those treated at PEPP, only patients aged 18 to 30 years with no previous history of neurological  
364 disease or head trauma causing loss of consciousness were eligible for the neuroimaging study; only those  
365 suffering from schizophrenia spectrum disorders were considered for this analysis. For complete program  
366 details see Malla et al. (2003).

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

372 Expert whole hippocampal manual segmentation of each subject is produced following a validated seg-  
373 mentation protocol (Pruessner et al., 2000).

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

379 **Experiment details** The same overall design as Experiment 1 is followed in this experiment: a Monte  
380 Carlo cross-validation (MCCV) is conducted using the pool of 81 first episode psychosis subject brain images  
381 and corresponding Pruessner-protocol manual segmentations. Five rounds of validation are conducted for  
382 each subject, and each atlas and template library size combination (1-9 atlases, 1-19 templates). In each  
383 round, images and their manual labels are randomly selected from the pool, and the remaining images are  
384 segmented using MAGeT-Brain with a random subset of the unlabelled images also serving as template  
385 images. Majority vote fusion, and the ANTS registration algorithm are used, as these have shown to behave  
386 favourably in previous experiments.

387 In addition to the MCCV, we segment the entire first episode psychosis dataset using MAGeT-Brain using  
388 two different atlases, as well as two popular automated segmentation packages, FSL FIRST and FreeSurfer.  
389 Specifically, MAGeT-Brain is run once with the five Winterburn atlas images and labels as atlases and a

Table 3: **First Episode Psychosis Subject Demographics.** ambi - ambidextrous. SES - Socioeconomic Status score. FSIQ - Full Scale IQ.

	N	FEP N = 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.

390 randomly selected subset of 19 target images as templates. MAGeT-Brain is run a second time using the  
 391 same template images, but we using five additional first episode psychosis subjects and corresponding manual  
 392 segmentations (not included above) as atlases. FSL FIRST and FreeSurfer are run with the default settings:  
 393 FSL FIRST `run_first_all` script was used according to the FIRST user guide <sup>4</sup>, and FreeSurfer was run  
 394 with the command `recon-all -all`.

395 **Evaluation method** Manual and automated segmentations are directly compared using Dice's similarity  
 396 coefficient (DSC). In the MCCV, the per-subject DSC value is computed as the average value over the five  
 397 rounds of validation for a given atlas and template library size. The reported average DSC value per given  
 398 atlas and template library size is the average DSC value over all subjects segmented.

399 The Pruessner segmentation protocol differs slightly from the Winterburn protocol, and those used by  
 400 FreeSurfer and FSL FIRST, in the inclusion of neuroanatomical features and the manner they are delineated (see Winterburn et al. (2013), and Table 9 in the Discussion below). This variation in protocol poses  
 401 a problem if an overlap measure is used for evaluation: since different protocols will necessarily produce  
 402 segmentations that do not perfectly overlap, the degree of overlap cannot be solely used to compare segmen-  
 403 tation methods using different protocols. In place of an overlap metric, we assess the degree of (Pearson)  
 404 correlation in average bilateral hippocampal volume produced by each method. Additionally, we evaluate the  
 405 volume-related fixed and proportional biases in all segmentation methods using Bland-Altman plots (Bland  
 406 and Altman, 1986).

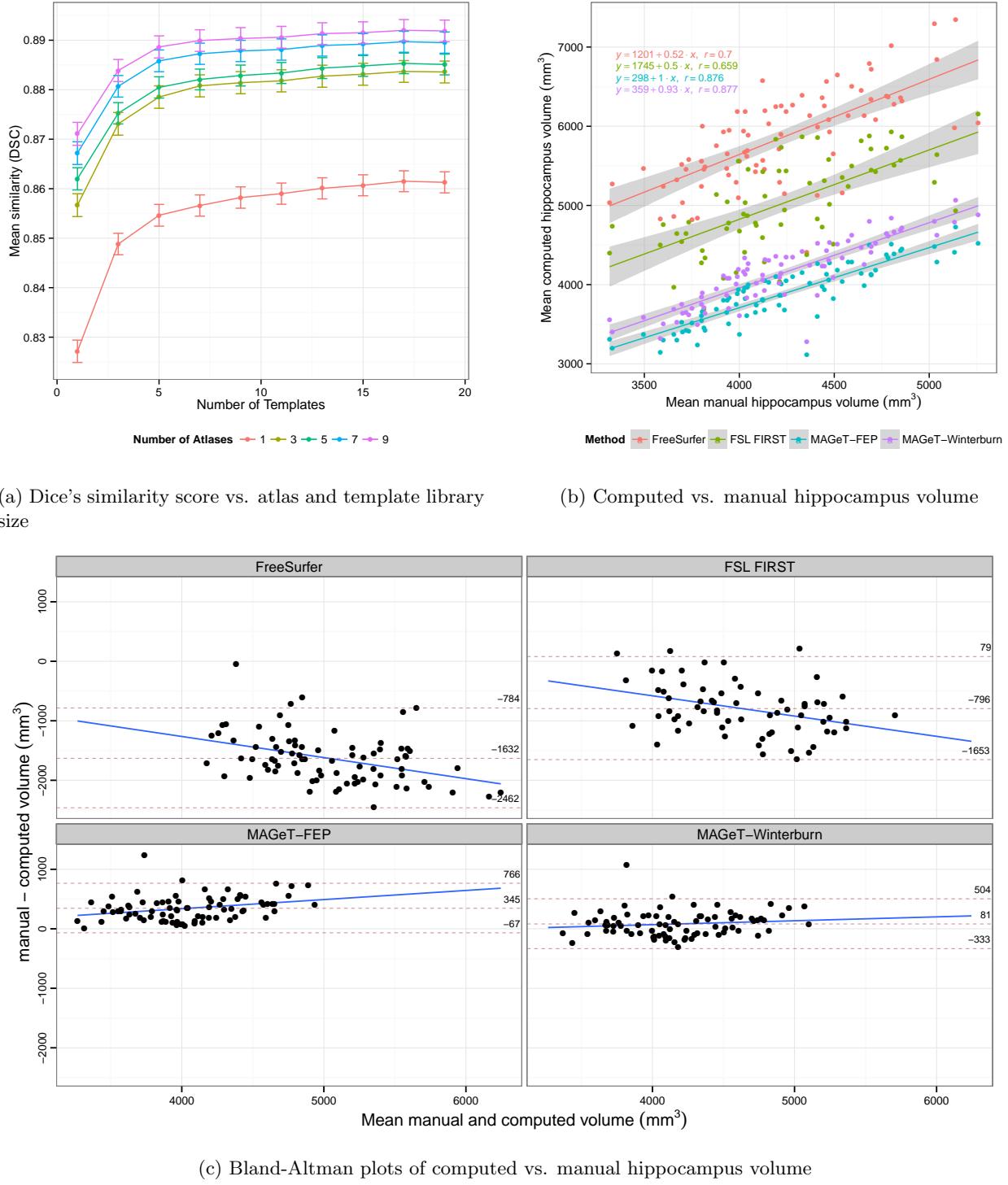
### 408 3.2.2 Experiment 2: Results

409 As in Experiment 1, we find that similarity score increases with a greater number of atlases or templates  
 410 but quickly plateaus (Figure 3a). A maximum similarity score of 0.892 is found when using 9 atlases, 19  
 411 templates, ANTS registration, and majority vote label fusion.

412 We found a close relationship in average hippocampal volume between the manual label volumes and  
 413 MAGeT-Brain when using the Winterburn atlases, or manually segmented FEP subjects as atlases (Figure  
 414 3b). Both sets of volumes are correlated with Pearson  $r > 0.88$ . FreeSurfer and FSL FIRST volumes are  
 415 both correlated with manual volumes at Pearson  $r > 0.7$ .

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

416 As Bland and Altman (1986) noted, high correlation amongst measures of the same quantity does not  
417 necessarily imply agreement (as correlation can be driven by a large range in true values, for instance).  
418 Figure 3c shows Bland-Altman plots illustrating the level of agreement of each method with manual vol-  
419 umes. All methods show an obvious proportional bias: FreeSurfer and FSL FIRST markedly underestimate  
420 smaller hippocampi and over-estimate large hippocampi (the limits of agreement are between  $-2482mm^3$   
421 and  $-784mm^3$ , and between  $-1653mm^3$  and  $79mm^3$ , respectively), whereas both MAGeT-Brain methods  
422 show a much less exaggerated, but conservative bias (limits of agreement between  $-67mm^3$  and  $766mm^3$   
423 when using FEP atlases, and between  $-333mm^3$  and  $504mm^3$  when using Winterburn atlases). On average,  
424 FreeSurfer and FSL FIRST overestimate hippocampal volume by about  $1600mm^3$  and  $800mm^3$ , respec-  
425 tively. In contrast, on average MAGeT-Brain underestimates volumes by about  $300mm^3$  when using FEP  
426 atlases and by about  $80mm^3$  when using Winterburn atlases (compared to the Pruessner-protocol manual  
427 segmentations).



**Figure 3: First Episode Patient dataset validation.** All manual segmentation of the 81 subjects is done with the Pruessner-protocol. MAGeT-Brain uses ANTS registration and majority vote label fusion. (3a) shows mean DSC score of MAGeT-Brain segmentations, as atlas and template library size is varied over a 5-fold validation. Error bars indicate standard error. (3b) shows segmentation volumes from FSL FIRST, FreeSurfer, MAGeT-Brain using the five Winterburn atlases (MAGeT-Winterburn), and MAGeT-Brain using five manually segmented FEP subjects as atlases (MAGeT-FEP). Linear fit lines are shown, with the shaded region showing standard error. (3c) shows the agreement between computed and manually volumes. 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 manual volume, and vice versa.

428 **3.3 Experiment 3: Whole Hippocampus Segmentation Comparison — ADNI1**  
429 **Complete 1Yr**

430 To validate MAGeT-Brain segmentation quality with respect to other established automated hippocampal  
431 segmentation methods, we apply MAGeT-Brain to a large dataset from the ADNI project. The resulting seg-  
432 mentations are compared to those produced by FreeSurfer, FSL FIRST, MAPER, as well as semi-automated  
433 whole hippocampal segmentations (SNT) provided by ADNI.

434 **3.3.1 Experiment 3: Materials and Methods**

435 **ADNI1:Complete 1Yr 1.5T dataset** The *ADNI1:Complete 1Yr 1.5T* standardized dataset contains  
436 1919 images in total. SNT, MAPER, and FreeSurfer hippocampal volumes for a subset of images were  
437 provided by ADNI, along with quality control data for each FreeSurfer segmentation (guidelines described  
438 in (Hartig et al., 2010)). See Section 3.1.1 for study details, inclusion criteria and imaging characteristics.

439 For a subset of the ADNI images, semi-automated segmentations of the left and right whole hippocampi  
440 generated using the SNT tool from Medtronic Surgical Navigation Technologies, Louisville, CO (see Sup-  
441 plementary Materials for detailed discussion of the segmentation process) are made available (Hsu et al.,  
442 2002). These labels are used as the reference labels in several other studies of (semi-)automated segmenta-  
443 tion methods (see Discussion). In addition, ADNI also distributes hippocampal segmentations and volumes  
444 determined using MAPER (Heckemann et al., 2011), a multi-atlas segmentation tool, and the FreeSurfer  
445 tool (including quality control data, with guidelines described in Hartig et al. (2010)).

446 **Experiment details** MAGeT-Brain was configured with an atlas library composed of the five Winterburn  
447 atlas images (Experiment 2, section 3.2) and segmentations. A template library of 19 images were randomly  
448 selected from the target dataset of ADNI subjects, and ANTS registration and majority vote label fusion  
449 were used as these were found to perform favourably in earlier experiments.

450 FSL FIRST segmentation was performed using the `run_first_all` script according to the FIRST user  
451 guide <sup>5</sup>. All images in the ADNI1:Complete 1Yr 1.5T dataset were segmented by both methods.

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

457 **Evaluation method** As in previous experiments, the Winterburn hippocampal segmentation protocol  
458 differs in the delineated neuroanatomical features (Winterburn et al. (2013), and Table 9, Discussion) and  
459 so we assess MAGeT-Brain by the degree of (Pearson) correlation of average hippocampal volume across  
460 subjects. We also computed the correlation in hippocampal volume between existing, established automated  
461 segmentation methods – FSL FIRST, FreeSurfer, and MAPER, and SNT semi-automated segmentations.  
462 Additionally, we evaluate the volume-related fixed and proportional biases in all segmentation methods using  
463 Bland-Altman plots (Bland and Altman, 1986).

---

<sup>5</sup><http://fsl.fmrib.ox.ac.uk/fsl/fslwiki/FIRST/UserGuide>

Table 4: **ADNI1 1.5T Complete 1Yr dataset demographics.** CN - Cognitively Normal. LMCI - Late-onset Mild Cognitive Impairment. AD - Alzheimer’s Disease. CDR-SB - Clinical Dementia Rating-Sum of Boxes. ADAS - Alzheimer’s Disease Assessment Scale. MMSE - Mini-Mental State Examination.

	N	CN N = 584	LMCI N = 931	AD N = 404	Combined N = 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.

Table 5: Number of segmented images and quality control failures of ADNI1:Complete 1Yr 1.5T dataset by method.label

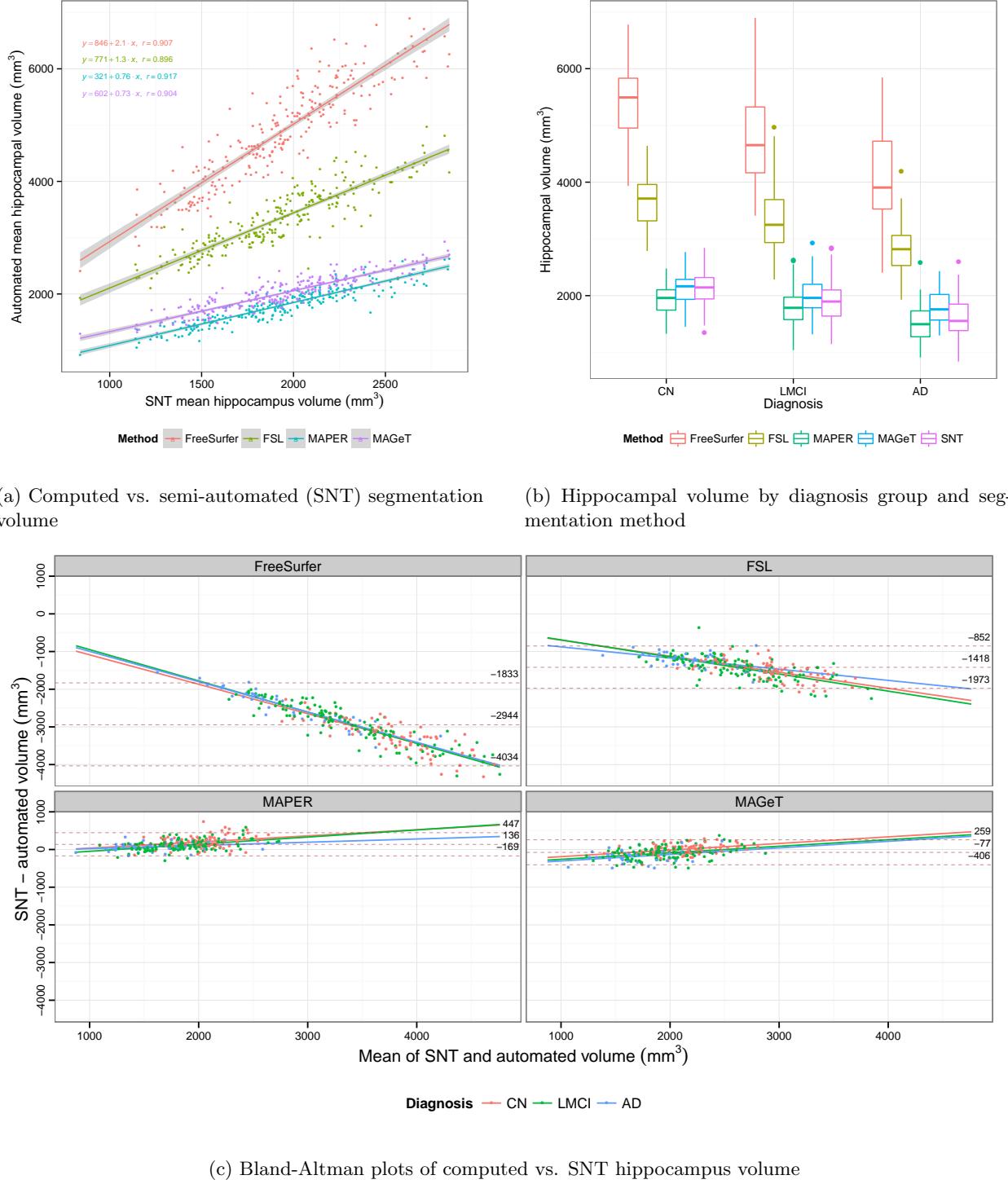
X	SNT	MAGeT	MAPER	FSL	FS
Images	368	368	368	368	368
Failures	n/a	30	n/a	20	88

### 464 3.3.2 Experiment 3: Results

465 We found a close relationship in total bilateral hippocampal volume between all methods and the SNT semi-  
 466 automated label volumes (Figure 4a). Volumes are well correlated ( $r > 0.78$ )c for all methods, and across  
 467 disease categories. Within disease categories (Figure 4b), MAGeT-Brain is consistently well correlated to  
 468 SNT volumes ( $r > 0.85$ ), but appears to slightly over-estimate the volume of the AD hippocampus compared  
 469 to the SNT segmentations.

470 Bland-Altman plots illustrate the level of agreement of each method with SNT segmentation hippocampal  
 471 volumes (Figure 4c). All methods show an obvious proportional bias: FreeSurfer and FSL FIRST markedly  
 472 under-estimate smaller hippocampi and over-estimate large hippocampi, whereas MAPER and MAGeT-  
 473 Brain show a reverse, conservative bias (Figure 4c). Additionally, all methods show a fixed volume bias,  
 474 with FreeSurfer and FSL FIRST most dramatically over-estimating hippocampal volume by  $2600mm^3$  and  
 475  $2800mm^3$  on average, respectively, and MAPER and MAGeT-Brain within  $250mm^3$  on average.

476 Figure 5 shows a qualitative comparison of MAGeT-Brain and SNT hippocampal segmentations for 10  
 477 randomly selected subjects in each disease category, and illustrates some of the common errors found during  
 478 visual inspection. Mostly frequently, we found MAGeT-Brain improperly includes the vestigial hippocampal  
 479 sulcus and, although not anatomically incorrect, MAGeT-Brain under-estimates the hippocampal body in  
 480 comparison to the SNT segmentation.



**Figure 4: ADNI1:Complete 1Yr 1.5T dataset segmentation.** (4a) Subject mean hippocampal volume as measured by each of the four automated methods (FreeSurfer (FS), FSL FIRST, MAPER, MAGeT-Brain) versus the semi-automated SNT segmentation volumes. Linear fit lines and Pearson correlations with SNT labels are shown for each method. (4b) Mean hippocampal volume by method and disease category. AD = Alzheimer’s disease, LMCI = late-onset mild cognitive impairment, and CN = cognitively normal. (4c) Bland-Altman plots shows the agreement between computed and SNT hippocampus volume. 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.

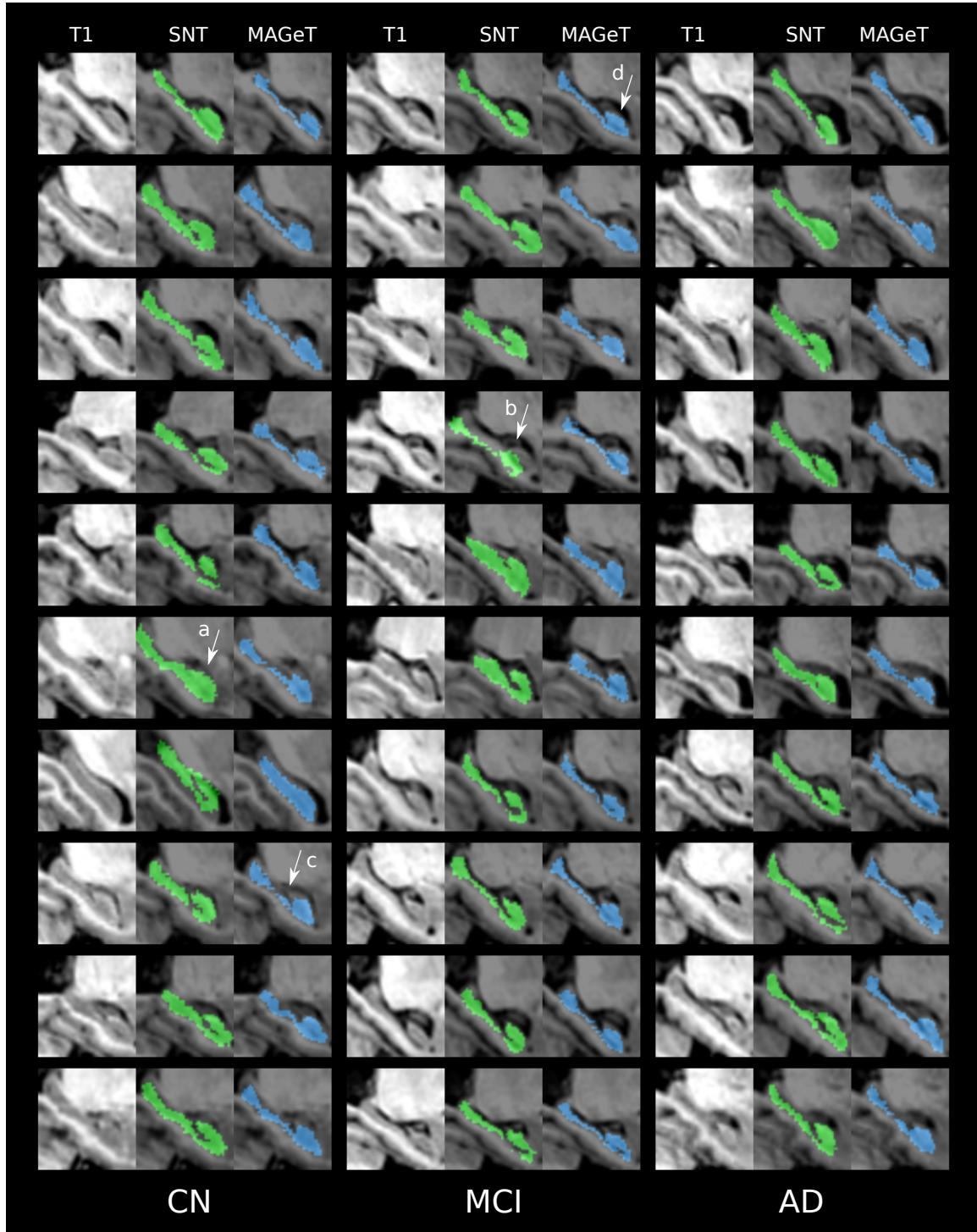


Figure 5: SNT and MAGeT-Brain segmentations for 30 ADNI subjects — 10 subjects randomly selected from each disease category in the subject pool used in Experiment 1 (Section 3.1). Sagittal slices are shown for each unlabelled T1-weighted anatomical image. SNT labels appear in green, and MAGeT-Brain labels 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.

**481 3.4 Experiment 4: Hippocampal Subfield Segmentation Cross-Validation**

482 The previous experiment assesses MAGeT-Brain performance on whole hippocampus segmentation. In this  
483 experiment, we ~~evaluate~~conduct a proof-of-concept evaluation of MAGeT-Brain hippocampal subfield  
484 segmentation of standard 3T T1-weighted images at 0.9mm-isotropic voxels. We use a modified leave-one-  
485 out cross-validation (LOOCV) design.

**486 3.4.1 Experiment 4: Materials and Methods**

487 **Healthy Control Dataset** T1 MR images of 14 subjects were acquired as a part of an ongoing study at  
488 the Centre for Addiction and Mental Health (Table 6). Subjects were known to be free of neuropsychiatric  
489 disorders and gave informed consent. These images were acquired on a 3T GE Discovery MR 750 system  
490 (General Electric, Milwaukee, WI) using an 8-channel head coil with the enhanced fast gradient recalled  
491 echo 3-dimensional acquisition protocol, FGRE-BRAVO, with the following parameters:  $TE/TR/TI =$   
492  $3.0ms/6.7ms/650ms$ , flip angle= $8^\circ$ ,  $FOV = 15.3cm$ , slice thickness=  $0.9mm$ , 170 in-plane steps for an  
493 approximate 0.9mm-isotropic voxel resolution.

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

498 In this experiment, the Winterburn atlases (Experiment 2, section 3.2) are resampled to 0.9mm-isotropic  
499 voxel resolution to simulate standard 3T T1-weighted resolution images. Image subsampling is performed  
500 using trilinear subsampling techniques. In each round of LOOCV, a single atlas image is selected and treated  
501 as a target image to be segmented by MAGeT-Brain. So as to have an odd-sized atlas library, atlas image  
502 is segmented once using each possible triple of atlas images, and corresponding manual segmentations, from  
503 the remaining four unselected atlases. Thus, for each of the five atlases, a total of  $\binom{3}{4} = 4$  segmentations are  
504 evaluated, resulting in a combined total of  $5 \times 4 = 20$  segmentations evaluated overall. We chose an atlas  
505 library with an odd number of images so as to ensure unbiased label fusion when using majority voting (see  
506 Discussion).

507 The template library used has a total of 19 images composed of all five resampled atlas images plus the  
508 additional 14 images from the healthy control dataset. The ANTS registration algorithm was used for image  
509 registration, and majority voting was used for label fusion, as these methods proved most favourable in the  
510 previous whole hippocampal validation experiments.

511 **Evaluation method** Evaluating the agreement of automated hippocampal subfield segmentations with  
512 manual segmentations for T1 images at 0.9mm-isotropic voxels is inherently ill-defined since there are no  
513 manual protocols for segmentation at this resolution. Instead, we must evaluate ~~the reliability, or precision,~~  
514 ~~with which how well the lower-resolution MAGeT-Brain produces hippocampal subfields segmentations at~~  
515 ~~this resolution that hippocampal subfield segmentations~~ correspond in form to the segmentation protocol  
516 used ~~by the given in the~~ high-resolution ~~atlas library images.~~

517 ~~images.~~ By directly resampling the Winterburn atlas segmentations to  $0.9mm^3$  voxels (using standard  
518 nearest-neighbour image resampling techniques) we obtain a subsampled version of the labels which preserve  
519 the original segmentation protocol within the limits of error from rounding and interpolation. Therefore,

Table 6: **Demographics for the hippocampal subfield cross-validation healthy control subject sample used in the template library (excluding the Winterburn atlas subjects).** Education is shown in years.

	N	Control N = 14
Age	14	34.5 53.0 62.0
Sex : Male	14	43% (6)
Education : 12	13	15% (2)
13		8% (1)
14		23% (3)
16		15% (2)
18		38% (5)
Handedness : R	14	93% (13)

$a$   $b$   $c$  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.

520 using the resampled Winterburn segmentations as definitive for the  $0.9mm^3$  resolution we evaluate **reliability**  
 521 **agreement** of MAGeT-Brain segmentations using DSC overlap scores and evaluate consistency across the  
 522 range of hippocampal sizes using Bland-Altman plots of subfield volumes.

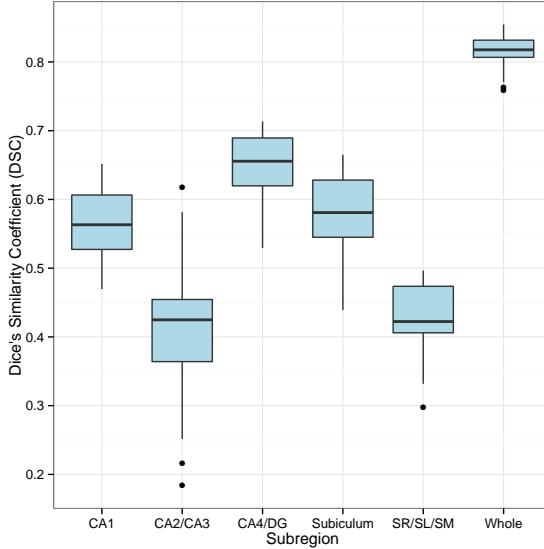
523 Additionally, by shifting the original manual  $0.3mm$ -isotropic voxel segmentations by one voxel in the x, y,  
 524 and z direction and then resampling it to  $0.9mm$ -isotropic voxels we obtain a simulated manual segmentation  
 525 having a small amount of error. We can compare the DSC overlap score of the shifted labels (relative to the  
 526 directly resampled labels) with the DSC score of the MAGeT-Brain generated labels in order to evaluate  
 527 their relevance.

### 528 3.4.2 Experiment 4: Results

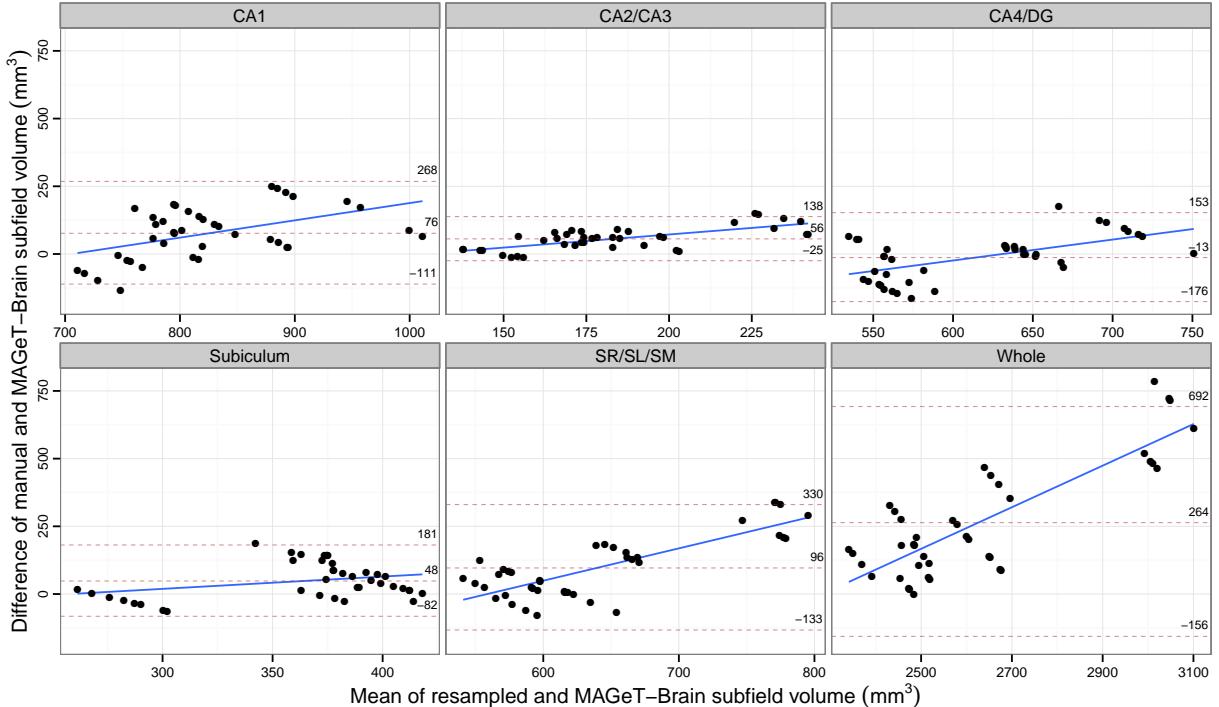
529 Figure 6a shows the overlap similarity scores between the MAGeT-Brain segmentations and the resampled  
 530 Winterburn atlases for each hippocampal subfield across all subjects and folds of the validation. Mean  
 531 and standard deviation DSC scores of the subfields are shown in Table 7, along with DSC scores for the  
 532 resampled atlas segmentations when perturbed slightly and compared to the originals. We find that the  
 533 CA4/DG subfield shows the highest mean DSC score of  $0.647 \pm 0.051$ , followed by the Subiculum and CA1  
 534 subfields having scores of  $0.563 \pm 0.046$  and  $0.58 \pm 0.057$ , respectively. Both the CA4/DG and molecular  
 535 regions score below 0.5. These scores may seem low but not when taken in context and compared to existing  
 536 (semi-)automated methods (see Discussion). The whole hippocampus is segmented with a mean DSC score  
 537 of  $0.816 \pm 0.023$ .

538 Figure 6b contains Bland-Altman plots comparing MAGeT-Brain volumes with manual volumes across  
 539 all validation folds. MAGeT-Brain displays a conservative proportional bias — small hippocampi are over-  
 540 estimated in volume, and larger hippocampi are underestimated (a mean maximum difference of approx-  
 541 imately  $200mm^3$  across all subfields). MAGeT-Brain display a slight conservative fixed bias, tending to  
 542 underestimate all subfields except CA4/DG (mean underestimation:  $CA1 = 76mm^3$ ,  $CA2/3 = 56mm^3$ ,  
 543  $CA4/DG = -16mm^3$ ,  $Subiculum = 48mm^3$ ,  $SR/SL/SM = 96mm^3$ ).

544 Figure 7 shows slices subfield segmentations for a single subject for qualitative inspection.

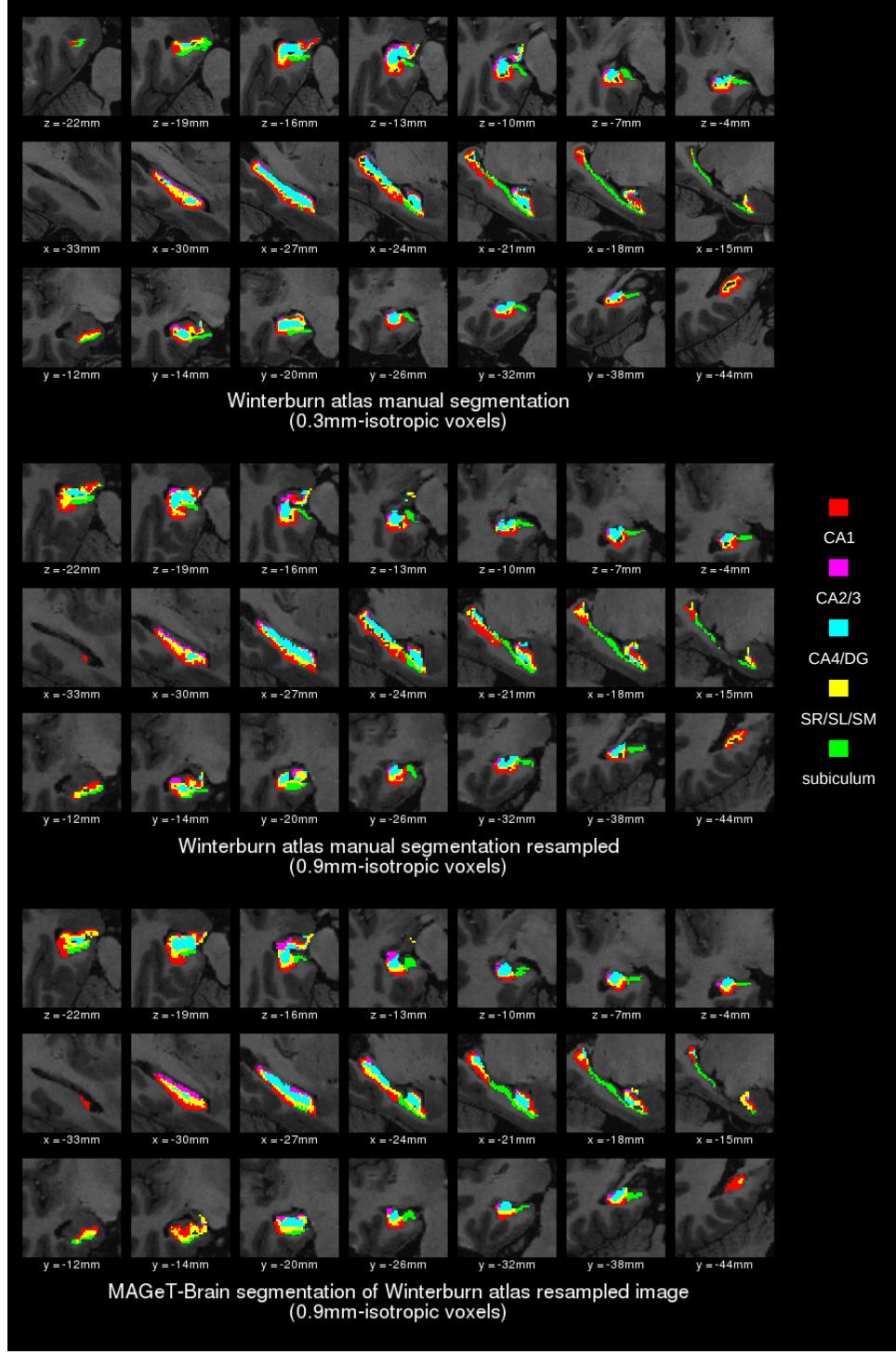


(a) DSC score by subfield



(b) Bland-Altman plots of computed vs. manual subfield volumes

**Figure 6: Hippocampal subfield cross-validation.** (6a) Similarity of MAGEt-Brain segmentation of subfields and the resampled Winterburn atlas segmentations at  $0.9\text{mm}^3$  voxel resolution, over all validation folds. Overlap score for each hemisphere is measured separately. (6b) shows the agreement, by subfield, of computed and manual volumes across all validation folds. The overall mean difference in volume, and limits of agreement ( $\pm 1.96SD$ ) are shown by dashed horizontal lines. Linear fit lines are shown. Note, points below the mean difference indicate overestimation of the volume with respect to the resampled volume, and vice versa.



**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 subsampled from 0.3mm- to 0.9mm-isotropic voxels; in the lower left corner is the MAGeT-Brain segmentation of the resampled Winterburn atlas image from a single fold of the cross-validation. 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.

Table 7: Overlap similarity results for the each of the subfields of the hippocampus. Simulated overlap similarity results are also given for manual labels that were translated by one voxel (i.e.:  $0.3\text{mm}$ ) in all directions and then resampled. Values are given as mean Dice’s Similarity Coefficient (DSC)  $\pm$  standard deviation.

Subfield	MAGeT	$0.9\text{mm}$ translation
CA1	$0.56 \pm 0.05$	$0.27 \pm 0.03$
CA2/CA3	$0.41 \pm 0.10$	$0.12 \pm 0.05$
CA4/DG	$0.65 \pm 0.05$	$0.42 \pm 0.05$
SR/SL/SM	$0.43 \pm 0.05$	$0.19 \pm 0.04$
Subiculum	$0.58 \pm 0.06$	$0.14 \pm 0.04$

## 545 4 Discussion

546 In this manuscript we have presented the implementation and validation of the MAGeT-Brain framework –  
 547 a methodology that requires very few input atlases in order to provide accurate and reliable segmentations  
 548 with respect to manual segmentations. Both Experiment 1 (Section 3.1) and Experiment 2 (Section 3.2)  
 549 compare MAGeT-Brain to basic-multi-atlas segmentation by characterising the change in segmentation qual-  
 550 ity with varying parameter settings (atlas and template library sizes, registration method, and label fusion  
 551 method) and differing age and neuropsychiatric populations. Together, these experiments allow us to choose  
 552 optimal MAGeT-Brain parameter settings for use in subsequent experiments. Experiment 3 (Section 3.3)  
 553 demonstrates that across 246 images from the ADNI1:Complete 1Yr 1.5T dataset, MAGeT-Brain performs  
 554 as well as, or better, than other established and popular methods, and has a much more conservative propor-  
 555 tional bias in segmentation volume. Finally, Experiment 4 (Section 3.4) demonstrates is a proof-of-concept  
 556 validation demonstrating the reliability of MAGeT-Brain in producing subfield segmentations which match  
 557 the segmentation protocol of the input atlases despite contrast and resolution limitations in standard T1-  
 558 weighted image volumes. All of these experiments together demonstrate that MAGeT-Brain’s algorithmic  
 559 performance is not dependent on a single definition of the hippocampus but is effective with differing hip-  
 560 pocampal definitions (Winterburn et al., 2013; Pruessner et al., 2000; Hsu et al., 2002), across image types,  
 561 and subject populations.

562 The core claim the MAGeT-Brain method is based on – that we can meaningfully bootstrap a template  
 563 library-a useful template library can be generated from a small set of labelled atlas images – is validated in  
 564 the cross-validation-cross-validation conducted in Experiment 1 (and the replication in Experiment 2 and  
 565 Experiment 5, Supplementary Materials). We find that both increasing the number of atlases and the number  
 566 of templates used improves MAGeT-Brain segmentation over and above basic-multi-atlas segmentations  
 567 using the same number of atlas images. That is, by taking the extra step of generating a template library  
 568 using target images, MAGeT-Brain is able to improve the overlap between the automatically generated  
 569 segmentations and manually generated “gold standard” segmentations. The magnitude of this improvement  
 570 is greatest with a small number of atlases, but even with larger atlas libraries we have found that generating  
 571 a template library reduces the variability in segmentation precision agreement (i.e. MAGeT-Brain more  
 572 consistently produces high-quality segmentations in greater agreement with manual segmentations than  
 573 does basic-multi-atlas segmentation method, over repeated randomized trials). These effects do not appear  
 574 dependant on the hippocampal segmentation protocol used.

575 Interestingly, previous work on multi-atlas segmentation methods (Aljabar et al., 2009; Collins and  
 576 Pruessner, 2010) has found that cross-correlation and normalized mutual information-based weighted la-

**Table 8: Automated segmentation accuracy of the Hippocampus Summary of automated segmentation methods of the Hippocampus.**  
This table shows the best summarizes published Dice's overlap measure between automated and ground-truth-manual segmentations of the hippocampus. Unless otherwise specified, validation datasets are composed equally of cases and control subjects, and use manual segmentation labels as ground truth in computing DSC scores. AD = Alzheimer's Disease; MCI = Mild Cognitive Impairment; CN = Cognitively Normal (CN); FEP = First Episode of Psychosis; LOOCV = Leave-one-out cross-validation; MCCC = Monte Carlo cross-validation; SNT = Surgical Medtronic Navigation Technologies semi-automated labels. 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 mean (AD; MCI; CN)	Reference	Validation	Dataset (Truth)
MAGeT-Brain	9	0.841		10-fold MCCC on 69 subjects	ADNI (SNT)
Patch-based label fusion	16	0.861 (0.838; —; 0.883)	Coupe et al. (2011)	LOOCV	ADNI (SNT)
Multi-atlas	20	0.848 (—; 0.798; 0.898)	Wang et al. (2011)	10-fold MCCC on 20 of 139 subjects	ADNI (SNT)
ACM (AdaBoost-based)	21	0.862	Morra et al. (2008)	LOOCV	ADNI (SNT)
LEAP	30	0.848	Wolz et al. (2010)	Segmentation of 182 subjects	ADNI (SNT)
Multi-atlas	30	0.885	Lötjönen et al. (2010)	Segmentation of 60 subjects	ADNI (SNT)
Multi-atlas (MAPS)	55	0.890	Leung et al. (2010)	Segmentation of 30 subjects (10 AD, MCI, and CN)	ADNI (SNT)
MAGeT-Brain	9	0.869		10-fold MCCC on 60 subjects	ADNI (Pruessner)
MAGeT-Brain	9	0.892		5-fold MCCC on 81 subjects	FEP subjects
Neural nets	10	0.740		Segmentation of 5 subjects	controls
Probabilistic atlas	11	0.852		11 atlases used in 100 rounds of LOOCV on 20 elderly subjects	elderly controls
Probabilistic Atlas	16	0.860	Powell et al. (2008)	LOOCV	AD subjects
Anatomically-guided EM	17	0.812	van der Lijn et al. (2008)	LOOCV	mixed diagnosis
Multi-atlas	30	0.820	Chupin et al. (2009)	LOOCV	controls
Multi-atlas	30	0.880	Pohl et al. (2007)	LOOCV	2yr old controls
Multi-atlas	80	0.890	Heckemann et al. (2006a)	30 adult atlas used, segmentation of 33 2yr old subjects	controls
Multi-atlas	55	0.860	Gousias et al. (2008)	LOOCV	Collins and Pruessner (2010)
Multi-atlas	275	0.835	Barnes et al. (2008)	LOOCV	Aljabar et al. (2009)
Multi-atlas				LOOCV	controls and AD controls

577 label fusion improves segmentation ~~accuracy-reliability~~ over simple majority vote label fusion, and yet we did  
578 not see a significant indication of this effect in the MAGeT-Brain segmentations. Selectively filtering out  
579 atlases with lower image similarity is believed to reduce sources of error from estimating deformations via  
580 nonlinear registration, partial volume effects from nearest neighbour image resampling, and neuroanatomical  
581 mismatch between atlases and subjects. That MAGeT-Brain does not see the same boost in performance  
582 from weighted voting may suggest that the neuroanatomical variability of a template library constructed  
583 from study subjects more closely matches any particular subject and thereby leaving less error to filter. From  
584 our previous work on the MAGeT-Brain algorithm we have shown that the reduction in error is not simply  
585 a smoothing or averaging effect (Chakravarty et al., 2013).

586 Although, the goal of this manuscript was not to exhaustively test or validate multiple different voting  
587 strategies in the context of our segmentation algorithm, it is important to note that other strategies for  
588 voting are available. For example, other groups have used the STAPLE algorithm (Warfield et al., 2004)  
589 (or variants of the STAPLE algorithm (Robitaille and Duchesne, 2012)) which weights each segmentation  
590 based upon its estimated performance level with respect to the other available candidate segmentations.  
591 Further, the sensitivity and specificity parameters can also be tuned to potentially improve segmentation  
592 ~~accuracy-and-reliability~~. It is likely that using more sophisticated voting methods would have a positive  
593 effect on the overall segmentation performance, as demonstrated by the STAPLE algorithm. However, it  
594 is also important to note that even in the absence of a more sophisticated label fusion algorithm, MAGeT  
595 Brain performs reasonably well in comparison to other groups that have tested new segmentation algorithm  
596 with Alzheimer disease, mild cognitive impairment, and cognitively normal data from the ADNI database  
597 (Table 8. In addition, our validation in Experiment 2 (with the first episode psychosis subjects) yields DSC's  
598 that are amongst the highest reported. Thus, more work is required to determine the extent to which label  
599 fusion will improve the ~~accuracy-reliability~~ of our algorithm.

600 More work is required to determine the source of the slight decrease in segmentation performance when  
601 the number of templates are set to an even number. Our initial concern was that this dip in performance  
602 was a by-product of the MAGeT-Brain algorithm itself. However, this pattern is also found in the results of  
603 the multi-atlas segmentations we used in our experiments. We believe that our majority voting methodology  
604 is biased towards labels with the lowest numeric values when breaking ties (by way of the implementation  
605 of the `mode` function used to determine majority), thus causing the slight bias observed when using an even  
606 number of templates. This is another area where the voting scheme could be used to improve performance.  
607 However, it is worth noting that this limitation was previously identified by Heckemann et al. (2006b) and,  
608 subsequently, other groups have not even considered the potential pitfalls of an even number of candidate  
609 labels (e.g. Leung et al. (2010)).

610 Despite MAGeT-Brain achieving segmentation results which are competitive with the rest of the field  
611 (Table 8), a concern may be raised over the modest improvement in segmentation agreement observed using  
612 MAGeT-Brain over multi-atlas, with the same number of atlases (Experiment 1). As we have shown in that  
613 same experiment, the benefit in using MAGeT-Brain is both an increase in the overlap agreement and also  
614 in the improved consistency of the labelling regardless of atlas or template choice. Reducing the variability  
615 in segmentation agreement is an important consideration that few have touched on previously. In addition,  
616 the Monte Carlo cross-validations that we present in Experiment 1 and Experiment 2 are amongst the most  
617 stringent performed in the multi-atlas segmentation literature. To the best of our knowledge, with the  
618 exception of (Wang et al., 2011), other groups do at most a single round of leave-one-out-validation (Table  
619 8). Thus, the thoroughness of our validation suggests that our results are reflective of a true average over

620 the choice of parameter settings and are independent of atlas or template choice.

621 On that note, one author (JW), an expert manual rater (Winterburn et al., 2013), identified regular in-  
622 consistencies in the SNT segmentations: occurrences of over- and under-estimation, as well as misalignments  
623 of the entire segmentation volume (Figure 5). Although the SNT segmentations are used as benchmarks  
624 for validation in many other studies (Table 8), these segmentation inconsistencies present the possibility  
625 that a more accurate and consistent benchmark segmentation protocol ought to be used in order to truly  
626 understand the results of such validations. Indeed, our replication of the 10-fold cross-validation using SNT  
627 segmentations (Experiment 5, Supplementary Materials) shows noticeably poorer mean similarity scores for  
628 both MAGeT-Brain and multi-atlas.

629 Thus, in comparison to other methodologies in the field MAGeT-Brain performs favourably. Table 8 sur-  
630 veys some of the most recent reported DSC values reported on ADNI dataset, using SNT segmentations for  
631 the atlas library and as gold standards for evaluation. While it is difficult to compare segmentation results  
632 across studies, gold standards, evaluation metrics, and algorithms it is worth noting that the methods sum-  
633 marized require more atlases (between 16-55) than our MAGeT-Brain implementation with the Winterburn  
634 atlases (Winterburn et al., 2013).

635 There are some important differences between our method and these specific methods. Others have  
636 reported the difficulty with mis-registrations in candidate segmentation (i.e. segmentations generated that  
637 are then input in the voxel-voting procedure (Collins and Pruessner, 2010)). The work of Leung et al.  
638 (2010) tackles this problem by using an intensity threshold that is estimated heuristically at the time of  
639 segmentation (this work also reports some of the highest DSC scores for the segmentation of ADNI data).  
640 While this method is effective for the ADNI dataset (which is partially homogenized with respect to image  
641 acquisition and pre-processing), it is unclear if this type of heuristic is applicable to other datasets. In all  
642 cases, these methods require more atlases than our implementation with the Winterburn atlases. Lötjönen  
643 et al. (2010) achieve highly accurate segmentation but correct their segmentations produced segmentations  
644 which strongly agree with manual segmentations by way of post-processing corrections using classifications  
645 derived using an expectation maximization framework. In their initial work, Chupin et al. (2009) develop  
646 their probabilistic methodology using a cohort of 8 healthy controls and 15 epilepsy patients, and then use  
647 this method to segment an ADNI sample, with a hierarchical experimentation protocol. These methods  
648 suggest that some post-processing of the final segmentations would improve agreement of the segmentation.  
649 While that may be true, there is little consensus regarding how to achieve this.

650 To the best of our knowledge, no other groups have validated their work using multiple atlas segmentation  
651 protocols, different acquisitions, and disease populations in order to demonstrate the robustness of their  
652 technique. This is one of the clear strengths of this work. Furthermore, unlike some of the algorithms  
653 mentioned, our implementation does not require retuning for new populations or datasets as it inherently  
654 models the variability of the dataset through the template library. However it should be noted that the  
655 increased agreement that follows increasing the number of atlases and templates comes at an increased  
656 computational cost ( $O(\log(n))$ ), as previously mentioned in other work (Heckemann et al., 2006a).

657 Among the automated segmentation methods we compared in this paper (FreeSurfer, MAPER, FSL  
658 FIRST), we find extremely variable performance of all methods. With the exception of FSL FIRST all  
659 methods correlate well with the semi-automated SNT volumes provided in the ADNI database. However,  
660 the FreeSurfer and FSL FIRST provide radically different definitions of the size of the hippocampus in  
661 comparison to the hippocampal segmentations are on average about twice the volume of those from all other  
662 methods. Furthermore, when estimating bias of these methods relative to the bias of FreeSurfer and

---

Table 9: **Summary of labelled subfields of the Hippocampus from recent MRI segmentation protocols.**

Protocol	Labelled Subfields
Winterburn et al. (2013)	CA1, CA2/CA3, CA4/dentate gyrus, strata radiatum/lacunosum/moleculare, subiculum
Wisse et al. (2012)	CA1, CA2, CA3, CA4/dentate gyrus, subiculum, entorhinal cortex
Van Leemput et al. (2009)	CA1, CA2/CA3, CA4/dentate gyrus, presubiculum, subiculum, hippocampal fissure, fimbria, hippocampal tail, inferior lateral ventricle, choroid plexus
Yushkevich et al. (2009)	CA1, CA2/CA3, dentate gyrus (hilus), dentate gyrus (stratum moleculare), strata radiatum/lacunosom/moleculare/vestigial hippocampal sulcus
Mueller et al. (2007)	CA1, CA2, CA3/CA4 & dentate gyrus, Sibicum, entorhinal cortex

663 FSL FIRST [relative to the](#) SNT hippocampal volumes we see that large hippocampi are over estimated  
 664 while small hippocampi are under estimated. By comparison, MAGE-T-Brain and MAPER are far more  
 665 conservative in volume estimation, suggesting these methods may be better suited for estimating true-  
 666 positives, especially in neurodegenerative disease subjects featuring smaller overall hippocampi. However, in  
 667 this analysis we have only compared methods by total hippocampal volume, and so more work is needed to  
 668 understand the full extent to which these methods differ.

669 Finally, we have provided evidence that using the Winterburn high-resolution hippocampal subfield at-  
 670 lases (Winterburn et al., 2013) our algorithmic framework is appropriate for the segmentation of hippocampal  
 671 subfields in standard T1-weighted data. Subfield segmentation is a burgeoning topic in the literature although  
 672 very few automated methods are available for the segmentation of 3T data (Yushkevich et al., 2009, 2010;  
 673 Van Leemput et al., 2009). Table 10 compares segmentation agreement from some of these methods and  
 674 MAGE-T-Brain. The overlap DSC scores for MAGE-T-Brain subfields are notably lower but a direct compar-  
 675 ison of overlap values must be done cautiously. In the present work, our overlap scores are computed on  
 676 0.9mm-isotropic voxel resolution images, whereas Yushkevich et al. (2010) uses focal  $0.4 \times 0.5 \times 2.0\text{mm}$  voxel  
 677 resolution images, and Van Leemput et al. (2009) use supersampled 0.380mm-isotropic voxel resolution im-  
 678 ages. The larger voxel images we use necessarily entail a greater change in DSC for each incorrectly labelled  
 679 voxel. In addition, our automated segmentations are compared to manual segmentations resampled from  
 680 0.3mm-isotropic voxel labels; the resampling process inevitably introduces noise which may lower overlap  
 681 scores. Lastly, as our method is aimed specifically at situations when manually produced atlases are scarce,  
 682 in our cross validation we are forced to use three rather than all five of the Winterburn atlases (which, based  
 683 on our findings with whole hippocampal segmentation, would have resulted in improved overlap similarity).  
 684 Although having more atlases would be ideal in this context, these atlases are very time consuming to gen-  
 685 erate (Winterburn et al., 2013). Nevertheless, the advantage of evaluating MAGE-T-Brain on standard 3T  
 686 T1-weighted resolution MR images with a publically available atlas library is that our results reflect typical  
 687 usage scenarios of researchers and clinicians.

688 Experiments 1, 2, and 5 have demonstrated that our algorithm flexibly accommodates different whole  
 689 hippocampus manual segmentation methodologies. We have not explicitly evaluated a subfield definition  
 690 other than the Winterburn protocol, and therefore it is possible that using an alternate subfield definition  
 691 could improve the reliability of our automated subfield definitions. For example, established definitions such  
 692 as those from Mueller et al. (2007) could be a prime candidate for further exploration. In addition, the  
 693 conservative nature of the Mueller definition (labelling of the 5 slices in the hippocampus body only) would  
 694 likely further aid in reliability measurement. However, there are two main logistical problems that we would  
 695 have to overcome prior to implementation. The first is that these definitions were developed for data that

---

Table 10: A comparison of subfield segmentation overlap similarity with manual raters.

Subfield	MAGeT-Brain	Van Leemput et al. (2009)	Yushkevich et al. (2010)
CA1	0.563	0.62	0.875
CA2/3	0.412	0.74	$CA2 = 0.538, CA3 = 0.618$
CA4/DG	0.647	0.68	$DG = 0.873$
presubiculum	—	0.68	—
subiculum	0.58	0.74	0.770
hippocampal fissure	—	0.53	—
SR/SL/SM	0.428	—	—
fimbria	—	0.51	—
head	—	—	0.902
tail	—	—	0.863

is highly anisotropic ( $0.4mm \times 0.5mm \times 2mm$ ), and it is unclear how our algorithms would deal with such atlases used as input. The second is that, since these atlases are not publicly available, we would have to re-implement the protocol using our atlases. At the present time it is unclear how we would adapt these protocol to data that we used, where subfield segmentations are defined on  $0.3mm^3$  voxels. However, the impact of subfield definitions in the context of our work is an important one and should be considered in subsequent studies.

One further complication common to all subfield segmentation evaluation is that, by its nature, the Dice's Similarity Coefficient score penalizes structures with high surface area-to-volume ratios. Therefore subfield DSC scores will generally be lower than whole hippocampal segmentations. We attempted to put this effect into perspective by comparing MAGeT-Brain subfield segmentation agreement with the agreement of voxel-shifted manual segmentations (Table 7). The results of this exercise show conclusively, despite the very limited number of atlases we had to work with, that MAGeT-Brain subfield segmentations are well within the bounds of error of a  $0.3mm^3$  voxel shift.

Our overlap DSC values demonstrates that we can reliably reproduce segmentations for the CA1, subiculum, and CA4/dentate subfields ( $DSC > 0.5$ ). That the CA2/CA3 and molecular layers are less well reproduced ( $DSC < 0.5$ ) should not be surprising as these are extremely thin and spatially convoluted regions that originally required high-resolution MRI for identification and so it is likely that the extents of these regions are well below the resolution and contrast offered by standard T1-weighted images.

This points to a larger issue of how to truly validate subfield segmentations, both in high resolution images and in standard T1-weighted images. There are several manual subfield segmentation methodologies, and they do not agree on which regions can be differentiated, even on high-resolution scans. See Table 9 for a comparison of MRI-based manual subfield segmentation methodologies. A further complication is that different researchers have differing operational definitions for the subfields and how they ought to be parcellated. The disagreement in the community has led to an international working group devoted to normalizing the ontology and segmentation rules for the hippocampal subfields (<http://www.hippocampalsubfields.com/>). In addition, there have been recent advances from the Yushkevich group to revise their MRI subfield segmentation protocol based on anatomy discerned from serial histological acquisitions (Adler et al., 2014). The definitional and operational disagreements suggest that direct comparison across automated methods using “ground truth”-based overlap similarity metrics, such as Dice's Similarity Coefficient, are not possible without carefully taking into account the differences in underlying segmentation protocols and image characteristics.

In conclusion, we have demonstrated the viability of leveraging a small number of input atlases to **bootstrap**-**generate** a large template library and thereby improve segmentation reliability when using multi-

---

729 atlas methods. We demonstrated that this method works robustly over hippocampal definitions, different  
730 disease populations, and different acquisition types. Finally, we also demonstrate that reliable reproduction  
731 of hippocampal subfield segmentations in standard 3T T1-weighted images is possible.

## 732 5 Acknowledgements

733 We wish acknowledge support from the CAMH Foundation, thanks to Michael and Sonja Koerner, the Kimel  
734 Family, and the Paul E. Garfinkel New Investigator Catalyst Award. MMC is funded by the W. Garfield  
735 Weston Foundation and ANV is funded by the Canadian Institutes of Health Research, Ontario Mental  
736 Health Foundation, NARSAD, and the National Institute of Mental Health (R01MH099167).

737 Computations were performed on the gpc supercomputer at the SciNet HPC Consortium (Loken et al.,  
738 2010). SciNet is funded by: the Canada Foundation for Innovation under the auspices of Compute Canada;  
739 the Government of Ontario; Ontario Research Fund - Research Excellence; and the University of Toronto.

740 In addition, computations were performed on the CAMH Specialized Computing Cluster. The SCC is  
741 funded by: The Canada Foundation for Innovation, Research Hospital Fund.

742 ADNI Acknowledgements: Data collection and sharing for this project was funded by the Alzheimer's  
743 Disease Neuroimaging Initiative (ADNI) (National Institutes of Health Grant U01 AG024904). ADNI is  
744 funded by the National Institute on Aging, the National Institute of Biomedical Imaging and Bioengineering,  
745 and through generous contributions from the following: Abbott; Alzheimer's Association; Alzheimer's  
746 Drug Discovery Foundation; Amorfix Life Sciences Ltd.; AstraZeneca; Bayer HealthCare; BioClinica, Inc.;  
747 Biogen Idec Inc.; Bristol-Myers Squibb Company; Eisai Inc.; Elan Pharmaceuticals Inc.; Eli Lilly and  
748 Company; F. Hoffmann-La Roche Ltd and its affiliated company Genentech, Inc.; GE Healthcare; Innogenetics,  
749 N.V.; IXICO Ltd.; Janssen Alzheimer Immunotherapy Research Development, LLC.; Johnson & Johnson  
750 Pharmaceutical Research Development LLC.; Medpace, Inc.; Merck & Co., Inc.; Meso Scale Diagnostics,  
751 LLC.; Novartis Pharmaceuticals Corporation; Pfizer Inc.; Servier; Synarc Inc.; and Takeda Pharmaceutical  
752 Company. The Canadian Institutes of Health Research is providing funds to support ADNI clinical sites in  
753 Canada. Private sector contributions are facilitated by the Foundation for the National Institutes of Health  
754 ([www.fnih.org](http://www.fnih.org)). The grantee organization is the Northern California Institute for Research and Education,  
755 and the study is Rev March 26, 2012 coordinated by the Alzheimer's disease Cooperative Study at the  
756 University of California, San Diego. ADNI data are disseminated by the Laboratory for NeuroImaging at  
757 the University of California, Los Angeles. This research was also supported by NIH grants P30 AG010129  
758 and K01 AG030514.

759 We would also like to thank G. Clinton, E. Hazel, and B. Worrell for inspiring this work.

## 760 6 Supplementary Materials

### 761 6.1 SNT Hippocampal Labels

762 Semi-automated hippocampal volumetry was carried out using a commercially available high dimensional  
763 brain mapping tool (Medtronic Surgical Navigation Technologies, Louisville, CO), that has previously been  
764 validated and compared to manual tracing of the hippocampus (Hsu et al., 2002). Measurement of hippocam-  
765 pal volume is achieved first by placing manually 22 control points as local landmarks for the hippocampus on  
766 the individual brain MRI data: one landmark at the hippocampal head, one at the tail, and four per image

## 6.2 Experiment 5: Whole Hippocampus Segmentation Cross-Validation — Alzheimer’s Disease, SNT Segmentations

---

Table 11: **ADNI1:Complete 1Yr 1.5T SNT 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 N = 23	LMCI N = 23	AD N = 23	Combined N = 69
Age at baseline Years	72.2 75.5 78.5	71.0 77.1 81.4	71.7 77.8 81.8	71.5 76.6 81.3
Sex : Female	43% (10)	43% (10)	43% (10)	43% (30)
Education	16.0 16.0 18.0	15.0 16.0 18.0	12.0 16.0 16.5	14.0 16.0 18.0
CDR-SB	0.00 0.00 0.00	0.75 1.50 1.50	4.00 4.50 5.00	0.00 1.50 4.00
ADAS 13	4.67 5.67 12.34	14.34 16.00 20.50	23.83 29.00 31.66	10.00 16.00 25.33
MMSE	28.5 29.0 30.0	25.0 27.0 28.0	21.0 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.

767 (i.e., at the superior, inferior, medial and lateral boundaries) on five equally spaced images perpendicular  
 768 to the long axis of the hippocampus. Second, fluid image transformation is used to match the individual  
 769 brains to a template brain (Christensen et al., 1997). The pixels corresponding to the hippocampus are then  
 770 labeled and counted to obtain volumes. This method of hippocampal voluming has a documented reliability  
 771 of an intraclass coefficient better than .94 (Hsu et al., 2002).

## 772 6.2 Experiment 5: Whole Hippocampus Segmentation Cross-Validation — Alzheimer’s 773 Disease, SNT Segmentations

774 This experiment is a replication of Experiment 1 using a pool of 69 images and SNT semi-automated  
 775 segmentations from the ADNI dataset (Hsu et al., 2002). See Experiment 1 for full details on the ADNI  
 776 dataset, and validation process.

### 777 6.2.1 Experiment 5: Materials and Methods

778 **Dataset** 69 1.5T images were arbitrarily selected from the baseline scans in the *ADNI1:Complete 1Yr 1.5T*  
 779 standardized dataset. 23 subjects were chosen from each disease category: cognitively normal (CN), mild  
 780 cognitive impairment (MCI) and Alzheimer’s disease (AD). Demographics for this subset are shown in Table  
 781 1. Each image has a corresponding semi-automated segmentation of the left and right whole hippocampus  
 782 made available with ADNI images (SNT; see Supplementary Materials).

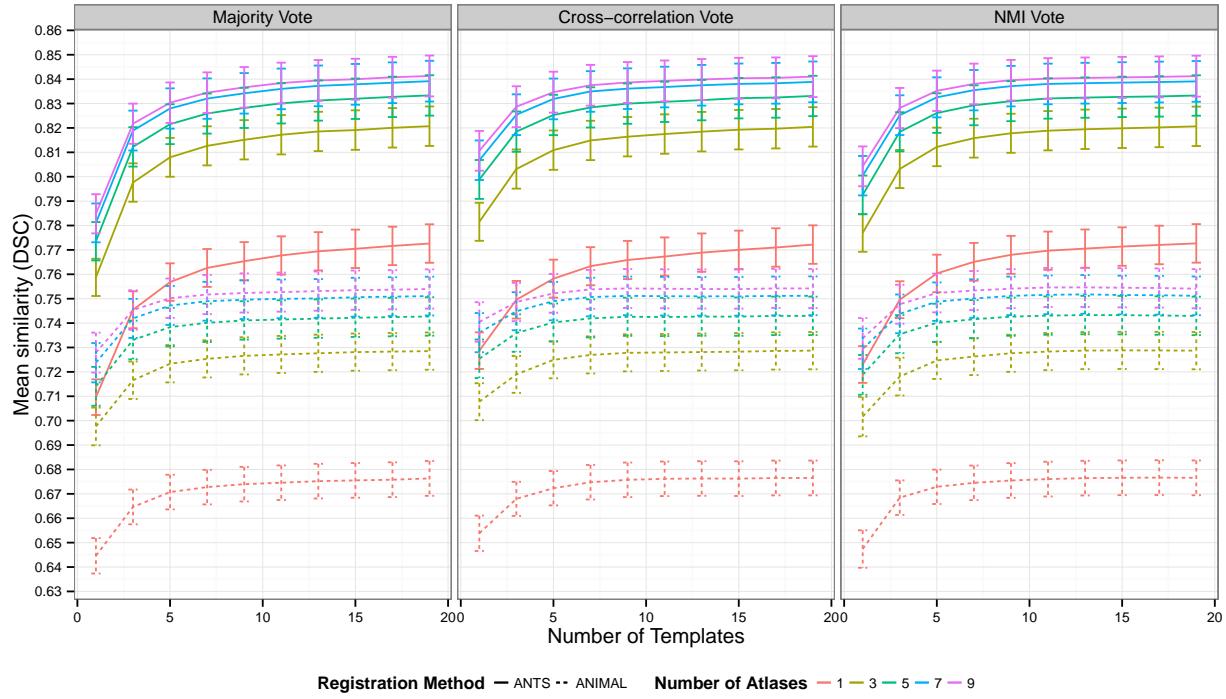
783 **Experiment details** A total of ten validation rounds were performed on each subject in the dataset, for  
 784 each combination of parameter settings: atlas library size (1-9), template library size (1-20), registration  
 785 method (ANTS or ANIMAL), and label fusion method (majority vote, cross-correlation weighted majority  
 786 vote, and normalized mutual information weighted majority vote). A total of  $10 \times 69 \times 9 \times 20 \times 2 \times 3 =$   
 787  $7.452 \times 10^5$  validation rounds are conducted. The computed segmentations for a subject are compared to  
 788 the SNT labels provided by ADNI using Dice’s Similarity Coefficient and the score is averaged over the  
 789 validation rounds.

### 790 6.2.2 Experiment 5: Results

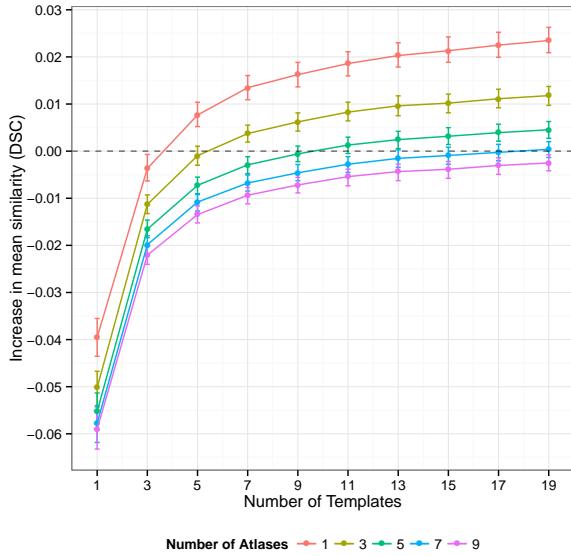
791 As when comparing against manual labels in Experiment 1, we find similar behaviour when comparing  
 792 MAGeT-Brain segmentations to SNT labels: similarity scores increase with increasing numbers of atlases

793 and templates, with diminishing increases in improvement trending towards a plateau (Figure 2a). As in  
794 Experiment 1, using ANTS registration leads to significantly increased similarity scores, and there is no  
795 significant difference in scores from any of the label fusion methods. Mean DSC score peaks at 0.841 when  
796 using 9 atlases, 19 templates, ANTS registration, and majority vote label fusion. Compared to multi-atlas  
797 segmentations, we find MAGeT-Brain segmentations show increasing improvement with larger atlas and  
798 template libraries when using more than 9 templates and 5 or fewer atlases (Figure 8b). Peak improvement  
799 (+0.023 DSC) is found with a single atlas and template library of 19 images. In addition to a mean increase in  
800 similarity score over multi-atlas-based segmentation, MAGeT-Brain also shows more consistency in similarity  
801 scores across all subjects and validation folds (Figure 8c) with a large enough template library.

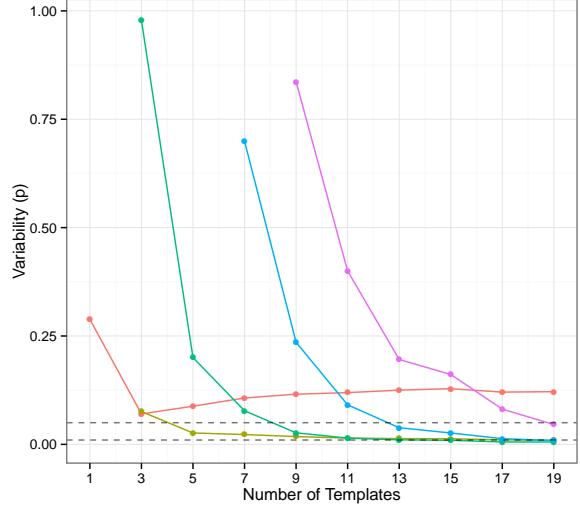
802 S



(a) DSC vs. atlas and template library size



(b) Increase in similarity score over multi-atlas



(c) Difference in variability with multi-atlas

Figure 8: **Whole hippocampus segmentation cross-validation on ADNI subjects with SNT segmentations.** (8a) Average DSC score of MAGeT-Brain with SNT segmentations for 69 ADNI subjects taken over 10 folds of cross-validation at each parameter setting. Error bars indicate standard error. (8b) Increase in DSC of MAGeT-Brain over multi-atlas segmentations. (8c) shows the significance of t-tests comparing the variability in DSC scores of MAGeT-Brain and multi-atlas across validation folds. Only points where MAGeT-Brain mean variability is lower than multi-atlas are shown. Dashed lines indicate p-values of 0.05 and 0.01.

**803 References**

- 804 D. H. Adler, J. Pluta, S. Kadivar, C. Craige, J. C. Gee, B. B. Avants, and P. a. Yushkevich. Histology-derived  
805 volumetric annotation of the human hippocampal subfields in postmortem MRI. *NeuroImage*, 84:505–23,  
806 Jan. 2014. ISSN 1095-9572. doi: 10.1016/j.neuroimage.2013.08.067.
- 807 P. Aljabar, R. a. Heckemann, a. Hammers, J. V. Hajnal, and D. Rueckert. Multi-atlas based segmentation  
808 of brain images: atlas selection and its effect on accuracy. *NeuroImage*, 46(3):726–38, July 2009. ISSN  
809 1095-9572. doi: 10.1016/j.neuroimage.2009.02.018.
- 810 B. B. Avants, C. L. Epstein, M. Grossman, and J. C. Gee. Symmetric diffeomorphic image registration with  
811 cross-correlation: evaluating automated labeling of elderly and neurodegenerative brain. *Medical image  
analysis*, 12(1):26–41, Feb. 2008. ISSN 1361-8423. doi: 10.1016/j.media.2007.06.004.
- 812 J. Barnes, J. Foster, R. G. Boyes, T. Pepple, E. K. Moore, J. M. Schott, C. Frost, R. I. Scahill, and N. C. Fox.  
813 A comparison of methods for the automated calculation of volumes and atrophy rates in the hippocampus.  
814 *NeuroImage*, 40(4):1655–71, May 2008. ISSN 1053-8119. doi: 10.1016/j.neuroimage.2008.01.012.
- 815 J. M. Bland and D. Altman. Statistical methods for assessing agreement between two methods of clinical  
816 measurement. *The lancet*, pages 307–310, 1986.
- 817 M. Boccardi, M. Bocchetta, L. G. Apostolova, G. Preboske, N. Robitaille, P. Pasqualetti, L. D. Collins,  
818 S. Duchesne, C. R. Jack, and G. B. Frisoni. Establishing Magnetic Resonance Images Orientation for the  
819 EADC-ADNI Manual Hippocampal Segmentation Protocol. *Journal of neuroimaging : official journal of  
the American Society of Neuroimaging*, pages 1–6, Nov. 2013a. ISSN 1552-6569. doi: 10.1111/jon.12065.
- 820 M. Boccardi, M. Bocchetta, R. Ganzola, N. Robitaille, A. Redolfi, S. Duchesne, C. R. Jack, and G. B. Frisoni.  
821 Operationalizing protocol differences for EADC-ADNI manual hippocampal segmentation. *Alzheimer's &  
dementia : the journal of the Alzheimer's Association*, pages 1–11, May 2013b. ISSN 1552-5279. doi:  
822 10.1016/j.jalz.2013.03.001.
- 823 M. M. Chakravarty, A. F. Sadikot, S. Mongia, G. Bertrand, and D. L. Collins. Towards a multi-modal  
824 atlas for neurosurgical planning. *Medical image computing and computer-assisted intervention : MICCAI  
... International Conference on Medical Image Computing and Computer-Assisted Intervention*, 9(Pt 2):  
825 389–96, Jan. 2006.
- 826 M. M. Chakravarty, A. F. Sadikot, J. Germann, G. Bertrand, and D. L. Collins. Towards a validation  
827 of atlas warping techniques. *Medical image analysis*, 12(6):713–26, Dec. 2008. ISSN 1361-8423. doi:  
828 10.1016/j.media.2008.04.003.
- 829 M. M. Chakravarty, A. F. Sadikot, J. Germann, P. Hellier, G. Bertrand, and D. L. Collins. Comparison  
830 of piece-wise linear, linear, and nonlinear atlas-to-patient warping techniques: analysis of the labeling of  
831 subcortical nuclei for functional neurosurgical applications. *Human brain mapping*, 30(11):3574–95, Nov.  
832 2009. ISSN 1097-0193. doi: 10.1002/hbm.20780.
- 833 M. M. Chakravarty, P. Steadman, M. C. van Eede, R. D. Calcott, V. Gu, P. Shaw, A. Raznahan, D. L. Collins,  
834 and J. P. Lerch. Performing label-fusion-based segmentation using multiple automatically generated tem-  
835 plates. *Human brain mapping*, 34(10):2635–54, Oct. 2013. ISSN 1097-0193. doi: 10.1002/hbm.22092.
- 836

## REFERENCES

---

- 840 G. E. Christensen, S. C. Joshi, and M. I. Miller. Volumetric transformation of brain anatomy. *IEEE*  
841 *transactions on medical imaging*, 16(6):864–77, Dec. 1997. ISSN 0278-0062. doi: 10.1109/42.650882.
- 842 M. Chupin, E. Gérardin, R. Cuingnet, C. Boutet, L. Lemieux, S. Lehéricy, H. Benali, L. Garnero, and  
843 O. Colliot. Fully automatic hippocampus segmentation and classification in Alzheimer’s disease and mild  
844 cognitive impairment applied on data from ADNI. *Hippocampus*, 19(6):579–87, June 2009. ISSN 1098-1063.  
845 doi: 10.1002/hipo.20626.
- 846 D. L. Collins and J. C. Pruessner. Towards accurate, automatic segmentation of the hippocampus and  
847 amygdala from MRI by augmenting ANIMAL with a template library and label fusion. *NeuroImage*, 52  
848 (4):1355–66, Oct. 2010. ISSN 1095-9572. doi: 10.1016/j.neuroimage.2010.04.193.
- 849 D. L. Collins, P. Neelin, T. M. Peters, and A. C. Evans. Automatic 3D intersubject registration of MR  
850 volumetric data in standardized Talairach space. *Journal of computer assisted tomography*, 18(2):192–205,  
851 1994. ISSN 0363-8715.
- 852 D. L. Collins, C. J. Holmes, T. M. Peters, and A. C. Evans. Automatic 3-D model-based neuroanatomical  
853 segmentation. *Human Brain Mapping*, 3(3):190–208, Oct. 1995. ISSN 10659471. doi: 10.1002/hbm.  
854 460030304.
- 855 P. Coupe, V. Fonov, S. Eskildsen, J. Manjón, D. Arnold, and L. Collins. Influence of the training library  
856 composition on a patch-based label fusion method: Application to hippocampus segmentation on the ADNI  
857 dataset. *Alzheimer’s & Dementia*, 7(4):S316, July 2011. ISSN 15525260. doi: 10.1016/j.jalz.2011.05.918.
- 858 P. Coupé, S. F. Eskildsen, J. V. Manjón, V. S. Fonov, and D. L. Collins. Simultaneous segmentation and  
859 grading of anatomical structures for patient’s classification: application to Alzheimer’s disease. *NeuroImage*,  
860 59(4):3736–47, Feb. 2012. ISSN 1095-9572. doi: 10.1016/j.neuroimage.2011.10.080.
- 861 J. G. Csernansky, S. Joshi, L. Wang, J. W. Haller, M. Gado, J. P. Miller, U. Grenander, and M. I. Miller. Hip-  
862 pocampal morphometry in schizophrenia by high dimensional brain mapping. *Proceedings of the National  
863 Academy of Sciences of the United States of America*, 95(19):11406–11411, 1998.
- 864 T. den Heijer, F. V. der Lijn, M. W. Vernooij, M. de Groot, P. J. Koudstaal, a. V. der Lugt, G. P.  
865 Krestin, a. Hofman, W. J. Niessen, and M. M. B. Breteler. Structural and diffusion MRI measures of  
866 the hippocampus and memory performance. *NeuroImage*, 63(4):1782–9, Dec. 2012. ISSN 1095-9572. doi:  
867 10.1016/j.neuroimage.2012.08.067.
- 868 B. Fischl, D. H. Salat, E. Busa, M. Albert, M. Dieterich, C. Haselgrove, A. van der Kouwe, R. Killiany,  
869 D. Kennedy, S. Klaveness, A. Montillo, N. Makris, B. Rosen, and A. M. Dale. Whole brain segmentation:  
870 automated labeling of neuroanatomical structures in the human brain. *Neuron*, 33(3):341–55, Jan. 2002.  
871 ISSN 0896-6273.
- 872 E. Geuze, E. Vermetten, and J. D. Bremner. MR-based in vivo hippocampal volumetrics: 2. Findings in  
873 neuropsychiatric disorders. *Molecular Psychiatry*, 10(2):160, Sept. 2004. doi: 10.1038/sj.mp.4001579.
- 874 I. S. Gousias, D. Rueckert, R. a. Heckemann, L. E. Dyet, J. P. Boardman, a. D. Edwards, and A. Hammers.  
875 Automatic segmentation of brain MRIs of 2-year-olds into 83 regions of interest. *NeuroImage*, 40(2):  
876 672–84, Apr. 2008. ISSN 1053-8119. doi: 10.1016/j.neuroimage.2007.11.034.

## REFERENCES

---

- 877 J. W. Haller, A. Banerjee, G. E. Christensen, M. Gado, S. Joshi, M. I. Miller, Y. Sheline, M. W. Van-  
878 nier, and J. G. Csernansky. Three-dimensional hippocampal MR morphometry with high-dimensional  
879 transformation of a neuroanatomic atlas. *Radiology*, 202(2):504–510, 1997.
- 880 M. Hartig, D. Truran-sacrey, S. Raptentsetsang, N. Schuff, and M. Weiner. USCF FreeSurfer Overview and  
881 QC Ratings. 2010.
- 882 R. A. Heckemann, J. V. Hajnal, P. Aljabar, D. Rueckert, and A. Hammers. Automatic anatomical brain  
883 MRI segmentation combining label propagation and decision fusion. *NeuroImage*, 46(3):726–38, July  
884 2006a. ISSN 1095-9572. doi: 10.1016/j.neuroimage.2009.02.018.
- 885 R. A. Heckemann, J. V. Hajnal, P. Aljabar, D. Rueckert, and A. Hammers. Automatic anatomical brain MRI  
886 segmentation combining label propagation and decision fusion. *NeuroImage*, 33(1):115–26, Oct. 2006b.  
887 ISSN 1053-8119. doi: 10.1016/j.neuroimage.2006.05.061.
- 888 R. A. Heckemann, S. Keihaninejad, P. Aljabar, K. R. Gray, C. Nielsen, D. Rueckert, J. V. Hajnal, and  
889 A. Hammers. Automatic morphometry in Alzheimer’s disease and mild cognitive impairment. *NeuroImage*,  
890 56(4):2024–37, July 2011. ISSN 1095-9572. doi: 10.1016/j.neuroimage.2011.03.014.
- 891 Y.-Y. Hsu, N. Schuff, A.-T. Du, K. Mark, X. Zhu, D. Hardin, and M. W. Weiner. Comparison of automated  
892 and manual MRI volumetry of hippocampus in normal aging and dementia. *Journal of magnetic resonance*  
893 *imaging : JMRI*, 16(3):305–10, Sept. 2002. ISSN 1053-1807. doi: 10.1002/jmri.10163.
- 894 C. R. Jack, M. A. Bernstein, N. C. Fox, P. Thompson, G. Alexander, D. Harvey, B. Borowski, P. J. Britson,  
895 J. L Whitwell, C. Ward, A. M. Dale, J. P. Felmlee, J. L. Gunter, D. L. G. Hill, R. Killiany, N. Schuff,  
896 S. Fox-Bosetti, C. Lin, C. Studholme, C. S. DeCarli, G. Krueger, H. A. Ward, G. J. Metzger, K. T.  
897 Scott, R. Mallozzi, D. Blezek, J. Levy, J. P. Debbins, A. S. Fleisher, M. Albert, R. Green, G. Bartzokis,  
898 G. Glover, J. Mugler, and M. W. Weiner. The Alzheimer’s Disease Neuroimaging Initiative (ADNI): MRI  
899 methods. *Journal of magnetic resonance imaging : JMRI*, 27(4):685–91, Apr. 2008. ISSN 1053-1807. doi:  
900 10.1002/jmri.21049.
- 901 C. R. Jack, F. Barkhof, M. A. Bernstein, M. Cantillon, P. E. Cole, C. Decarli, B. Dubois, S. Duchesne,  
902 N. C. Fox, G. B. Frisoni, H. Hampel, D. L. G. Hill, K. Johnson, J.-F. Mangin, P. Scheltens, A. J. Schwarz,  
903 R. Sperling, J. Suhy, P. M. Thompson, M. Weiner, and N. L. Foster. Steps to standardization and validation  
904 of hippocampal volumetry as a biomarker in clinical trials and diagnostic criterion for Alzheimer’s disease.  
905 *Alzheimer’s & dementia : the journal of the Alzheimer’s Association*, 7(4):474–485.e4, July 2011. ISSN  
906 1552-5279. doi: 10.1016/j.jalz.2011.04.007.
- 907 A. Jeneson and L. Squire. Working memory, long-term memory, and medial temporal lobe function. *Learning*  
908 *& Memory*, 19(1):15–25, 2012. doi: 10.1101/lm.024018.111.
- 909 M. S. Karnik-Henry, L. Wang, D. M. Barch, M. P. Harms, C. Campanella, and J. G. Csernansky. Medial  
910 temporal lobe structure and cognition in individuals with schizophrenia and in their non-psychotic siblings.  
911 *Schizophrenia research*, 138(2-3):128–35, July 2012. ISSN 1573-2509. doi: 10.1016/j.schres.2012.03.015.
- 912 K. K. Leung, J. Barnes, G. R. Ridgway, J. W. Bartlett, M. J. Clarkson, K. Macdonald, N. Schuff, N. C. Fox,  
913 and S. Ourselin. Automated cross-sectional and longitudinal hippocampal volume measurement in mild  
914 cognitive impairment and Alzheimer’s disease. *NeuroImage*, 51(4):1345–59, July 2010. ISSN 1095-9572.  
915 doi: 10.1016/j.neuroimage.2010.03.018.

- 916 C. Loken, D. Gruner, L. Groer, R. Peltier, N. Bunn, M. Craig, T. Henriques, J. Dempsey, C.-H. Yu, J. Chen,  
917 L. J. Dursi, J. Chong, S. Northrup, J. Pinto, N. Knecht, and R. V. Zon. SciNet: Lessons Learned from  
918 Building a Power-efficient Top-20 System and Data Centre. *Journal of Physics: Conference Series*, 256:  
919 012026, Nov. 2010. ISSN 1742-6596. doi: 10.1088/1742-6596/256/1/012026.
- 920 J. M. Lötjönen, R. Wolz, J. R. Koikkalainen, L. Thurfjell, G. Waldemar, H. Soininen, and D. Rueckert. Fast  
921 and robust multi-atlas segmentation of brain magnetic resonance images. *NeuroImage*, 49(3):2352–65,  
922 Mar. 2010. ISSN 1095-9572. doi: 10.1016/j.neuroimage.2009.10.026.
- 923 A. Malla, R. Norman, T. McLean, D. Scholten, and L. Townsend. A Canadian programme for early inter-  
924 vention in non-affective psychotic disorders. *The Australian and New Zealand journal of psychiatry*, 37  
925 (4):407–13, Aug. 2003. ISSN 0004-8674.
- 926 J. Mazziotta, A. Toga, A. Evans, P. Fox, J. Lancaster, K. Zilles, R. Woods, T. Paus, G. Simpson, B. Pike,  
927 C. Holmes, L. Collins, P. Thompson, D. MacDonald, M. Iacoboni, T. Schormann, K. Amunts, N. Palomero-  
928 Gallagher, S. Geyer, L. Parsons, K. Narr, N. Kabani, G. Le Goualher, J. Feidler, K. Smith, D. Boomsma,  
929 H. Hulshoff Pol, T. Cannon, R. Kawashima, and B. Mazoyer. A four-dimensional probabilistic atlas of  
930 the human brain. *Journal of the American Medical Informatics Association : JAMIA*, 8(5):401–30. ISSN  
931 1067-5027.
- 932 J. Mazziotta, A. Toga, A. Evans, P. Fox, J. Lancaster, K. Zilles, R. Woods, T. Paus, G. Simpson,  
933 B. Pike, C. Holmes, L. Collins, P. Thompson, D. MacDonald, M. Iacoboni, T. Schormann, K. Amunts,  
934 N. Palomero-Gallagher, S. Geyer, L. Parsons, K. Narr, N. Kabani, G. Le Goualher, D. Boomsma,  
935 T. Cannon, R. Kawashima, and B. Mazoyer. A probabilistic atlas and reference system for the hu-  
936 man brain: International Consortium for Brain Mapping (ICBM). *Philosophical transactions of the Royal  
937 Society of London. Series B, Biological sciences*, 356(1412):1293–322, Aug. 2001. ISSN 0962-8436. doi:  
938 10.1098/rstb.2001.0915.
- 939 J. C. Mazziotta, A. W. Toga, A. Evans, P. Fox, and J. Lancaster. A probabilistic atlas of the human brain:  
940 theory and rationale for its development. The International Consortium for Brain Mapping (ICBM).  
941 *NeuroImage*, 2(2):89–101, June 1995. ISSN 1053-8119.
- 942 J. H. Morra, Z. Tu, L. G. Apostolova, A. E. Green, C. Avedissian, S. K. Madsen, N. Parikshak, X. Hua, A. W.  
943 Toga, C. R. Jack, M. W. Weiner, and P. M. Thompson. Validation of a fully automated 3D hippocampal  
944 segmentation method using subjects with Alzheimer's disease mild cognitive impairment, and elderly  
945 controls. *NeuroImage*, 43(1):59–68, Oct. 2008. ISSN 1095-9572. doi: 10.1016/j.neuroimage.2008.07.003.
- 946 S. Mueller, L. Stables, A. Du, and N. Schuff. Measurement of hippocampal subfields and age-related changes  
947 with high resolution MRI at 4T. *Neurobiology of ...*, 1(5):719–726, 2007. doi: 10.1016/j.neurobiolaging.  
948 2006.03.007.
- 949 S. G. Mueller and M. W. Weiner. Selective effect of age, Apo e4, and Alzheimer's disease on hippocampal  
950 subfields. *Hippocampus*, 19(6):558–64, June 2009. ISSN 1098-1063. doi: 10.1002/hipo.20614.
- 951 K. L. Narr, P. M. Thompson, P. Szeszko, D. Robinson, S. Jang, R. P. Woods, S. Kim, K. M. Hayashi,  
952 D. Asunction, A. W. Toga, and R. M. Bilder. Regional specificity of hippocampal volume reductions  
953 in first-episode schizophrenia. *NeuroImage*, 21(4):1563–75, Apr. 2004. ISSN 1053-8119. doi: 10.1016/j.  
954 neuroimage.2003.11.011.

## REFERENCES

---

- 955 S. M. Nestor, E. Gibson, F.-Q. Gao, A. Kiss, and S. E. Black. A Direct Morphometric Comparison of Five  
956 Labeling Protocols for Multi-Atlas Driven Automatic Segmentation of the Hippocampus in Alzheimer's  
957 Disease. *NeuroImage*, Nov. 2012. ISSN 1095-9572. doi: 10.1016/j.neuroimage.2012.10.081.
- 958 Z. Pausova, T. Paus, M. Abrahamowicz, J. Almerigi, N. Arbour, M. Bernard, D. Gaudet, P. Hanzalek,  
959 P. Hamet, A. C. Evans, M. Kramer, L. Laberge, S. M. Leal, G. Leonard, J. Lerner, R. M. Lerner, J. Math-  
960 ieu, M. Perron, B. Pike, A. Pitiot, L. Richer, J. R. Séguin, C. Syme, R. Toro, R. E. Tremblay, S. Veillette,  
961 and K. Watkins. Genes, maternal smoking, and the offspring brain and body during adolescence: design  
962 of the Saguenay Youth Study. *Human brain mapping*, 28(6):502–18, June 2007. ISSN 1065-9471. doi:  
963 10.1002/hbm.20402.
- 964 K. M. Pohl, S. Bouix, M. Nakamura, T. Rohlfing, R. W. McCarley, R. Kikinis, W. E. L. Grimson, M. E.  
965 Shenton, and W. M. Wells. A hierarchical algorithm for MR brain image parcellation. *IEEE transactions  
966 on medical imaging*, 26(9):1201–12, Sept. 2007. ISSN 0278-0062. doi: 10.1109/TMI.2007.901433.
- 967 J. Poppenk and M. Moscovitch. A Hippocampal Marker of Recollection Memory Ability among Healthy  
968 Young Adults: Contributions of Posterior and Anterior Segments. *Neuron*, 72(6):931–937, Dec. 2011.  
969 ISSN 0896-6273. doi: 10.1016/j.neuron.2011.10.014.
- 970 S. Powell, V. A. Magnotta, H. Johnson, V. K. Jammalamadaka, R. Pierson, and N. C. Andreasen. Registration  
971 and machine learning-based automated segmentation of subcortical and cerebellar brain structures.  
972 *NeuroImage*, 39(1):238–47, Jan. 2008. ISSN 1053-8119. doi: 10.1016/j.neuroimage.2007.05.063.
- 973 J. C. Pruessner, L. M. Li, W. Serles, M. Pruessner, D. L. Collins, N. Kabani, S. Lupien, and A. C. Evans.  
974 Volumetry of hippocampus and amygdala with high-resolution MRI and three-dimensional analysis soft-  
975 ware: minimizing the discrepancies between laboratories. *Cerebral cortex (New York, N.Y. : 1991)*, 10  
976 (4):433–42, Apr. 2000. ISSN 1047-3211.
- 977 S. Robbins, A. C. Evans, D. L. Collins, and S. Whitesides. Tuning and comparing spatial normalization  
978 methods. *Medical image analysis*, 8(3):311–23, Sept. 2004. ISSN 1361-8415. doi: 10.1016/j.media.2004.  
979 06.009.
- 980 N. Robitaille and S. Duchesne. Label fusion strategy selection. *International journal of biomedical imaging*,  
981 2012:431095, Jan. 2012. ISSN 1687-4196. doi: 10.1155/2012/431095.
- 982 M. R. Sabuncu, R. S. Desikan, J. Sepulcre, B. T. T. Yeo, H. Liu, N. J. Schmansky, M. Reuter, M. W.  
983 Weiner, R. L. Buckner, R. a. Sperling, and B. Fischl. The dynamics of cortical and hippocampal atrophy  
984 in Alzheimer disease. *Archives of neurology*, 68(8):1040–8, Aug. 2011. ISSN 1538-3687. doi: 10.1001/  
985 archneurol.2011.167.
- 986 W. B. Scoville and B. Milner. Loss of recent memory after bilateral hippocampal lesions. 1957. *The Journal  
987 of neuropsychiatry and clinical neurosciences*, 12(1):103–113, 2000.
- 988 J. Shao. Linear Model Selection by Cross-validation. *Journal of the American Statistical Association*, 88  
989 (422):486–494, June 1993. ISSN 0162-1459. doi: 10.1080/01621459.1993.10476299.
- 990 J. G. Sled, a. P. Zijdenbos, and a. C. Evans. A nonparametric method for automatic correction of intensity  
991 nonuniformity in MRI data. *IEEE transactions on medical imaging*, 17(1):87–97, Feb. 1998. ISSN 0278-  
992 0062. doi: 10.1109/42.668698.

## REFERENCES

---

- 993 C. Studholme, D. Hill, and D. Hawkes. An overlap invariant entropy measure of 3D medical image alignment.  
994 *Pattern Recognition*, 32(1):71–86, Jan. 1999. ISSN 00313203. doi: 10.1016/S0031-3203(98)00091-0.
- 995 C. Studholme, E. Novotny, I. G. Zubal, and J. S. Duncan. Estimating tissue deformation between functional  
996 images induced by intracranial electrode implantation using anatomical MRI. *NeuroImage*, 13(4):561–76,  
997 Apr. 2001. ISSN 1053-8119. doi: 10.1006/nimg.2000.0692.
- 998 F. van der Lijn, T. den Heijer, M. M. B. Breteler, and W. J. Niessen. Hippocampus segmentation in MR  
999 images using atlas registration, voxel classification, and graph cuts. *NeuroImage*, 43(4):708–20, Dec. 2008.  
1000 ISSN 1095-9572. doi: 10.1016/j.neuroimage.2008.07.058.
- 1001 K. Van Leemput, A. Bakkour, T. Benner, G. Wiggins, L. L. Wald, J. Augustinack, B. C. Dickerson, P. Gol-  
1002 land, and B. Fischl. Automated segmentation of hippocampal subfields from ultra-high resolution *in vivo*  
1003 MRI. *Hippocampus*, 19(6):549–57, June 2009. ISSN 1098-1063. doi: 10.1002/hipo.20615.
- 1004 H. Wang, J. W. Suh, J. Pluta, M. Altinay, and P. Yushkevich. Optimal weights for multi-atlas label fusion.  
1005 *Information processing in medical imaging : proceedings of the ... conference*, 22:73–84, Jan. 2011. ISSN  
1006 1011-2499.
- 1007 S. K. Warfield, K. H. Zou, and W. M. Wells. Simultaneous truth and performance level estimation (STAPLE):  
1008 an algorithm for the validation of image segmentation. *IEEE transactions on medical imaging*, 23(7):903–  
1009 21, July 2004. ISSN 0278-0062. doi: 10.1109/TMI.2004.828354.
- 1010 J. L. Winterburn, J. C. Pruessner, S. Chavez, M. M. Schira, N. J. Lobaugh, A. N. Voineskos, and M. M.  
1011 Chakravarty. A novel *in vivo* atlas of human hippocampal subfields using high-resolution 3 T magnetic  
1012 resonance imaging. *NeuroImage*, 74:254–65, July 2013. ISSN 1095-9572. doi: 10.1016/j.neuroimage.2013.  
1013 02.003.
- 1014 L. E. M. Wisse, L. Gerritsen, J. J. M. Zwanenburg, H. J. Kuijf, P. R. Luijten, G. J. Biessels, and M. I.  
1015 Geerlings. Subfields of the hippocampal formation at 7 T MRI: *in vivo* volumetric assessment. *NeuroImage*,  
1016 61(4):1043–9, July 2012. ISSN 1095-9572. doi: 10.1016/j.neuroimage.2012.03.023.
- 1017 J. Wixted and L. Squire. The medial temporal lobe and the attributes of memory. *Trends in cognitive  
1018 sciences*, 15(5):210–217, 2011. doi: 10.1016/j.tics.2011.03.005.
- 1019 R. Wolz, P. Aljabar, J. V. Hajnal, A. Hammers, and D. Rueckert. LEAP: learning embeddings for atlas  
1020 propagation. *NeuroImage*, 49(2):1316–25, Jan. 2010. ISSN 1095-9572. doi: 10.1016/j.neuroimage.2009.09.  
1021 069.
- 1022 B. T. Wyman, D. J. Harvey, K. Crawford, M. A. Bernstein, O. Carmichael, P. E. Cole, P. K. Crane,  
1023 C. Decarli, N. C. Fox, J. L. Gunter, D. Hill, R. J. Killiany, C. Pachai, A. J. Schwarz, N. Schuff, M. L.  
1024 Senjem, J. Suhy, P. M. Thompson, M. Weiner, and C. R. Jack. Standardization of analysis sets for reporting  
1025 results from ADNI MRI data. *Alzheimer's & dementia : the journal of the Alzheimer's Association*, Oct.  
1026 2012. ISSN 1552-5279. doi: 10.1016/j.jalz.2012.06.004.
- 1027 J. Yelnik, E. Bardinet, D. Dormont, G. Malandain, S. Ourselin, D. Tandé, C. Karachi, N. Ayache, P. Cornu,  
1028 and Y. Agid. A three-dimensional, histological and deformable atlas of the human basal ganglia. I. Atlas  
1029 construction based on immunohistochemical and MRI data. *NeuroImage*, 34(2):618–38, Jan. 2007. ISSN  
1030 1053-8119. doi: 10.1016/j.neuroimage.2006.09.026.

## *REFERENCES*

---

- 1031 P. A. Yushkevich, B. B. Avants, J. Pluta, S. Das, D. Minkoff, D. Mechanic-Hamilton, S. Glynn, S. Pickup,  
1032 W. Liu, J. C. Gee, M. Grossman, and J. A. Detre. A high-resolution computational atlas of the human  
1033 hippocampus from postmortem magnetic resonance imaging at 9.4 T. *NeuroImage*, 44(2):385–98, Jan.  
1034 2009. ISSN 1095-9572. doi: 10.1016/j.neuroimage.2008.08.042.
- 1035 P. A. Yushkevich, H. Wang, J. Pluta, S. R. Das, C. Craige, B. B. Avants, M. W. Weiner, and S. Mueller.  
1036 Nearly automatic segmentation of hippocampal subfields in in vivo focal T2-weighted MRI. *NeuroImage*,  
1037 53(4):1208–24, Dec. 2010. ISSN 1095-9572. doi: 10.1016/j.neuroimage.2010.06.040.

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

Jon Pipitone<sup>1</sup>, Min Tae M. Park<sup>1</sup>, Julie Winterburn<sup>1</sup>, Tristram A. Lett<sup>1,9</sup>, Jason P. Lerch<sup>2,3</sup>, Jens C. Pruessner<sup>4</sup>, Martin Lepage<sup>4,5</sup>, Aristotle N. Voineskos<sup>1,6,9</sup>, M. Mallar Chakravarty<sup>1,6,7,8</sup> and the Alzheimer's Disease Neuroimaging Initiative\*

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

<sup>2</sup> *Neurosciences and Mental Health Laboratory, Hospital for Sick Children, Toronto, ON, Canada*

<sup>3</sup> *Department of Medical Biophysics, University of Toronto, Toronto, ON, Canada*

<sup>4</sup> *Douglas Mental Health University Institute, Verdun, QC, Canada*

<sup>5</sup> *Department of Psychiatry, McGill University, Montreal, QC, Canada*

<sup>6</sup> *Department of Psychiatry, University of Toronto, Toronto, ON, Canada*

<sup>7</sup> *Institute of Biomaterials and Biomedical Engineering, University of Toronto, Toronto, ON, Canada*

<sup>8</sup> *Rotman Research Institute, Baycrest, Toronto, ON, Canada*

<sup>9</sup> *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 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 (**M**ultiple **A**utomatically **G**enerated **T**emplates), for the automatic segmentation of the hippocampus that minimizes the number of atlases needed while still achieving similar agreement to multi-atlas approaches. Thus, our method acts as an reliable multi-atlas approach when using special or 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 nonlinear image registration. The resulting segmentations are then propagated to each target image and fused using a label fusion method.

We conduct two separate Monte Carlo cross-validation experiments comparing MAGeT-Brain and multi-atlas whole hippocampal segmentation using differing atlas and template library sizes, and registration and label fusion methods. The first experiment is a 10-fold validation (per parameter setting) over 60 subjects taken from the Alzheimer's Disease Neuroimaging Database (ADNI), and the second is a five-fold validation over 81 subjects having had a first episode of psychosis. In both cases, automated segmentations are compared with manual segmentations following the Pruessner-protocol. Using

\*Data used in preparation of this article were obtained from the Alzheimer's Disease Neuroimaging Initiative (ADNI) database (adni.loni.usc.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.usc.edu/wp-content/uploads/how\\_to\\_apply/ADNI\\_Acknowledgement\\_List.pdf](http://adni.loni.usc.edu/wp-content/uploads/how_to_apply/ADNI_Acknowledgement_List.pdf)

20 the best settings found from these experiments, we segment 246 images of the ADNI1:Complete 1Yr  
21 1.5T dataset and compare these with segmentations from existing automated methods: FSL FIRST,  
22 FreeSurfer, MAPER, and SNT. Finally, we conduct a leave-one-out cross-validation (LOOCV) of hip-  
23 pocampal subfield segmentation in standard 3T T1-weighted images, using five high-resolution manually  
24 segmented atlases (Winterburn et al., 2013).

25 **Results:** In the ADNI cross-validation, using 9 atlases MAGeT-Brain achieves a mean Dice's Sim-  
26 ilarity Coefficient (DSC) score of 0.869 with respect to manual whole hippocampus segmentations, and  
27 also exhibits significantly lower variability in DSC scores than multi-atlas segmentation. In the younger,  
28 psychosis dataset, MAGeT-Brain achieves a mean DSC score of 0.892 and produces volumes which  
29 agree with manual segmentation volumes better than those produced by the FreeSurfer and FSL FIRST  
30 methods (mean difference in volume:  $80mm^3$ ,  $1600mm^3$ , and  $800mm^3$ , respectively). Similarly, in the  
31 ADNI1:Complete 1Yr 1.5T dataset, MAGeT-Brain produces hippocampal segmentations well correlated  
32 ( $r > 0.85$ ) with SNT semi-automated reference volumes within disease categories, and shows a conserva-  
33 tive bias and a mean difference in volume of  $250mm^3$  across the entire dataset, compared with FreeSurfer  
34 and FSL FIRST which both overestimate volume differences by  $2600mm^3$  and  $2800mm^3$  on average, re-  
35 spectively. Finally, MAGeT-Brain segments the CA1, CA4/DG and subiculum subfields on standard 3T  
36 T1-weighted resolution images with DSC overlap scores of 0.56, 0.65, and 0.58, respectively, relative to  
37 manual segmentations.

38 **Conclusion:** We demonstrate that MAGeT-Brain produces consistent whole hippocampal segmen-  
39 tations using only 9 atlases, or fewer, with various hippocampal definitions, disease populations, and  
40 image acquisition types. Additionally, we show that MAGeT-Brain identifies hippocampal subfields in  
41 standard 3T T1-weighted images with overlap scores comparable to competing methods.

#### 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](mailto:jon.pipitone@camh.ca); [mallar.chakravarty@camh.ca](mailto:mallar.chakravarty@camh.ca)

## 43 1 Introduction

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

52 The gold standard for neuroanatomical image segmentation is manual delineation by an expert human  
53 rater. However, with the availability of increasingly large MRI datasets the time and expertise required  
54 for manual segmentation becomes prohibitive (Mazziotta et al., 1995, 2001; Mazziotta et al.; Pausova et al.,  
55 2007). This effort is complicated by the fact that there is significant variation between segmentation protocols

---

56 with respect to specific anatomical boundaries of the hippocampus (Geuze et al., 2004) and this has led to  
57 efforts to create an unified hippocampal segmentation protocol (Jack et al., 2011; Boccardi et al., 2013b,a).  
58 In addition, there is controversy over the appropriate manual segmentation protocol to use in a particular  
59 imaging study (Nestor et al., 2012). Thus, a segmentation algorithm that can easily adapt to different  
60 manual segmentation definitions would be of significant benefit to the neuroimaging community.

61 Automated segmentation techniques that are reliable, objective, and reproducible can be considered  
62 complementary to manual segmentation. In the case of classical model-based segmentation methods (Haller  
63 et al., 1997; Csernansky et al., 1998), an MRI atlas that was previously manually labelled by an expert rater  
64 is matched to target images using nonlinear registration methods. The resulting nonlinear transformation is  
65 applied to the manual labels (i.e. *label propagation*) to warp them into the target image space. While this  
66 methodology has been used successfully in several contexts (Chakravarty et al., 2008, 2009; Collins et al.,  
67 1995; Haller et al., 1997), it is limited by the error in the estimated nonlinear transformation itself, partial  
68 volume effects in label resampling, and irreconcilable differences between the neuroanatomy represented  
69 within the atlas and target images.

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

74 More recently, many groups have used multiple atlases to improve overall segmentation reliability (i.e.  
75 multi-atlas segmentation) over model-based approaches (Heckemann et al., 2006a, 2011; Collins and Pruessner,  
76 2010; Lötjönen et al., 2010; Aljabar et al., 2009; Leung et al., 2010; Wolz et al., 2010). Each atlas  
77 image is registered to a target image, and label propagation is performed to produce several labellings of  
78 the target image (one from each atlas). A *label fusion* technique, such as voxel-wise voting, is used to merge  
79 these labels into the definitive segmentation for the target. In addition, weighted voting procedures that  
80 use *atlas selection* techniques are often used to exclude atlases from label fusion that are dissimilar to a  
81 target image in order to reduce error from unrepresentative anatomy (Aljabar et al., 2009). This involves  
82 the selection of a subset of atlases using a similarity metric such as cross-correlation (Aljabar et al., 2009) or  
83 normalized mutual information. Such selection has the added benefit of significantly reducing the number  
84 of nonlinear registrations. For example Collins and Pruessner (2010) demonstrated that only 14 atlases, se-  
85 lected based on highest similarity between medial temporal lobe neuroanatomy as evaluated by normalized  
86 mutual information (Studholme et al., 1999) from a library of 80 atlases, were required to achieve favourable  
87 segmentations of the hippocampus. Also, several methods have been explored for label fusion. For example,  
88 the STAPLE algorithm (Simultaneous Truth And Performance Level Estimation; Warfield et al. (2004)) uses  
89 an expectation-maximization framework to compute a probabilistic segmentation from a set of competing  
90 segmentations, or the work of Coupé et al. (2012) who show that a subset of segmentations can be estimated  
91 using metrics, such as the sum of squared differences in the regions of interest to be segmented.

92 However, many of these methods require significant investment of time and resources for the creation  
93 of the atlas library ranging between 30 (Heckemann et al., 2006a) and 80 (Collins and Pruessner, 2010)  
94 manually segmented atlases. This strategy has the main drawback of being inflexible as it does not easily  
95 accommodate varying the definition of the hippocampal anatomy (such as the commonly used heuristic of  
96 subdividing the hippocampus into head, body, and tail (Poppenk and Moscovitch, 2011; Pruessner et al.,  
97 2000)). Furthermore, none of these methods have demonstrated sufficient flexibility to accommodate atlases  
98 that are somehow exceptional such as those derived from serial histological data (Chakravarty et al., 2006;

---

99 Yelnik et al., 2007) or high-resolution MRI data that enables robust identification of hippocampal subfields  
100 (Winterburn et al., 2013; Yushkevich et al., 2009; Mueller and Weiner, 2009; Van Leemput et al., 2009;  
101 Wisse et al., 2012). Due to the recent availability of the latter, there has been increased interest in the use  
102 of probabilistic methods for the identification of the hippocampal subfields on standard T1-weighted images.  
103 Our group recently demonstrated that through use of an intermediary automated segmentation stage, robust  
104 and reliable segmentation of the striatum, pallidum, and thalamus using a single atlas derived from serial  
105 histological data is possible (Chakravarty et al., 2013). The novelty of this manuscript is the extension of  
106 our multi-atlas methodology to the segmentation of hippocampus. Additionally, in this paper we rigorously  
107 explore the effects of using multiple input atlases, of varying the size of the template library constructed,  
108 and registration and label fusion methods. As a result, we aim to demonstrate that it is indeed possible to  
109 reliably apply the segmentation represented in a very small set of segmented input atlases to an unlabelled  
110 target image set.

111 Of particular relevance to the present work is the LEAP algorithm (Learning Embeddings for Atlas  
112 Propagation; Wolz et al. (2010)) because of its focus on performing multi-atlas segmentation with a limited  
113 number of input atlases. The LEAP algorithm is a clever modification to the basic multi-atlas strategy in  
114 which an atlas library is grown, beginning with a set of manually labelled atlases, by successively incorpo-  
115 rating unlabelled target images once they themselves have been labelled using multi-atlas techniques. The  
116 sequence in which target images are labelled is chosen so that the similarity between the atlas images and the  
117 target images is minimised at each step, effectively allowing for deformations between very dissimilar images  
118 to be broken up into sequences of smaller deformations. Although Wolz et al. (2010) begin with an atlas  
119 library of 30 MR images, this method could theoretically work using a much smaller atlas library. In their  
120 validation, LEAP was used to segment the whole hippocampus in the ADNI1 baseline dataset, achieving a  
121 mean Dice score of 0.85 against semi-automated segmentations.

122 Also of interest to this manuscript are the methods that attempt to define hippocampal subfields using  
123 standard T1- or T2-weighted data, of which there are few. Van Leemput et al. (2009) demonstrate that  
124 the applicability of hippocampal subfield segmentation in T1-weighted images by Bayesian techniques using  
125 Markov random field shape priors learned from 10 manual segmentations. This work, available as part of  
126 the FreeSurfer package, is limited as the segmentation omits the tail of the hippocampus and the protocol  
127 has yet to be fully validated. Yushkevich et al. (2009) manually segment hippocampal subfields on high-  
128 resolution (either 0.2mm-isotropic or 0.2mm × 0.3mm × 0.2mm resolution voxels) T2-weighted MR images  
129 acquired from five post-mortem medial temporal lobe samples. Then, using nonlinear registration guided by  
130 shape-based models of the subfield segmentations and manually derived hippocampus masks of the target  
131 images, the authors demonstrate accurate parcellation of hippocampal subfields, with respect to manual  
132 segmentations, in clinical 3T T1-weighted MRI volumes. Using multi-atlas with bias correction techniques,  
133 Yushkevich et al. (2010) demonstrate a semi-automated method of subfield segmentation on in vivo focal  
134 T2-weighted MR acquisitions of the temporal lobe. Manual input is only needed to mark divisions between  
135 the head, body and tail of the hippocampus on target images.

136 In this paper we describe a thorough validation of the MAGeT-Brain algorithm for the fully automatic  
137 segmentation of the hippocampus and a proof-of-concept validation of its application to the segmentation  
138 of hippocampal subfields in standard T1-weighted images. First, we address the very idea of generating a  
139 template library from a limited number of input atlases (Chakravarty et al., 2013) for whole hippocampus  
140 segmentation by conducting a multi-fold validation experiment over a range of atlas and template library  
141 sizes, registration and label fusion methods. This type of validation is done first on a subset of the Alzheimer's

---

<sup>142</sup> Disease Neuroimaging Initiative (ADNI) dataset with manual segmentations Pruessner-protocol, and then  
<sup>143</sup> replicated on a first episode psychosis patient dataset to determine the behaviour of MAGeT-Brain when  
<sup>144</sup> segmenting younger and differently diseased subjects. Next, we compare MAGeT-Brain with other popular  
<sup>145</sup> segmentation algorithms (FreeSurfer, FSL FIRST, MAPER, and SNT) on all the images available in the  
<sup>146</sup> ADNI1:Complete 1Yr 1.5T sample. Lastly, using the optimal parameter settings for MAGeT-Brain found  
<sup>147</sup> from the previous experiments, we investigate hippocampal subfield segmentation by conducting a leave-  
<sup>148</sup> one-out validation using the Winterburn et al. (2013) manually segmented high-resolution MR atlases.

<sup>149</sup> **2 The MAGeT-Brain Algorithm**

<sup>150</sup> In this paper, we use the term *label* to mean any segmentation (manual or derived) of an MR image. *Label*  
<sup>151</sup> *propagation* is the process by which two images are registered and the resulting transformation is applied  
<sup>152</sup> to the labels from one image to bring them into alignment with the other image. We use the term *atlas*  
<sup>153</sup> to mean a manually segmented image, and the term *template* to mean an automatically segmented image  
<sup>154</sup> (i.e. via label propagation). The terms *atlas library* and *template library* describe any set of such images.  
<sup>155</sup> Additionally, we use the term *target* to refer to an unlabelled image that is undergoing segmentation.

<sup>156</sup> The simplest form of multi-atlas segmentation, which we call *basic multi-atlas segmentation*, involves three  
<sup>157</sup> steps. First, each labelled input image (i.e. atlas or template) is registered to an unlabelled target image.  
<sup>158</sup> Second, the labels from each image are propagated to the target image space. Third, the labels are combined  
<sup>159</sup> into a single label by label fusion (Heckemann et al., 2006a, 2011). The basic multi-atlas segmentation method  
<sup>160</sup> is described in detail in other publications (Collins and Pruessner, 2010; Heckemann et al., 2011; Aljabar  
<sup>161</sup> et al., 2009). When only a single atlas is used, basic multi-atlas segmentation degenerates into model-based  
<sup>162</sup> segmentation: labels are propagated from the atlas to a target, and no label fusion is needed.

<sup>163</sup> The MAGeT-Brain (**M**ultiple **A**utomatically **G**enerated **T**emplates) algorithm creates a large template  
<sup>164</sup> library given a much smaller sized input atlas library and then uses this template library in basic multi-atlas  
<sup>165</sup> segmentation to segment a set of input target images. The images used in the template library are selected  
<sup>166</sup> from the input images, either arbitrarily or so as to reflect the neuroanatomy or demographics of the target  
<sup>167</sup> set as a whole (for instance, by sampling equally from cases and controls). The template library images are  
<sup>168</sup> automatically labelled by each of the atlases via label propagation. Effectively, basic multi-atlas segmentation  
<sup>169</sup> is then conducted using the template library to segment the entire set of target images (including the target  
<sup>170</sup> images used in the construction of the template library). Since each template library image has multiple  
<sup>171</sup> labels (one from each atlas), the final number of labels to be fused for each target may be quite large (i.e.  
<sup>172</sup> # of atlas  $\times$  # of templates).

<sup>173</sup> Figure 1 illustrates the MAGeT-Brain algorithm graphically. Source code for MAGeT-Brain can be found  
<sup>174</sup> at <http://github.com/pipitone/MAGeTbrain>.

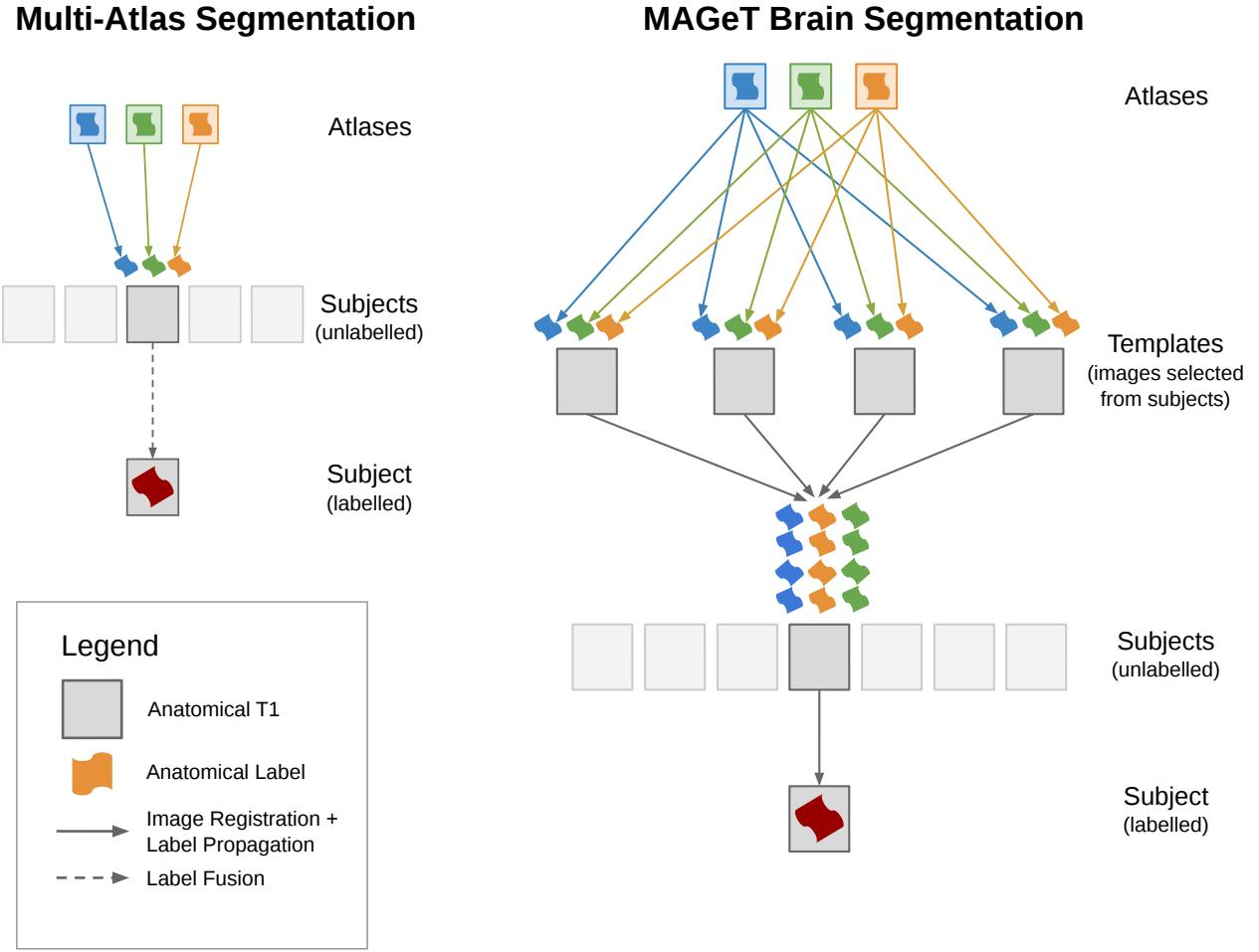


Figure 1: A schematic illustration of basic multi-atlas segmentation and MAGeT-Brain segmentation. In multi-atlas segmentation, manual labels from atlas images are warped (propagated) into subject space by applying the transformations estimated from nonlinear image registration. The resulting candidate labels from all atlas images are then fused to create a final segmentation. In MAGeT-Brain segmentation, a template library is created by sampling (either randomly or representatively) from the subject images. Atlas labels are propagated to all template images and then to each subject image (including those used in the template library). The candidate labels for a subject are then fused into a final segmentation.

---

### 175    3 Experiments

176    The following section describes experiments conducted to assess the segmentation quality of the MAGeT-  
177    Brain algorithm:

- 178    • Experiment 1 investigates MAGeT-Brain whole hippocampus segmentation of aging and Alzheimer's  
179    diseased subjects over a wide range of parameter settings using a Monte Carlo cross-validation design.  
180    The results of this experiment enable us to choose the parameter settings offering the best performance  
181    for use in subsequent experiments.
- 182    • Experiment 2 is a similar cross-validation to explore MAGeT-Brain segmentations on the brain im-  
183    ages of young, first episode psychosis patients. In addition, MAGeT-Brain segmentations with two  
184    different atlas segmentation protocols are compared to automated segmentations by the FSL FIRST  
185    and FreeSurfer algorithms. The results of this experiment combined with the previous experiment  
186    establishes parameter settings that do not overfit to the neuroanatomical features of a specific patient  
187    cohort.
- 188    • Experiment 3 bridges MAGeT-Brain with the existing segmentation literature by comparing MAGeT-  
189    Brain whole hippocampus segmentations with those of several well-known automated methods (FreeSurfer,  
190    FSL FIRST, MAPER, SNT) on the entire ADNI1:Complete 1Yr 1.5T image dataset consisting of  
191    246 brain images of subjects diagnosed as cognitively normal, having mild cognitive impairment, or  
192    Alzheimer's disease.
- 193    • Experiment 4 assesses hippocampal subfield segmentation quality in a leave-one-out cross-validation  
194    on the five high-resolution manually segmented Winterburn MR atlases (Winterburn et al., 2013).

#### 195    3.1    Experiment 1: Whole Hippocampus Segmentation Cross-Validation — Alzheimer's 196    Disease

197    In this experiment we explore the very idea of generating a template library for multi-atlas-based segmenta-  
198    tion from a small number of input atlases. To do so, we conduct repeated cross-validations of MAGeT-Brain  
199    whilst varying the composition and sizes of the atlas and template libraries used, as well as varying the  
200    registration algorithm and label fusion method. The dataset used in this experiment is images from the  
201    ADNI dataset (Jack et al., 2008) along with whole hippocampus labels manually segmented following the  
202    Pruessner-protocol (Pruessner et al., 2000).

203    Note, in the Supplementary Materials we have replicated this experiment using the SNT semi-automated  
204    segmentations included as part of the ADNI dataset.

##### 205    3.1.1    Experiment 1: Materials and Methods

206    **ADNI1:Complete 1Yr 1.5T dataset** Data used in the preparation of this article were obtained from the  
207    Alzheimer's Disease Neuroimaging Initiative (ADNI) database ([adni.loni.usc.edu](http://adni.loni.usc.edu)). The ADNI was launched  
208    in 2003 by the National Institute on Aging (NIA), the National Institute of Biomedical Imaging and Bio-  
209    engineering (NIBIB), the Food and Drug Administration (FDA), private pharmaceutical companies and  
210    non-profit organizations, as a \$60 million, 5-year public-private partnership. The primary goal of ADNI has  
211    been to test whether serial magnetic resonance imaging (MRI), positron emission tomography (PET), other

Table 1: **ADNI1 cross-validation subset demographics.** CN - Cognitively Normal. LMCI - Late-onset Mild Cognitive Impairment. AD - Alzheimer’s Disease. CDR-SB - Clinical Dementia Rating-Sum of Boxes. ADAS - Alzheimer’s Disease Assessment Scale. MMSE - Mini-Mental State Examination.

	CN N = 20	LMCI N = 20	AD N = 20	Combined 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.

212 biological markers, and clinical and neuropsychological assessment can be combined to measure the progression  
 213 of mild cognitive impairment (MCI) and early Alzheimer’s disease (AD). Determination of sensitive and  
 214 specific markers of very early AD progression is intended to aid researchers and clinicians to develop new  
 215 treatments and monitor their effectiveness, as well as lessen the time and cost of clinical trials.

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

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

230 Clinical, demographic and pre-processed T1-weighted MRI were downloaded by the authors from the  
 231 ADNI database ([adni.loni.usc.edu](http://adni.loni.usc.edu)) between March 2012 and August 2012. The image dataset used was  
 232 the *ADNI1:Complete 1Yr 1.5T* standardized dataset available from ADNI <sup>1</sup> (Wyman et al., 2012). This  
 233 image collection contains uniformly pre-processed images which have been designated to be the “best” after  
 234 quality control. All images were acquired using 1.5T scanners (General Electric Healthcare, Philips Medical  
 235 Systems or Siemens Medical Solutions) at multiple sites using the protocol described in Jack et al. (2008).  
 236 Representative 1.5T imaging parameters were TR = 2400ms, TI = 1000ms, TE = 3.5ms, flip angle = 8°,  
 237 field of view = 240 x 240mm, a 192 × 192 × 166 matrix (*x*, *y*, and *z* directions) yielding voxel dimensions of  
 238 1.25mm × 1.25mm × 1.2mm.

239 **Experiment details** Monte Carlo Cross-Validation (MCCV), also known as repeated random sub-sampling  
 240 cross-validation, consists of repeated rounds of validation conducted on a fixed dataset (Shao, 1993). In each

<sup>1</sup><http://adni.loni.usc.edu/methods/mri-analysis/adni-standardized-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

241 round, the dataset is randomly partitioned into a training set and a validation set. The method to be  
 242 validated is then given the training data, and its output is compared with the validation set.

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

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

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

257 **Automatic Normalization and Image Matching and Anatomical Labeling (ANIMAL)** The  
 258 ANIMAL algorithm carries out image registration in two phases. In the first, a 12-parameter linear trans-  
 259 formation (3 translations, rotations, scales, shears) is estimated between images using an algorithm that  
 260 maximizes the correlation between blurred MR intensities and gradient magnitude over the whole brain  
 261 (Collins et al., 1994). In the second phase, nonlinear registration is completed using the ANIMAL algorithm  
 262 (Collins et al., 1995): an iterative procedure that estimates a 3D deformation field between two MR images.  
 263 At first, large deformations are estimated using a blurred version of the input data. These larger deforma-  
 264 tions are then input to subsequent steps where the fit is refined by estimating smaller deformations on data  
 265 blurred with a Gaussian kernel with a smaller full width at half maximum (FWHM). The final transfor-  
 266 mation is a set of local translations defined on a bed of equally spaced nodes that were estimated through  
 267 the optimization of the correlation coefficient. For the purposes of this work we used the regularization  
 268 parameters optimized in Robbins et al. (2004), displayed in Table 2.

269 **Automatic Normalization Tools (ANTS)** ANTS is a diffeomorphic registration algorithm which  
 270 provides great flexibility over the choice of transformation model, objective function, and the consistency of  
 271 the final transformation (Avants et al., 2008). The transformation is estimated in a hierarchical fashion where  
 272 the MRI data is subsampled, allowing large deformations to be estimated and successively refined at later  
 273 hierarchical stages (where the data is subsampled to a finer grid). The deformation field and the objective

274 function are regularized with a Gaussian kernel at each level of the hierarchy. The ANTS algorithm is  
 275 freely available <http://www.picsl.upenn.edu/ANTS/>. We used an implementation of the ANTS algorithm  
 276 compatible with the MINC data format, mincANTS <https://github.com/vfonov/mincANTS>.

277 We used the following command line when running ANTS:

```
278 mincANTS 3 -m PR[target_file.mnc,source_file.mnc,1,4]
279   --number-of-affine-iterations 10000x10000x10000x10000x10000
280   --affine-gradient-descent-option 0.5x0.95x1.e-4x1.e-4
281   --use-Histogram-Matching --MI-option 32x16000
282   -r Gauss[3,0] -t SyN[0.5] -i 100x100x100x20
283   -o transformation.xfm
284
```

285 These settings were adapted from the "reasonable starting point" given in the ANTS manual <sup>2</sup>.

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

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

291 **Cross-correlation Weighted Majority Vote** An optimal combination of targets from the template li-  
 292 brary has previously been shown to improve segmentation accuracy with respect to manual segmenta-  
 293 tions (Aljabar et al., 2009; Collins and Pruessner, 2010). In this method, each template library image  
 294 is ranked in similarity to each unlabelled image by the normalized cross-correlation (CC) of image  
 295 intensities after linear registration, over a region of interest (ROI) generously encompassing the hip-  
 296 pocampus. Only the top ranked template library image labels are used in a voxel-wise majority vote.  
 297 The ROI is heuristically defined as the extent of all atlas labels after linear registration to the template,  
 298 dilated by three voxels (Chakravarty et al., 2013). The number of top ranked template library image  
 299 labels is a configurable parameter and displayed as the size of the template library in the rest of the  
 300 paper.

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

303 **Normalised Mutual Information Weighted Majority Vote** This method is similar to cross-correlation  
 304 weighted voting except that image similarity is calculated by the normalised mutual information score  
 305 over the region of interest (Studholme et al., 2001). The `itk_similarity` utility from the EZMinc  
 306 toolkit<sup>3</sup> is used to calculate the normalised mutual information measure between two images.

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

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

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

<sup>3</sup><https://github.com/vfonov/EZminc>

310 where  $A$  and  $B$  are the regions being compared, and the cardinality is the volume measured in voxels. The  
311 labels produced by MAGeT-Brain and multi-atlas segmentation are compared to the manual labels using  
312 the Dice similarity coefficient, and the recorded value for each subject at each parameter setting explored in  
313 this experiment is the average over ten validation rounds.

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

320 **3.1.2 Experiment 1: Results**

321 We find that for MAGeT-Brain segmentations, similarity score increases as atlas and template library size  
322 is increased, although with diminishing returns and an eventual trend towards a plateau (Figure 2a). For  
323 instance, with 9 atlases and using ANTS for registration and majority vote fusion, the mean DSC scores for  
324 1, 5, 9 and 17 templates are 0.845, 0.865, 0.867, 0.869, respectively. A maximum similarity score of 0.869 is  
325 found when using 9 atlases, 19 templates, ANTS registration, and majority vote label fusion.

326 The ANTS registration method consistently outperforms ANIMAL registration over all variable settings  
327 we tested (mean increase in DSC is 0.079). Pearson correlations of MAGeT-Brain DSC scores when using  
328 weighted voting and when using non-weighted majority vote label fusion (with ANTS registration) for all  
329 combinations of atlases and templates are  $r > 0.899$ ,  $p < 0.001$ , with a mean difference in DSC score of 0.002.  
330 This result suggests that using a weighted voting strategy does not significantly improve MAGeT-Brain  
331 segmentation agreement, contrary to the findings of Aljabar et al. (2009) for basic multi-atlas segmentation.  
332 Thus, in the remainder of our experiments only results using the ANTS registration algorithm and majority  
333 vote fusion will be shown.

334 With at least five templates, MAGeT-Brain consistently shows a higher DSC score than multi-atlas  
335 segmentation wth the same number of atlases:  $r = 0.94$ ,  $p < 0.001$ , mean DSC increase = 0.008 (Figure 2b).  
336 The magnitude of DSC increase grows with template library size but shows diminishing returns with larger  
337 atlas libraries. Peak increase (+0.025 DSC) is found with a single atlas and template library of 19 images.

338 In addition to a mean increase in similarity score over multi-atlas-based segmentation, MAGeT-Brain also  
339 shows more consistency in similarity scores across all subjects and validation folds (Figure 2c). A template  
340 library of at least 13 images is sufficient to show significant ( $p < 0.05$ ) decrease in variance for all sizes of  
341 atlas library tested (1-9 images).

342 We find similar behaviour with respect to optimal parameter settings and increased consistency of  
343 MAGeT-Brain segmentations in the replication of this experiment (Experiment 5, Supplementary Mate-  
344 rials) where a different hippocampal definition is used (SNT labels available with the ADNI datasets). This  
345 strongly suggests that these results are independent of the segmentation protocol used and are, instead,  
346 features of the MAGeT-Brain algorithm.

347 We have omitted results obtained when using an even number of atlases or templates since with these  
348 configurations we found significantly decreased performance. We believe this results from an inherent bias  
349 in the majority vote fusion method used (see Discussion).

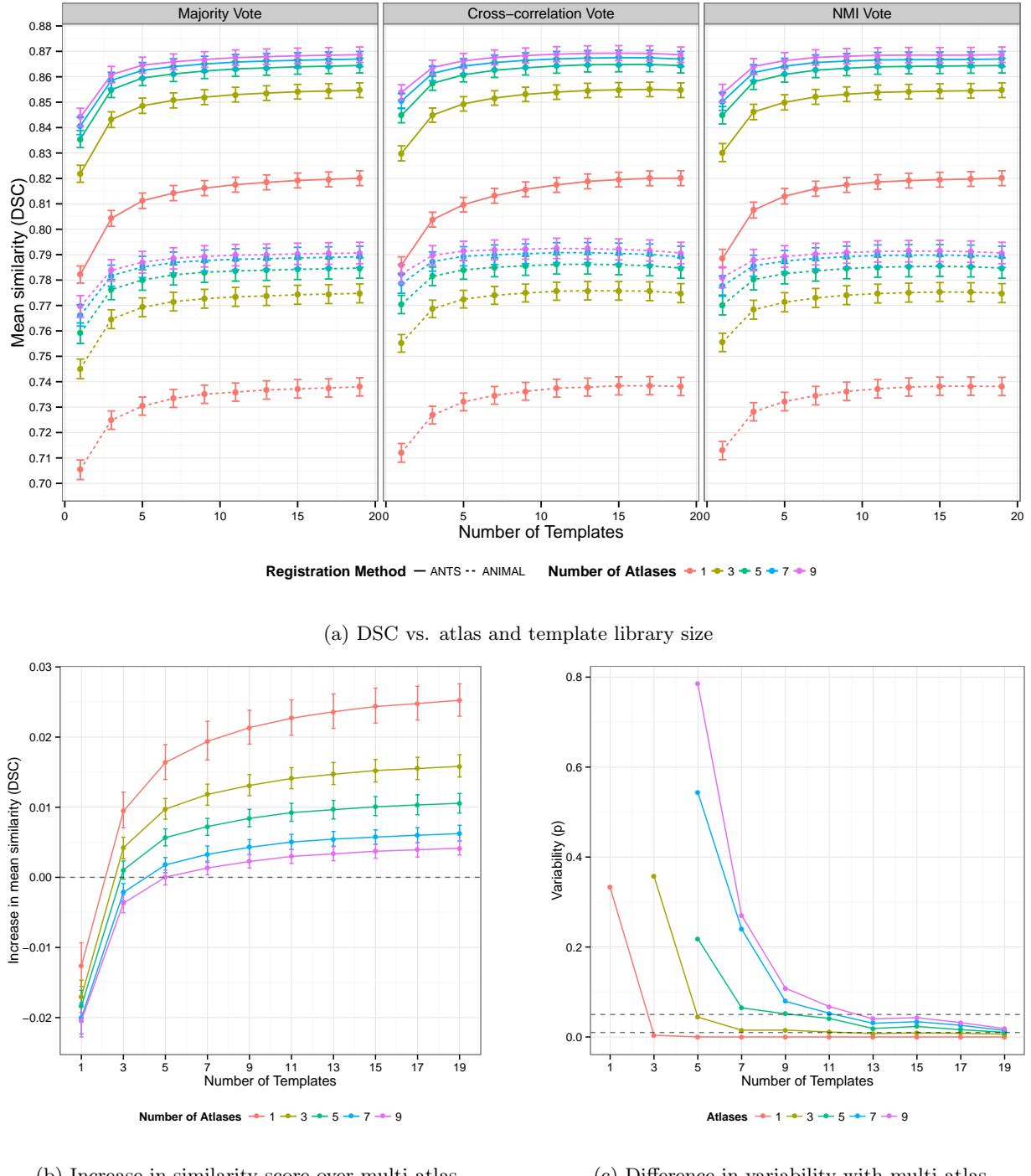


Figure 2: **Whole hippocampus segmentation cross-validation on ADNI subjects with Pruessner-protocol manual segmentations.** (2a) Average DSC score of MAGeT-Brain with manual segmentations for 60 ADNI subjects taken over 10 folds of cross-validation at each parameter setting. Error bars indicate standard error. (2b) Increase in DSC of MAGeT-Brain over multi-atlas segmentations. (2c) shows the significance of t-tests comparing the variability in DSC scores of MAGeT-Brain and multi-atlas across validation folds. Only points where MAGeT-Brain mean variability is lower than multi-atlas are shown. Dashed lines indicate p-values of 0.05 and 0.01.

350 **3.2 Experiment 2: Whole Hippocampus Segmentation Cross-Validation — First**  
351 **Episode of Psychosis**

352 To validate that the MAGeT-Brain works effectively in the context of other neurological disorders, in this  
353 experiment we replicate the cross-validation done in Experiment 1 with a dataset of patients having had a  
354 single episode of psychosis. We also compare MAGeT-Brain segmentations with those of two well-known  
355 automated segmentation methods, FSL FIRST and FreeSurfer.

356 **3.2.1 Experiment 2: Materials and Methods**

357 **First Episode Psychosis (FEP) Dataset** All patients were recruited and treated through the Prevention  
358 and Early Intervention Program for Psychoses (PEPP-Montreal), a specialized early intervention service at  
359 the Douglas Mental Health University Institute in Montreal, Canada. People aged 14 to 35 years from the  
360 local catchment area suffering from either affective or non-affective psychosis who had not taken antipsychotic  
361 medication for more than one month with an IQ above 70 were consecutively admitted as either in- or out-  
362 patients. Of those treated at PEPP, only patients aged 18 to 30 years with no previous history of neurological  
363 disease or head trauma causing loss of consciousness were eligible for the neuroimaging study; only those  
364 suffering from schizophrenia spectrum disorders were considered for this analysis. For complete program  
365 details see Malla et al. (2003).

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

371 Expert whole hippocampal manual segmentation of each subject is produced following a validated seg-  
372 mentation protocol (Pruessner et al., 2000).

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

378 **Experiment details** The same overall design as Experiment 1 is followed in this experiment: a Monte  
379 Carlo cross-validation (MCCV) is conducted using the pool of 81 first episode psychosis subject brain images  
380 and corresponding Pruessner-protocol manual segmentations. Five rounds of validation are conducted for  
381 each subject, and each atlas and template library size combination (1-9 atlases, 1-19 templates). In each  
382 round, images and their manual labels are randomly selected from the pool, and the remaining images are  
383 segmented using MAGeT-Brain with a random subset of the unlabelled images also serving as template  
384 images. Majority vote fusion, and the ANTS registration algorithm are used, as these have shown to behave  
385 favourably in previous experiments.

386 In addition to the MCCV, we segment the entire first episode psychosis dataset using MAGeT-Brain using  
387 two different atlases, as well as two popular automated segmentation packages, FSL FIRST and FreeSurfer.  
388 Specifically, MAGeT-Brain is run once with the five Winterburn atlas images and labels as atlases and a

Table 3: **First Episode Psychosis Subject Demographics.** ambi - ambidextrous. SES - Socioeconomic Status score. FSIQ - Full Scale IQ.

	N	FEP N = 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.

389 randomly selected subset of 19 target images as templates. MAGeT-Brain is run a second time using the  
 390 same template images, but we using five additional first episode psychosis subjects and corresponding manual  
 391 segmentations (not included above) as atlases. FSL FIRST and FreeSurfer are run with the default settings:  
 392 FSL FIRST `run_first_all` script was used according to the FIRST user guide <sup>4</sup>, and FreeSurfer was run  
 393 with the command `recon-all -all`.

394 **Evaluation method** Manual and automated segmentations are directly compared using Dice's similarity  
 395 coefficient (DSC). In the MCCV, the per-subject DSC value is computed as the average value over the five  
 396 rounds of validation for a given atlas and template library size. The reported average DSC value per given  
 397 atlas and template library size is the average DSC value over all subjects segmented.

398 The Pruessner segmentation protocol differs slightly from the Winterburn protocol, and those used by  
 399 FreeSurfer and FSL FIRST, in the inclusion of neuroanatomical features and the manner they are delineated (see Winterburn et al. (2013), and Table 9 in the Discussion below). This variation in protocol poses  
 400 a problem if an overlap measure is used for evaluation: since different protocols will necessarily produce  
 401 segmentations that do not perfectly overlap, the degree of overlap cannot be solely used to compare segmen-  
 402 tation methods using different protocols. In place of an overlap metric, we assess the degree of (Pearson)  
 403 correlation in average bilateral hippocampal volume produced by each method. Additionally, we evaluate the  
 404 volume-related fixed and proportional biases in all segmentation methods using Bland-Altman plots (Bland  
 405 and Altman, 1986).

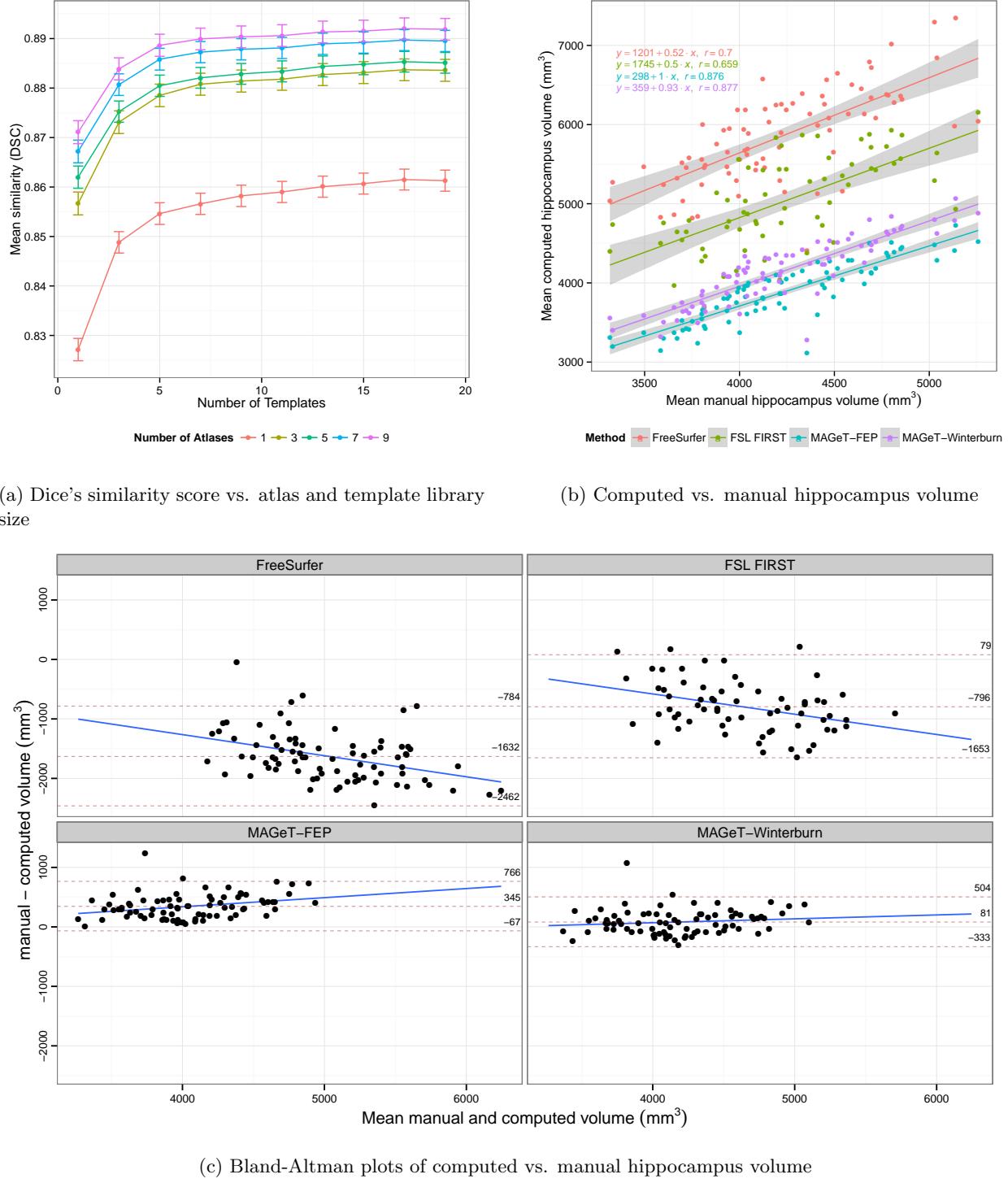
### 407 3.2.2 Experiment 2: Results

408 As in Experiment 1, we find that similarity score increases with a greater number of atlases or templates  
 409 but quickly plateaus (Figure 3a). A maximum similarity score of 0.892 is found when using 9 atlases, 19  
 410 templates, ANTS registration, and majority vote label fusion.

411 We found a close relationship in average hippocampal volume between the manual label volumes and  
 412 MAGeT-Brain when using the Winterburn atlases, or manually segmented FEP subjects as atlases (Figure  
 413 3b). Both sets of volumes are correlated with Pearson  $r > 0.88$ . FreeSurfer and FSL FIRST volumes are  
 414 both correlated with manual volumes at Pearson  $r > 0.7$ .

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

415 As Bland and Altman (1986) noted, high correlation amongst measures of the same quantity does not  
416 necessarily imply agreement (as correlation can be driven by a large range in true values, for instance).  
417 Figure 3c shows Bland-Altman plots illustrating the level of agreement of each method with manual vol-  
418 umes. All methods show an obvious proportional bias: FreeSurfer and FSL FIRST markedly underestimate  
419 smaller hippocampi and over-estimate large hippocampi (the limits of agreement are between  $-2482mm^3$   
420 and  $-784mm^3$ , and between  $-1653mm^3$  and  $79mm^3$ , respectively), whereas both MAGeT-Brain methods  
421 show a much less exaggerated, but conservative bias (limits of agreement between  $-67mm^3$  and  $766mm^3$   
422 when using FEP atlases, and between  $-333mm^3$  and  $504mm^3$  when using Winterburn atlases). On average,  
423 FreeSurfer and FSL FIRST overestimate hippocampal volume by about  $1600mm^3$  and  $800mm^3$ , respec-  
424 tively. In contrast, on average MAGeT-Brain underestimates volumes by about  $300mm^3$  when using FEP  
425 atlases and by about  $80mm^3$  when using Winterburn atlases (compared to the Pruessner-protocol manual  
426 segmentations).



**Figure 3: First Episode Patient dataset validation.** All manual segmentation of the 81 subjects is done with the Pruessner-protocol. MAGeT-Brain uses ANTS registration and majority vote label fusion. (3a) shows mean DSC score of MAGeT-Brain segmentations, as atlas and template library size is varied over a 5-fold validation. Error bars indicate standard error. (3b) shows segmentation volumes from FSL FIRST, FreeSurfer, MAGeT-Brain using the five Winterburn atlases (MAGeT-Winterburn), and MAGeT-Brain using five manually segmented FEP subjects as atlases (MAGeT-FEP). Linear fit lines are shown, with the shaded region showing standard error. (3c) shows the agreement between computed and manually volumes. 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 manual volume, and vice versa.

427 **3.3 Experiment 3: Whole Hippocampus Segmentation Comparison — ADNI1**  
428 **Complete 1Yr**

429 To validate MAGeT-Brain segmentation quality with respect to other established automated hippocampal  
430 segmentation methods, we apply MAGeT-Brain to a large dataset from the ADNI project. The resulting seg-  
431 mentations are compared to those produced by FreeSurfer, FSL FIRST, MAPER, as well as semi-automated  
432 whole hippocampal segmentations (SNT) provided by ADNI.

433 **3.3.1 Experiment 3: Materials and Methods**

434 **ADNI1:Complete 1Yr 1.5T dataset** The *ADNI1:Complete 1Yr 1.5T* standardized dataset contains  
435 1919 images in total. SNT, MAPER, and FreeSurfer hippocampal volumes for a subset of images were  
436 provided by ADNI, along with quality control data for each FreeSurfer segmentation (guidelines described  
437 in (Hartig et al., 2010)). See Section 3.1.1 for study details, inclusion criteria and imaging characteristics.

438 For a subset of the ADNI images, semi-automated segmentations of the left and right whole hippocampi  
439 generated using the SNT tool from Medtronic Surgical Navigation Technologies, Louisville, CO (see Sup-  
440 plementary Materials for detailed discussion of the segmentation process) are made available (Hsu et al.,  
441 2002). These labels are used as the reference labels in several other studies of (semi-)automated segmenta-  
442 tion methods (see Discussion). In addition, ADNI also distributes hippocampal segmentations and volumes  
443 determined using MAPER (Heckemann et al., 2011), a multi-atlas segmentation tool, and the FreeSurfer  
444 tool (including quality control data, with guidelines described in Hartig et al. (2010)).

445 **Experiment details** MAGeT-Brain was configured with an atlas library composed of the five Winterburn  
446 atlas images (Experiment 2, section 3.2) and segmentations. A template library of 19 images were randomly  
447 selected from the target dataset of ADNI subjects, and ANTS registration and majority vote label fusion  
448 were used as these were found to perform favourably in earlier experiments.

449 FSL FIRST segmentation was performed using the `run_first_all` script according to the FIRST user  
450 guide <sup>5</sup>. All images in the ADNI1:Complete 1Yr 1.5T dataset were segmented by both methods.

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

456 **Evaluation method** As in previous experiments, the Winterburn hippocampal segmentation protocol  
457 differs in the delineated neuroanatomical features (Winterburn et al. (2013), and Table 9, Discussion) and  
458 so we assess MAGeT-Brain by the degree of (Pearson) correlation of average hippocampal volume across  
459 subjects. We also computed the correlation in hippocampal volume between existing, established automated  
460 segmentation methods – FSL FIRST, FreeSurfer, and MAPER, and SNT semi-automated segmentations.  
461 Additionally, we evaluate the volume-related fixed and proportional biases in all segmentation methods using  
462 Bland-Altman plots (Bland and Altman, 1986).

---

<sup>5</sup><http://fsl.fmrib.ox.ac.uk/fsl/fslwiki/FIRST/UserGuide>

Table 4: **ADNI1 1.5T Complete 1Yr dataset demographics.** CN - Cognitively Normal. LMCI - Late-onset Mild Cognitive Impairment. AD - Alzheimer’s Disease. CDR-SB - Clinical Dementia Rating-Sum of Boxes. ADAS - Alzheimer’s Disease Assessment Scale. MMSE - Mini-Mental State Examination.

	N	CN N = 584	LMCI N = 931	AD N = 404	Combined N = 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.

Table 5: Number of segmented images and quality control failures of ADNI1:Complete 1Yr 1.5T dataset by method.label

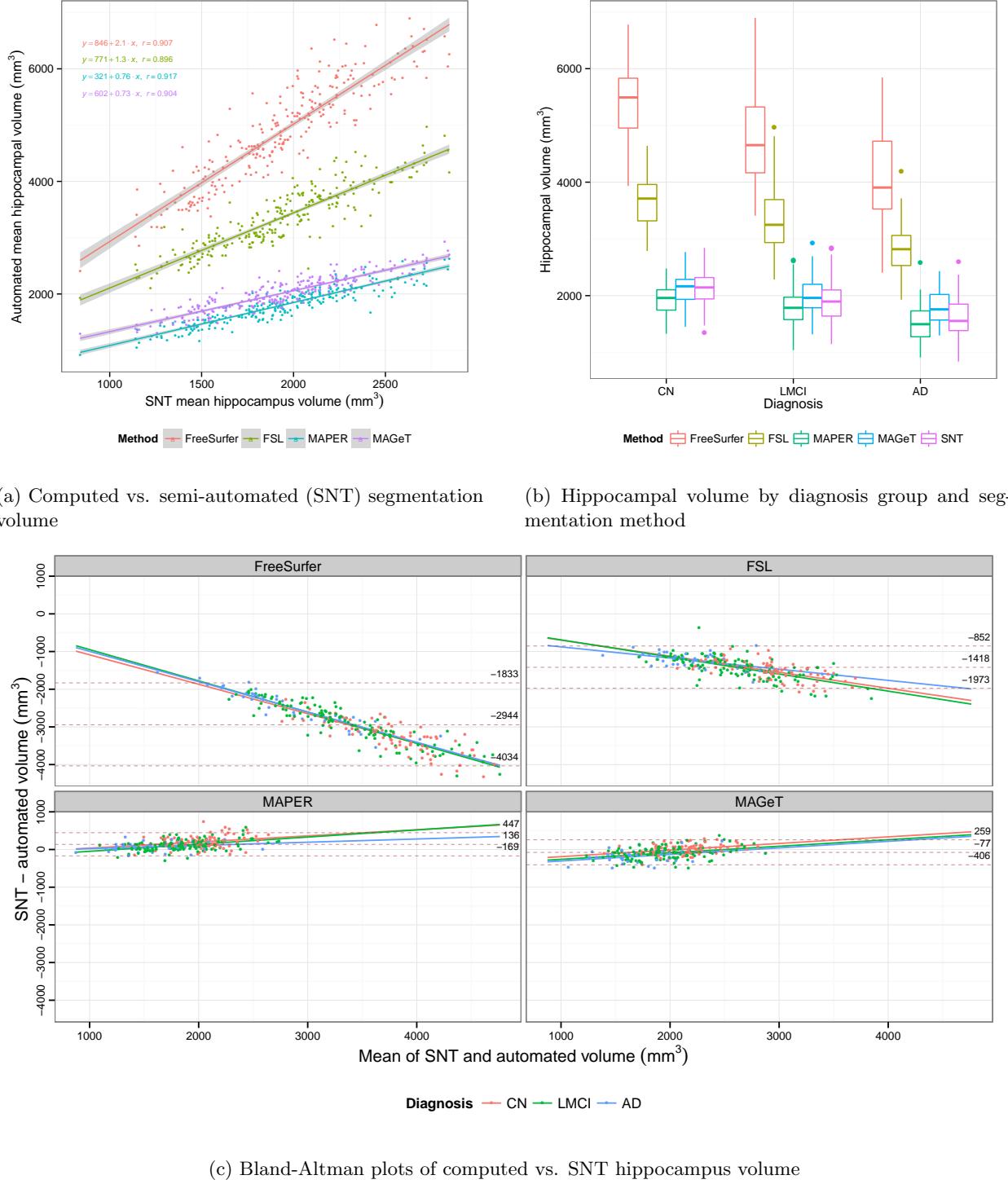
X	SNT	MAGeT	MAPER	FSL	FS
Images	368	368	368	368	368
Failures	n/a	30	n/a	20	88

### 463 3.3.2 Experiment 3: Results

464 We found a close relationship in total bilateral hippocampal volume between all methods and the SNT semi-  
 465 automated label volumes (Figure 4a). Volumes are well correlated ( $r > 0.78$ )c for all methods, and across  
 466 disease categories. Within disease categories (Figure 4b), MAGeT-Brain is consistently well correlated to  
 467 SNT volumes ( $r > 0.85$ ), but appears to slightly over-estimate the volume of the AD hippocampus compared  
 468 to the SNT segmentations.

469 Bland-Altman plots illustrate the level of agreement of each method with SNT segmentation hippocampal  
 470 volumes (Figure 4c). All methods show an obvious proportional bias: FreeSurfer and FSL FIRST markedly  
 471 under-estimate smaller hippocampi and over-estimate large hippocampi, whereas MAPER and MAGeT-  
 472 Brain show a reverse, conservative bias (Figure 4c). Additionally, all methods show a fixed volume bias,  
 473 with FreeSurfer and FSL FIRST most dramatically over-estimating hippocampal volume by  $2600mm^3$  and  
 474  $2800mm^3$  on average, respectively, and MAPER and MAGeT-Brain within  $250mm^3$  on average.

475 Figure 5 shows a qualitative comparison of MAGeT-Brain and SNT hippocampal segmentations for 10  
 476 randomly selected subjects in each disease category, and illustrates some of the common errors found during  
 477 visual inspection. Mostly frequently, we found MAGeT-Brain improperly includes the vestigial hippocampal  
 478 sulcus and, although not anatomically incorrect, MAGeT-Brain under-estimates the hippocampal body in  
 479 comparison to the SNT segmentation.



**Figure 4: ADNI1:Complete 1Yr 1.5T dataset segmentation.** (4a) Subject mean hippocampal volume as measured by each of the four automated methods (FreeSurfer (FS), FSL FIRST, MAPER, MAGeT-Brain) versus the semi-automated SNT segmentation volumes. Linear fit lines and Pearson correlations with SNT labels are shown for each method. (4b) Mean hippocampal volume by method and disease category. AD = Alzheimer’s disease, LMCI = late-onset mild cognitive impairment, and CN = cognitively normal. (4c) Bland-Altman plots shows the agreement between computed and SNT hippocampus volume. 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.

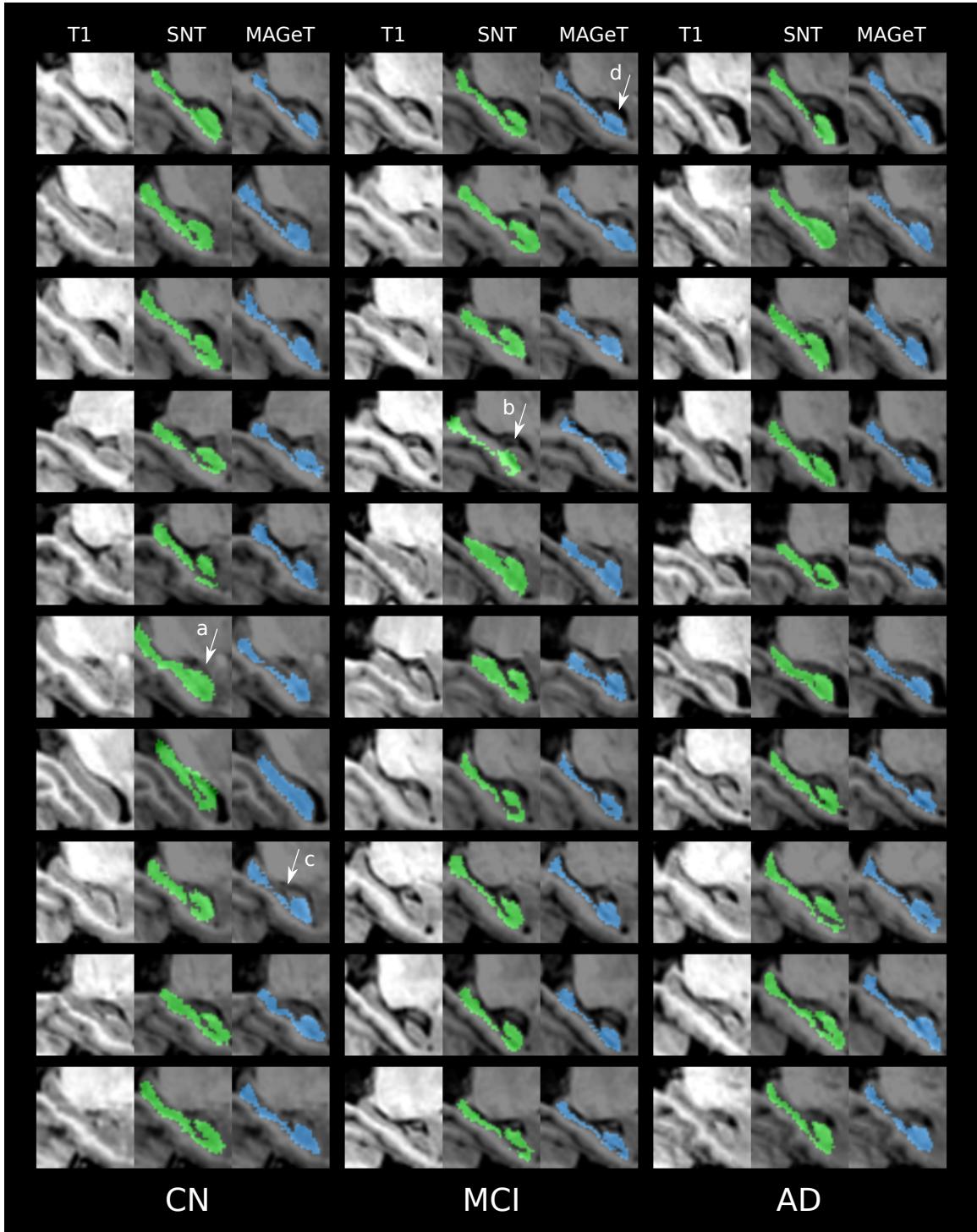


Figure 5: SNT and MAGeT-Brain segmentations for 30 ADNI subjects — 10 subjects randomly selected from each disease category in the subject pool used in Experiment 1 (Section 3.1). Sagittal slices are shown for each unlabelled T1-weighted anatomical image. SNT labels appear in green, and MAGeT-Brain labels 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.

480 **3.4 Experiment 4: Hippocampal Subfield Segmentation Cross-Validation**

481 The previous experiment assesses MAGeT-Brain performance on whole hippocampus segmentation. In this  
482 experiment, we conduct a proof-of-concept evaluation of MAGeT-Brain hippocampal subfield segmentation  
483 of standard 3T T1-weighted images at  $0.9\text{mm}$ -isotropic voxels. We use a modified leave-one-out cross-  
484 validation (LOOCV) design.

485 **3.4.1 Experiment 4: Materials and Methods**

486 **Healthy Control Dataset** T1 MR images of 14 subjects were acquired as a part of an ongoing study at  
487 the Centre for Addiction and Mental Health (Table 6). Subjects were known to be free of neuropsychiatric  
488 disorders and gave informed consent. These images were acquired on a 3T GE Discovery MR 750 system  
489 (General Electric, Milwaukee, WI) using an 8-channel head coil with the enhanced fast gradient recalled  
490 echo 3-dimensional acquisition protocol, FGRE-BRAVO, with the following parameters:  $TE/TR/TI =$   
491  $3.0ms/6.7ms/650ms$ , flip angle= $8^\circ$ ,  $FOV = 15.3cm$ , slice thickness=  $0.9mm$ , 170 in-plane steps for an  
492 approximate  $0.9mm$ -isotropic voxel resolution.

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

497 In this experiment, the Winterburn atlases (Experiment 2, section 3.2) are resampled to  $0.9mm$ -isotropic  
498 voxel resolution to simulate standard 3T T1-weighted resolution images. Image subsampling is performed  
499 using trilinear subsampling techniques. In each round of LOOCV, a single atlas image is selected and treated  
500 as a target image to be segmented by MAGeT-Brain. So as to have an odd-sized atlas library, atlas image  
501 is segmented once using each possible triple of atlas images, and corresponding manual segmentations, from  
502 the remaining four unselected atlases. Thus, for each of the five atlases, a total of  $\binom{3}{4} = 4$  segmentations are  
503 evaluated, resulting in a combined total of  $5 \times 4 = 20$  segmentations evaluated overall. We chose an atlas  
504 library with an odd number of images so as to ensure unbiased label fusion when using majority voting (see  
505 Discussion).

506 The template library used has a total of 19 images composed of all five resampled atlas images plus the  
507 additional 14 images from the healthy control dataset. The ANTS registration algorithm was used for image  
508 registration, and majority voting was used for label fusion, as these methods proved most favourable in the  
509 previous whole hippocampal validation experiments.

510 **Evaluation method** Evaluating the agreement of automated hippocampal subfield segmentations with  
511 manual segmentations for T1 images at  $0.9mm$ -isotropic voxels is inherently ill-defined since there are no  
512 manual protocols for segmentation at this resolution. Instead, we must evaluate how well the lower-resolution  
513 MAGeT-Brain hippocampal subfield segmentations correspond in form to the segmentation protocol used  
514 in the high-resolution images. By directly resampling the Winterburn atlas segmentations to  $0.9mm^3$  voxels  
515 (using standard nearest-neighbour image resampling techniques) we obtain a subsampled version of the labels  
516 which preserve the original segmentation protocol within the limits of error from rounding and interpolation.  
517 Therefore, using the resampled Winterburn segmentations as definitive for the  $0.9mm^3$  resolution we evaluate  
518 agreement of MAGeT-Brain segmentations using DSC overlap scores and evaluate consistency across the  
519 range of hippocampal sizes using Bland-Altman plots of subfield volumes.

Table 6: **Demographics for the hippocampal subfield cross-validation healthy control subject sample used in the template library (excluding the Winterburn atlas subjects).** Education is shown in years.

	N	Control <i>N</i> = 14
Age	14	34.5 53.0 62.0
Sex : Male	14	43% (6)
Education : 12	13	15% (2)
13		8% (1)
14		23% (3)
16		15% (2)
18		38% (5)
Handedness : R	14	93% (13)

*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.

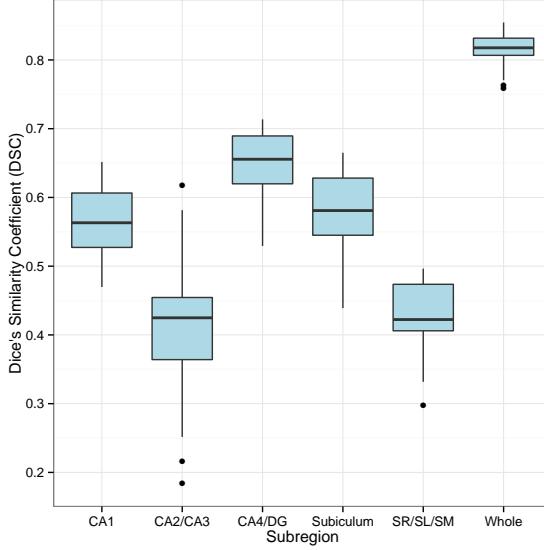
520     Additionally, by shifting the original manual 0.3mm-isotropic voxel segmentations by one voxel in the x, y,  
 521 and z direction and then resampling it to 0.9mm-isotropic voxels we obtain a simulated manual segmentation  
 522 having a small amount of error. We can compare the DSC overlap score of the shifted labels (relative to the  
 523 directly resampled labels) with the DSC score of the MAGeT-Brain generated labels in order to evaluate  
 524 their relevance.

### 525 3.4.2 Experiment 4: Results

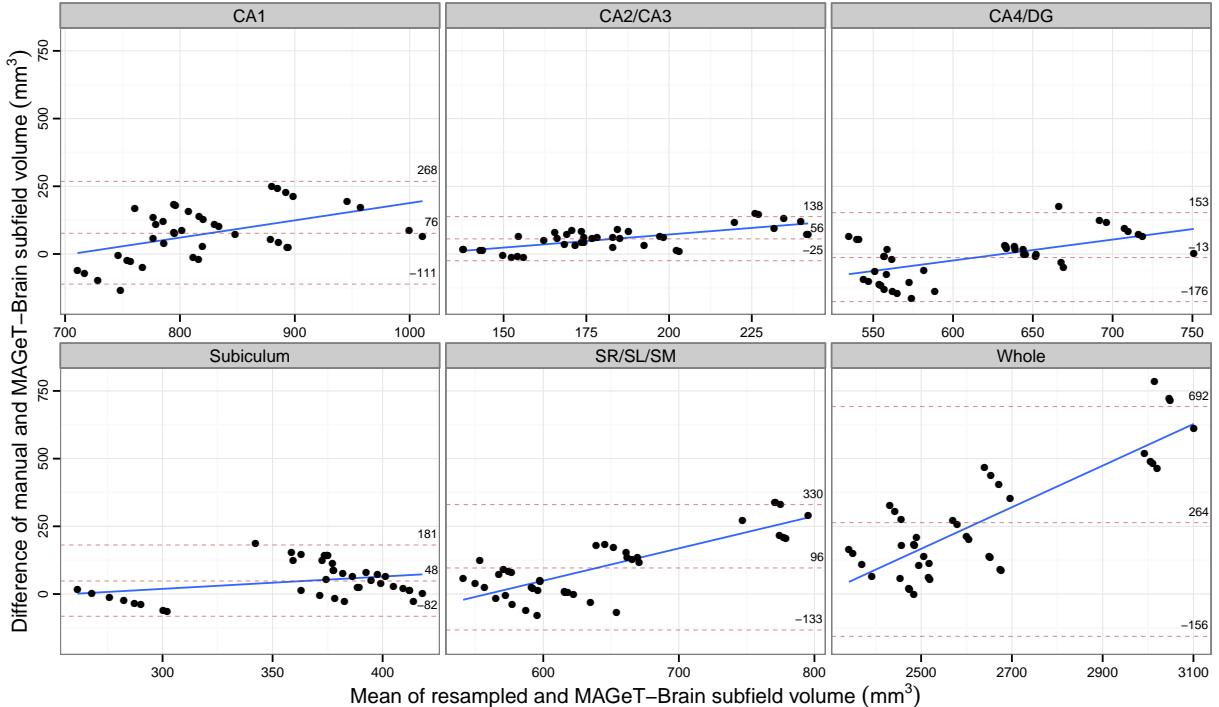
526     Figure 6a shows the overlap similarity scores between the MAGeT-Brain segmentations and the resampled  
 527 Winterburn atlases for each hippocampal subfield across all subjects and folds of the validation. Mean  
 528 and standard deviation DSC scores of the subfields are shown in Table 7, along with DSC scores for the  
 529 resampled atlas segmentations when perturbed slightly and compared to the originals. We find that the  
 530 CA4/DG subfield shows the highest mean DSC score of  $0.647 \pm 0.051$ , followed by the Subiculum and CA1  
 531 subfields having scores of  $0.563 \pm 0.046$  and  $0.58 \pm 0.057$ , respectively. Both the CA4/DG and molecular  
 532 regions score below 0.5. These scores may seem low but not when taken in context and compared to existing  
 533 (semi-)automated methods (see Discussion). The whole hippocampus is segmented with a mean DSC score  
 534 of  $0.816 \pm 0.023$ .

535     Figure 6b contains Bland-Altman plots comparing MAGeT-Brain volumes with manual volumes across  
 536 all validation folds. MAGeT-Brain displays a conservative proportional bias — small hippocampi are over-  
 537 estimated in volume, and larger hippocampi are underestimated (a mean maximum difference of approx-  
 538 imately  $200\text{mm}^3$  across all subfields). MAGeT-Brain display a slight conservative fixed bias, tending to  
 539 underestimate all subfields except CA4/DG (mean underestimation:  $\text{CA1} = 76\text{mm}^3$ ,  $\text{CA2/3} = 56\text{mm}^3$ ,  
 540  $\text{CA4/DG} = -16\text{mm}^3$ ,  $\text{Subiculum} = 48\text{mm}^3$ ,  $\text{SR/SL/SM} = 96\text{mm}^3$ ).

541     Figure 7 shows slices subfield segmentations for a single subject for qualitative inspection.

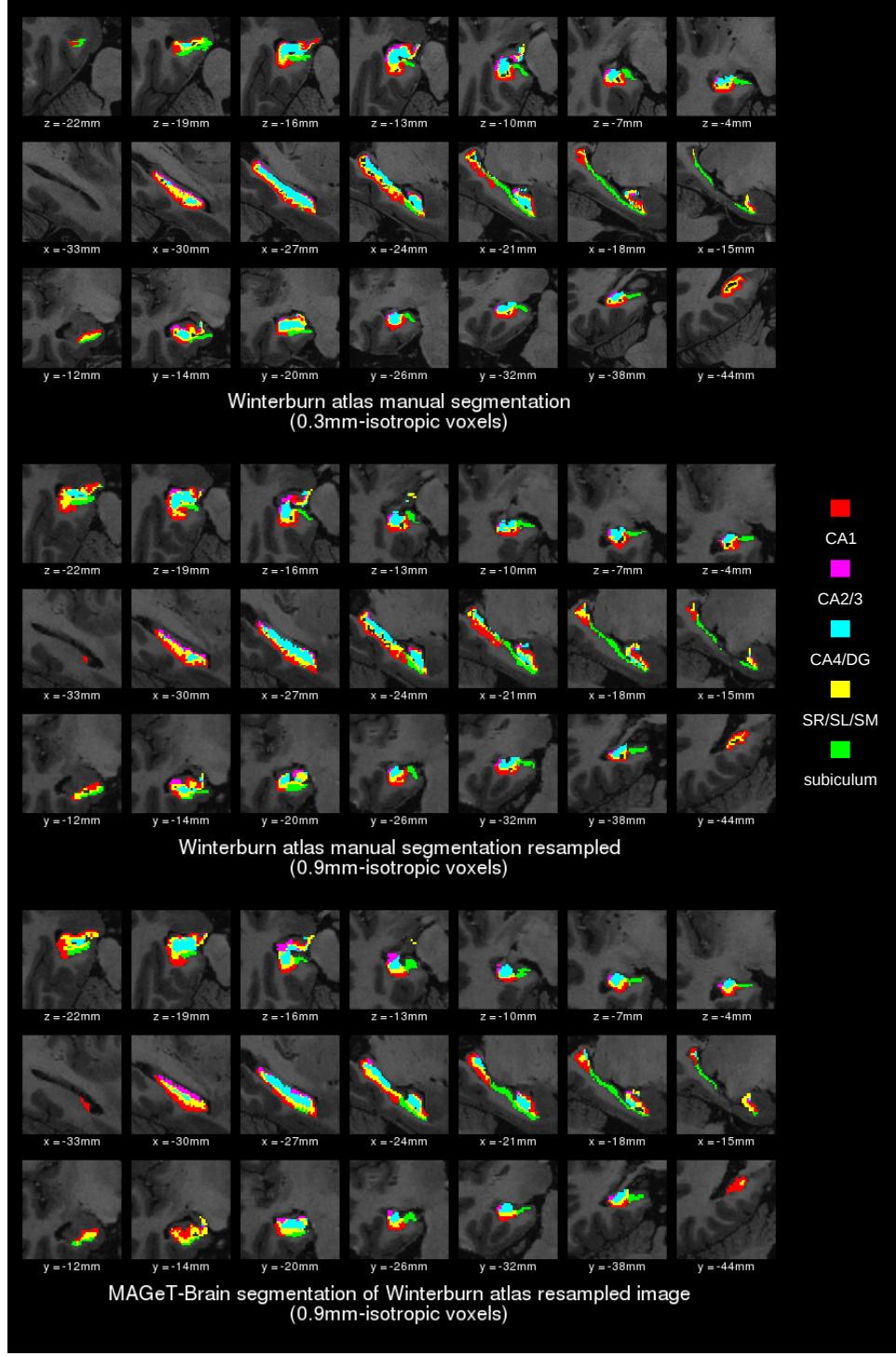


(a) DSC score by subfield



(b) Bland-Altman plots of computed vs. manual subfield volumes

**Figure 6: Hippocampal subfield cross-validation.** (6a) Similarity of MAGEt-Brain segmentation of subfields and the resampled Winterburn atlas segmentations at  $0.9\text{mm}^3$  voxel resolution, over all validation folds. Overlap score for each hemisphere is measured separately. (6b) shows the agreement, by subfield, of computed and manual volumes across all validation folds. The overall mean difference in volume, and limits of agreement ( $\pm 1.96\text{SD}$ ) are shown by dashed horizontal lines. Linear fit lines are shown. Note, points below the mean difference indicate overestimation of the volume with respect to the resampled volume, and vice versa.



**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 subsampled from 0.3mm- to 0.9mm-isotropic voxels; in the lower left corner is the MAGeT-Brain segmentation of the resampled Winterburn atlas image from a single fold of the cross-validation. 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.

Table 7: Overlap similarity results for the each of the subfields of the hippocampus. Simulated overlap similarity results are also given for manual labels that were translated by one voxel (i.e.:  $0.3\text{mm}$ ) in all directions and then resampled. Values are given as mean Dice’s Similarity Coefficient (DSC)  $\pm$  standard deviation.

Subfield	MAGeT	$0.9\text{mm}$ translation
CA1	$0.56 \pm 0.05$	$0.27 \pm 0.03$
CA2/CA3	$0.41 \pm 0.10$	$0.12 \pm 0.05$
CA4/DG	$0.65 \pm 0.05$	$0.42 \pm 0.05$
SR/SL/SM	$0.43 \pm 0.05$	$0.19 \pm 0.04$
Subiculum	$0.58 \pm 0.06$	$0.14 \pm 0.04$

## 542 4 Discussion

543 In this manuscript we have presented the implementation and validation of the MAGeT-Brain framework  
 544 – a methodology that requires very few input atlases in order to provide accurate and reliable segmentations  
 545 with respect to manual segmentations. Both Experiment 1 (Section 3.1) and Experiment 2 (Section  
 546 3.2) compare MAGeT-Brain to basic-multi-atlas segmentation by characterising the change in segmentation  
 547 quality with varying parameter settings (atlas and template library sizes, registration method, and label  
 548 fusion method) and differing age and neuropsychiatric populations. Together, these experiments allow us  
 549 to choose optimal MAGeT-Brain parameter settings for use in subsequent experiments. Experiment 3 (Sec-  
 550 tion 3.3) demonstrates that across 246 images from the ADNI1:Complete 1Yr 1.5T dataset, MAGeT-Brain  
 551 performs as well as, or better, than other established and popular methods, and has a much more conser-  
 552 vative proportional bias in segmentation volume. Finally, Experiment 4 (Section 3.4) is a proof-of-concept  
 553 validation demonstrating the reliability of MAGeT-Brain in producing subfield segmentations which match  
 554 the segmentation protocol of the input atlases despite contrast and resolution limitations in standard T1-  
 555 weighted image volumes. All of these experiments together demonstrate that MAGeT-Brain’s algorithmic  
 556 performance is not dependent on a single definition of the hippocampus but is effective with differing hip-  
 557 pocampal definitions (Winterburn et al., 2013; Pruessner et al., 2000; Hsu et al., 2002), across image types,  
 558 and subject populations.

559 The core claim the MAGeT-Brain method is based on – that a useful template library can be generated  
 560 from a small set of labelled atlas images – is validated in the cross-validation conducted in Experiment 1 (and  
 561 the replication in Experiment 2 and Experiment 5, Supplementary Materials). We find that both increasing  
 562 the number of atlases and the number of templates used improves MAGeT-Brain segmentation over and above  
 563 basic-multi-atlas segmentations using the same number of atlas images. That is, by taking the extra step  
 564 of generating a template library using target images, MAGeT-Brain is able to improve the overlap between  
 565 the automatically generated segmentations and manually generated “gold standard” segmentations. The  
 566 magnitude of this improvement is greatest with a small number of atlases, but even with larger atlas libraries  
 567 we have found that generating a template library reduces the variability in segmentation agreement (i.e.  
 568 MAGeT-Brain more consistently produces segmentations in greater agreement with manual segmentations  
 569 than does basic-multi-atlas method, over repeated randomized trials). These effects do not appear dependant  
 570 on the hippocampal segmentation protocol used.

571 Interestingly, previous work on multi-atlas segmentation methods (Aljabar et al., 2009; Collins and  
 572 Pruessner, 2010) has found that cross-correlation and normalized mutual information-based weighted la-  
 573 bel fusion improves segmentation reliability over simple majority vote label fusion, and yet we did not see

**Table 8: Summary of automated segmentation methods of the Hippocampus.** This table summarizes published Dice’s overlap measure between automated and manual segmentations of the hippocampus. Unless otherwise specified, validation datasets are composed equally of cases and control subjects, and use manual segmentation labels as ground truth in computing DSC scores. AD = Alzheimer’s Disease; MCI = Mild Cognitive Impairment; CN = Cognitively Normal (CN); FEP = First Episode of Psychosis; LOOCV = Leave-one-out cross-validation; MCCV = Monte Carlo cross-validation; SNT = Surgical Medtronic Navigation Technologies semi-automated labels. 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 mean (AD; MCI; CN)	Reference	Validation	Dataset (Truth)
MAGeT-Brain	9	0.841		10-fold MCCV on 69 subjects	ADNI (SNT)
Patch-based label fusion	16	0.861 (0.838; →; 0.883)	Coupe et al. (2011)	LOOCV	ADNI (SNT)
Multi-atlas	20	0.848 (→; 0.798; 0.898)	Wang et al. (2011)	10-fold MCCV on 20 of 139 subjects	ADNI (SNT)
ACM (Ada Boost-based)	21	0.862	Morra et al. (2008)	LOOCV	ADNI (SNT)
LEAP	30	0.848	Wolz et al. (2010)	Segmentation of 182 subjects	ADNI (SNT)
Multi-atlas	30	0.885	Lötjönen et al. (2010)	Segmentation of 60 subjects	ADNI (SNT)
Multi-atlas (MAPS)	55	0.890	Leung et al. (2010)	Segmentation of 30 subjects (10 AD, MCI, and CN)	ADNI (SNT)
MAGeT-Brain	9	0.869		10-fold MCCV on 60 subjects	ADNI (Pruessner)
MAGeT-Brain	9	0.892		5-fold MCCV on 81 subjects	FEP subjects
Neural nets	10	0.740		Segmentation of 5 subjects	controls
Probabilistic atlas	11	0.852		11 atlases used in 100 rounds of LOOCV on 20 elderly subjects	elderly controls
Probabilistic Atlas	16	0.860	Powell et al. (2008) van der Lijn et al. (2008) Chupin et al. (2009)	LOOCV	AD subjects
Anatomically-guided EM	17	0.812	Pohl et al. (2007)	LOOCV on 17 controls, segmentation of 33 mixed subjects	mixed diagnosis
Multi-atlas	30	0.820	Heckemann et al. (2006a)	LOOCV	controls
Multi-atlas	30	0.880	Gousias et al. (2008)	30 adult atlases used, segmentation of 33 2yr old subjects	2yr old controls
Multi-atlas	80	0.890	Collins and Pruessner (2010)	LOOCV	controls
Multi-atlas	55	0.860	Barnes et al. (2008)	LOOCV	controls and AD
Multi-atlas	275	0.835	Aljabar et al. (2009)	LOOCV	controls

---

574 a significant indication of this effect in the MAGeT-Brain segmentations. Selectively filtering out atlases  
575 with lower image similarity is believed to reduce sources of error from estimating deformations via non-  
576 linear registration, partial volume effects from nearest neighbour image resampling, and neuroanatomical  
577 mismatch between atlases and subjects. That MAGeT-Brain does not see the same boost in performance  
578 from weighted voting may suggest that the neuroanatomical variability of a template library constructed  
579 from study subjects more closely matches any particular subject and thereby leaving less error to filter. From  
580 our previous work on the MAGeT-Brain algorithm we have shown that the reduction in error is not simply  
581 a smoothing or averaging effect (Chakravarty et al., 2013).

582 Although, the goal of this manuscript was not to exhaustively test or validate multiple different voting  
583 strategies in the context of our segmentation algorithm, it is important to note that other strategies for  
584 voting are available. For example, other groups have used the STAPLE algorithm (Warfield et al., 2004)  
585 (or variants of the STAPLE algorithm (Robitaille and Duchesne, 2012)) which weights each segmentation  
586 based upon its estimated performance level with respect to the other available candidate segmentations.  
587 Further, the sensitivity and specificity parameters can also be tuned to potentially improve segmentation  
588 reliability. It is likely that using more sophisticated voting methods would have a positive effect on the overall  
589 segmentation performance, as demonstrated by the STAPLE algorithm. However, it is also important to note  
590 that even in the absence of a more sophisticated label fusion algorithm, MAGeT Brain performs reasonably  
591 well in comparison to other groups that have tested new segmentation algorithm with Alzheimer disease,  
592 mild cognitive impairment, and cognitively normal data from the ADNI database (Table 8. In addition, our  
593 validation in Experiment 2 (with the first episode psychosis subjects) yields DSC's that are amongst the  
594 highest reported. Thus, more work is required to determine the extent to which label fusion will improve  
595 the reliability of our algorithm.

596 More work is required to determine the source of the slight decrease in segmentation performance when  
597 the number of templates are set to an even number. Our initial concern was that this dip in performance  
598 was a by-product of the MAGeT-Brain algorithm itself. However, this pattern is also found in the results of  
599 the multi-atlas segmentations we used in our experiments. We believe that our majority voting methodology  
600 is biased towards labels with the lowest numeric values when breaking ties (by way of the implementation  
601 of the `mode` function used to determine majority), thus causing the slight bias observed when using an even  
602 number of templates. This is another area where the voting scheme could be used to improve performance.  
603 However, it is worth noting that this limitation was previously identified by Heckemann et al. (2006b) and,  
604 subsequently, other groups have not even considered the potential pitfalls of an even number of candidate  
605 labels (e.g. Leung et al. (2010)).

606 Despite MAGeT-Brain achieving segmentation results which are competitive with the rest of the field  
607 (Table 8), a concern may be raised over the modest improvement in segmentation agreement observed using  
608 MAGeT-Brain over multi-atlas, with the same number of atlases (Experiment 1). As we have shown in that  
609 same experiment, the benefit in using MAGeT-Brain is both an increase in the overlap agreement and also  
610 in the improved consistency of the labelling regardless of atlas or template choice. Reducing the variability  
611 in segmentation agreement is an important consideration that few have touched on previously. In addition,  
612 the Monte Carlo cross-validations that we present in Experiment 1 and Experiment 2 are amongst the most  
613 stringent performed in the multi-atlas segmentation literature. To the best of our knowledge, with the  
614 exception of (Wang et al., 2011), other groups do at most a single round of leave-one-out-validation (Table  
615 8). Thus, the thoroughness of our validation suggests that our results are reflective of a true average over  
616 the choice of parameter settings and are independent of atlas or template choice.

617 On that note, one author (JW), an expert manual rater (Winterburn et al., 2013), identified regular in-  
618 consistencies in the SNT segmentations: occurrences of over- and under-estimation, as well as misalignments  
619 of the entire segmentation volume (Figure 5). Although the SNT segmentations are used as benchmarks  
620 for validation in many other studies (Table 8), these segmentation inconsistencies present the possibility  
621 that a more accurate and consistent benchmark segmentation protocol ought to be used in order to truly  
622 understand the results of such validations. Indeed, our replication of the 10-fold cross-validation using SNT  
623 segmentations (Experiment 5, Supplementary Materials) shows noticeably poorer mean similarity scores for  
624 both MAGeT-Brain and multi-atlas.

625 Thus, in comparison to other methodologies in the field MAGeT-Brain performs favourably. Table 8 sur-  
626 veys some of the most recent reported DSC values reported on ADNI dataset, using SNT segmentations for  
627 the atlas library and as gold standards for evaluation. While it is difficult to compare segmentation results  
628 across studies, gold standards, evaluation metrics, and algorithms it is worth noting that the methods sum-  
629 marized require more atlases (between 16-55) than our MAGeT-Brain implementation with the Winterburn  
630 atlases (Winterburn et al., 2013).

631 There are some important differences between our method and these specific methods. Others have re-  
632 ported the difficulty with mis-registrations in candidate segmentation (i.e. segmentations generated that are  
633 then input in the voxel-voting procedure (Collins and Pruessner, 2010)). The work of Leung et al. (2010)  
634 tackles this problem by using an intensity threshold that is estimated heuristically at the time of segmen-  
635 tation (this work also reports some of the highest DSC scores for the segmentation of ADNI data). While  
636 this method is effective for the ADNI dataset (which is partially homogenized with respect to image acqui-  
637 sition and pre-processing), it is unclear if this type of heuristic is applicable to other datasets. In all cases,  
638 these methods require more atlases than our implementation with the Winterburn atlases. Lötjönen et al.  
639 (2010) produced segmentations which strongly agree with manual segmentations by way of post-processing  
640 corrections using classifications derived using an expectation maximization framework. In their initial work,  
641 Chupin et al. (2009) develop their probabilistic methodology using a cohort of 8 healthy controls and 15  
642 epilepsy patients, and then use this method to segment an ADNI sample, with a hierarchical experimenta-  
643 tion protocol. These methods suggest that some post-processing of the final segmentations would improve  
644 agreement of the segmentation. While that may be true, there is little consensus regarding how to achieve  
645 this.

646 To the best of our knowledge, no other groups have validated their work using multiple atlas segmentation  
647 protocols, different acquisitions, and disease populations in order to demonstrate the robustness of their  
648 technique. This is one of the clear strengths of this work. Furthermore, unlike some of the algorithms  
649 mentioned, our implementation does not require retuning for new populations or datasets as it inherently  
650 models the variability of the dataset through the template library. However it should be noted that the  
651 increased agreement that follows increasing the number of atlases and templates comes at an increased  
652 computational cost ( $O(\log(n))$ ), as previously mentioned in other work (Heckemann et al., 2006a).

653 Among the automated segmentation methods we compared in this paper (FreeSurfer, MAPER, FSL  
654 FIRST), we find extremely variable performance of all methods. With the exception of FSL FIRST all  
655 methods correlate well with the semi-automated SNT volumes provided in the ADNI database. However, the  
656 FreeSurfer and FSL FIRST hippocampal segmentations are on average about twice the volume of those from  
657 all other methods. Furthermore, when estimating the bias of FreeSurfer and FSL FIRST relative to the SNT  
658 hippocampal volumes we see that large hippocampi are over estimated while small hippocampi are under  
659 estimated. By comparison, MAGeT-Brain and MAPER are far more conservative in volume estimation,

---

Table 9: **Summary of labelled subfields of the Hippocampus from recent MRI segmentation protocols.**

Protocol	Labelled Subfields
Winterburn et al. (2013)	CA1, CA2/CA3, CA4/dentate gyrus, strata radiatum/lacunosum/moleculare, subiculum
Wisse et al. (2012)	CA1, CA2, CA3, CA4/dentate gyrus, subiculum, entorhinal cortex
Van Leemput et al. (2009)	CA1, CA2/CA3, CA4/dentate gyrus, presubiculum, subiculum, hippocampal fissure, fimbria, hippocampal tail, inferior lateral ventricle, choroid plexus
Yushkevich et al. (2009)	CA1, CA2/CA3, dentate gyrus (hilus), dentate gyrus (stratum moleculare), strata radiatum/lacunosum/moleculare/vestigial hippocampal sulcus
Mueller et al. (2007)	CA1, CA2, CA3/CA4 & dentate gyrus, Sibicum, entorhinal cortex

660 suggesting these methods may be better suited for estimating true-positives, especially in neurodegenerative  
 661 disease subjects featuring smaller overall hippocampi. However, in this analysis we have only compared  
 662 methods by total hippocampal volume, and so more work is needed to understand the full extent to which  
 663 these methods differ.

664 Finally, we have provided evidence that using the Winterburn high-resolution hippocampal subfield at-  
 665 lases (Winterburn et al., 2013) our algorithmic framework is appropriate for the segmentation of hippocampal  
 666 subfields in standard T1-weighted data. Subfield segmentation is a burgeoning topic in the literature although  
 667 very few automated methods are available for the segmentation of 3T data (Yushkevich et al., 2009, 2010;  
 668 Van Leemput et al., 2009). Table 10 compares segmentation agreement from some of these methods and  
 669 MAGeT-Brain. The overlap DSC scores for MAGeT-Brain subfields are notably lower but a direct compar-  
 670 ison of overlap values must be done cautiously. In the present work, our overlap scores are computed on  
 671 0.9mm-isotropic voxel resolution images, whereas Yushkevich et al. (2010) uses focal  $0.4 \times 0.5 \times 2.0\text{mm}$  voxel  
 672 resolution images, and Van Leemput et al. (2009) use supersampled 0.380mm-isotropic voxel resolution im-  
 673 ages. The larger voxel images we use necessarily entail a greater change in DSC for each incorrectly labelled  
 674 voxel. In addition, our automated segmentations are compared to manual segmentations resampled from  
 675 0.3mm-isotropic voxel labels; the resampling process inevitably introduces noise which may lower overlap  
 676 scores. Lastly, as our method is aimed specifically at situations when manually produced atlases are scarce,  
 677 in our cross validation we are forced to use three rather than all five of the Winterburn atlases (which, based  
 678 on our findings with whole hippocampal segmentation, would have resulted in improved overlap similarity).  
 679 Although having more atlases would be ideal in this context, these atlases are very time consuming to gen-  
 680 erate (Winterburn et al., 2013). Nevertheless, the advantage of evaluating MAGeT-Brain on standard 3T  
 681 T1-weighted resolution MR images with a publically available atlas library is that our results reflect typical  
 682 usage scenarios of researchers and clinicians.

683 Experiments 1, 2, and 5 have demonstrated that our algorithm flexibly accommodates different whole  
 684 hippocampus manual segmentation methodologies. We have not explicitly evaluated a subfield definition  
 685 other than the Winterburn protocol, and therefore it is possible that using an alternate subfield definition  
 686 could improve the reliability of our automated subfield definitions. For example, established definitions such  
 687 as those from Mueller et al. (2007) could be a prime candidate for further exploration. In addition, the  
 688 conservative nature of the Mueller definition (labelling of the 5 slices in the hippocampus body only) would  
 689 likely further aid in reliability measurement. However, there are two main logistical problems that we would  
 690 have to overcome prior to implementation. The first is that these definitions were developed for data that  
 691 is highly anisotropic ( $0.4\text{mm} \times 0.5\text{mm} \times 2\text{mm}$ ), and it is unclear how our algorithms would deal with such  
 692 atlases used as input. The second is that, since these atlases are not publicly available, we would be have

---

Table 10: A comparison of subfield segmentation overlap similarity with manual raters.

Subfield	MAGeT-Brain	Van Leemput et al. (2009)	Yushkevich et al. (2010)
CA1	0.563	0.62	0.875
CA2/3	0.412	0.74	$CA2 = 0.538, CA3 = 0.618$
CA4/DG	0.647	0.68	$DG = 0.873$
presubiculum	—	0.68	—
subiculum	0.58	0.74	0.770
hippocampal fissure	—	0.53	—
SR/SL/SM	0.428	—	—
fimbria	—	0.51	—
head	—	—	0.902
tail	—	—	0.863

693 to re-implement the protocol using our atlases. At the present time it is unclear how we would adapt these  
 694 protocol to data that we used, where subfield segmentations are defined on  $0.3mm^3$  voxels. However, the  
 695 impact of subfield definitions in the context of our work is an important one and should be considered in  
 696 subsequent studies.

697 One further complication common to all subfield segmentation evaluation is that, by its nature, the  
 698 Dice’s Similarity Coefficient score penalizes structures with high surface area-to-volume ratios. Therefore  
 699 subfield DSC scores will generally be lower than whole hippocampal segmentations. We attempted to put  
 700 this effect into perspective by comparing MAGeT-Brain subfield segmentation agreement with the agreement  
 701 of voxel-shifted manual segmentations (Table 7). The results of this exercise show conclusively, despite the  
 702 very limited number of atlases we had to work with, that MAGeT-Brain subfield segmentations are well  
 703 within the bounds of error of a  $0.3mm^3$  voxel shift.

704 Our overlap DSC values demonstrates that we can reliably reproduce segmentations for the CA1, subicu-  
 705 lum, and CA4/dentate subfields ( $DSC > 0.5$ ). That the CA2/CA3 and molecular layers are less well  
 706 reproduced ( $DSC < 0.5$ ) should not be surprising as these are extremely thin and spatially convoluted  
 707 regions that originally required high-resolution MRI for identification and so it is likely that the extents of  
 708 these regions are well below the resolution and contrast offered by standard T1-weighted images.

709 This points to a larger issue of how to truly validate subfield segmentations, both in high resolution  
 710 images and in standard T1-weighted images. There are several manual subfield segmentation methodologies,  
 711 and they do not agree on which regions can be differentiated, even on high-resolution scans. See Table 9 for a  
 712 comparison of MRI-based manual subfield segmentation methodologies. A further complication is that differ-  
 713 ent researchers have differing operational definitions for the subfields and how they ought to be parcellated.  
 714 The disagreement in the community has led to an international working group devoted to normalizing the  
 715 ontology and segmentation rules for the hippocampal subfields (<http://www.hippocampalsubfields.com/>).  
 716 In addition, there have been recent advances from the Yushkevich group to revise their MRI subfield seg-  
 717 mentation protocol based on anatomy discerned from serial histological acquisitions (Adler et al., 2014).  
 718 The definitional and operational disagreements suggest that direct comparison across automated methods  
 719 using “ground truth”-based overlap similarity metrics, such as Dice’s Similarity Coefficient, are not possi-  
 720 ble without carefully taking into account the differences in underlying segmentation protocols and image  
 721 characteristics.

722 In conclusion, we have demonstrated the viability of leveraging a small number of input atlases to generate  
 723 a large template library and thereby improve segmentation reliability when using multi-atlas methods. We  
 724 demonstrated that this method works robustly over hippocampal definitions, different disease populations,  
 725 and different acquisition types. Finally, we also demonstrate that reliable reproduction of hippocampal

---

726 subfield segmentations in standard 3T T1-weighted images is possible.

727 **5 Acknowledgements**

728 We wish acknowledge support from the CAMH Foundation, thanks to Michael and Sonja Koerner, the Kimel  
729 Family, and the Paul E. Garfinkel New Investigator Catalyst Award. MMC is funded by the W. Garfield  
730 Weston Foundation and ANV is funded by the Canadian Institutes of Health Research, Ontario Mental  
731 Health Foundation, NARSAD, and the National Institute of Mental Health (R01MH099167).

732 Computations were performed on the gpc supercomputer at the SciNet HPC Consortium (Loken et al.,  
733 2010). SciNet is funded by: the Canada Foundation for Innovation under the auspices of Compute Canada;  
734 the Government of Ontario; Ontario Research Fund - Research Excellence; and the University of Toronto.

735 In addition, computations were performed on the CAMH Specialized Computing Cluster. The SCC is  
736 funded by: The Canada Foundation for Innovation, Research Hospital Fund.

737 ADNI Acknowledgements: Data collection and sharing for this project was funded by the Alzheimer's  
738 Disease Neuroimaging Initiative (ADNI) (National Institutes of Health Grant U01 AG024904). ADNI is  
739 funded by the National Institute on Aging, the National Institute of Biomedical Imaging and Bioengineering,  
740 and through generous contributions from the following: Abbott; Alzheimer's Association; Alzheimer's  
741 Drug Discovery Foundation; Amorfix Life Sciences Ltd.; AstraZeneca; Bayer HealthCare; BioClinica, Inc.;  
742 Biogen Idec Inc.; Bristol-Myers Squibb Company; Eisai Inc.; Elan Pharmaceuticals Inc.; Eli Lilly and  
743 Company; F. Hoffmann-La Roche Ltd and its affiliated company Genentech, Inc.; GE Healthcare; Innogenetics,  
744 N.V.; IXICO Ltd.; Janssen Alzheimer Immunotherapy Research Development, LLC.; Johnson & Johnson  
745 Pharmaceutical Research Development LLC.; Medpace, Inc.; Merck & Co., Inc.; Meso Scale Diagnostics,  
746 LLC.; Novartis Pharmaceuticals Corporation; Pfizer Inc.; Servier; Synarc Inc.; and Takeda Pharmaceutical  
747 Company. The Canadian Institutes of Health Research is providing funds to support ADNI clinical sites in  
748 Canada. Private sector contributions are facilitated by the Foundation for the National Institutes of Health  
749 ([www.fnih.org](http://www.fnih.org)). The grantee organization is the Northern California Institute for Research and Education,  
750 and the study is Rev March 26, 2012 coordinated by the Alzheimer's disease Cooperative Study at the  
751 University of California, San Diego. ADNI data are disseminated by the Laboratory for NeuroImaging at  
752 the University of California, Los Angeles. This research was also supported by NIH grants P30 AG010129  
753 and K01 AG030514.

754 We would also like to thank G. Clinton, E. Hazel, and B. Worrell for inspiring this work.

755 **6 Supplementary Materials**

756 **6.1 SNT Hippocampal Labels**

757 Semi-automated hippocampal volumetry was carried out using a commercially available high dimensional  
758 brain mapping tool (Medtronic Surgical Navigation Technologies, Louisville, CO), that has previously been  
759 validated and compared to manual tracing of the hippocampus (Hsu et al., 2002). Measurement of hippocam-  
760 pal volume is achieved first by placing manually 22 control points as local landmarks for the hippocampus on  
761 the individual brain MRI data: one landmark at the hippocampal head, one at the tail, and four per image  
762 (i.e., at the superior, inferior, medial and lateral boundaries) on five equally spaced images perpendicular  
763 to the long axis of the hippocampus. Second, fluid image transformation is used to match the individual

## 6.2 Experiment 5: Whole Hippocampus Segmentation Cross-Validation — Alzheimer’s Disease, SNT Segmentations

---

Table 11: **ADNI1:Complete 1Yr 1.5T SNT 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 N = 23	LMCI N = 23	AD N = 23	Combined N = 69
Age at baseline Years	72.2 75.5 78.5	71.0 77.1 81.4	71.7 77.8 81.8	71.5 76.6 81.3
Sex : Female	43% (10)	43% (10)	43% (10)	43% (30)
Education	16.0 16.0 18.0	15.0 16.0 18.0	12.0 16.0 16.5	14.0 16.0 18.0
CDR-SB	0.00 0.00 0.00	0.75 1.50 1.50	4.00 4.50 5.00	0.00 1.50 4.00
ADAS 13	4.67 5.67 12.34	14.34 16.00 20.50	23.83 29.00 31.66	10.00 16.00 25.33
MMSE	28.5 29.0 30.0	25.0 27.0 28.0	21.0 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.

764 brains to a template brain (Christensen et al., 1997). The pixels corresponding to the hippocampus are then  
 765 labeled and counted to obtain volumes. This method of hippocampal voluming has a documented reliability  
 766 of an intraclass coefficient better than .94 (Hsu et al., 2002).

## 767 6.2 Experiment 5: Whole Hippocampus Segmentation Cross-Validation — Alzheimer’s 768 Disease, SNT Segmentations

769 This experiment is a replication of Experiment 1 using a pool of 69 images and SNT semi-automated  
 770 segmentations from the ADNI dataset (Hsu et al., 2002). See Experiment 1 for full details on the ADNI  
 771 dataset, and validation process.

### 772 6.2.1 Experiment 5: Materials and Methods

773 **Dataset** 69 1.5T images were arbitrarily selected from the baseline scans in the *ADNI1:Complete 1Yr 1.5T*  
 774 standardized dataset. 23 subjects were chosen from each disease category: cognitively normal (CN), mild  
 775 cognitive impairment (MCI) and Alzheimer’s disease (AD). Demographics for this subset are shown in Table  
 776 1. Each image has a corresponding semi-automated segmentation of the left and right whole hippocampus  
 777 made available with ADNI images (SNT; see Supplementary Materials).

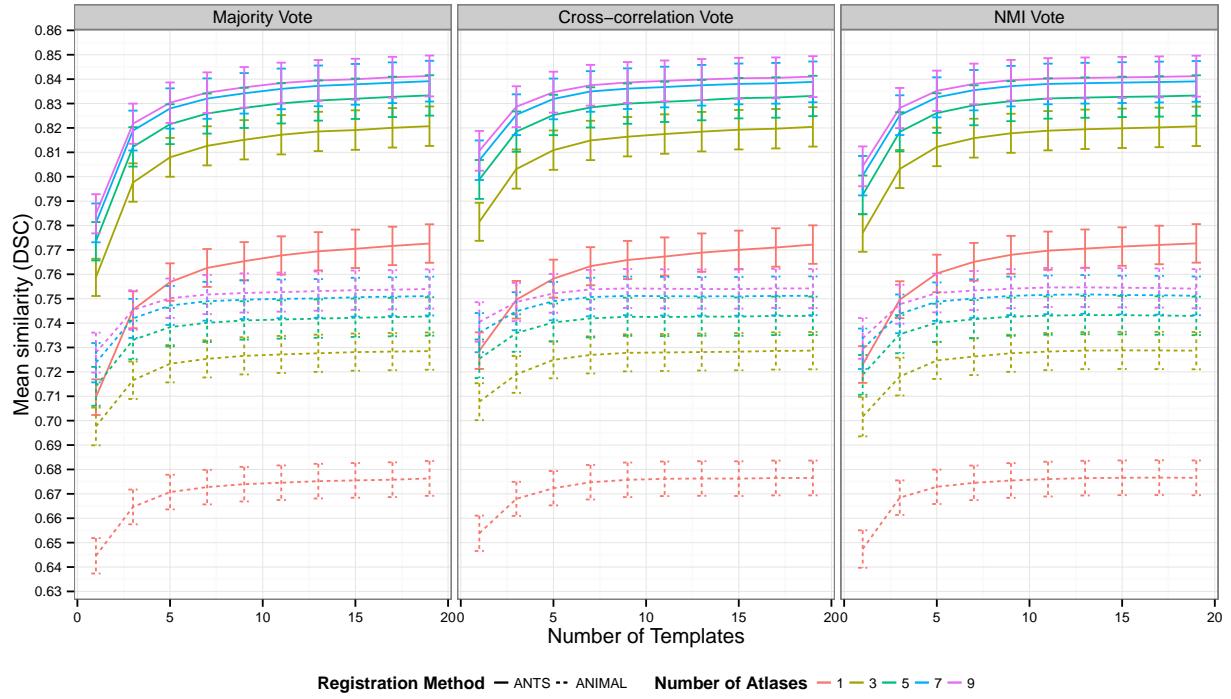
778 **Experiment details** A total of ten validation rounds were performed on each subject in the dataset, for  
 779 each combination of parameter settings: atlas library size (1-9), template library size (1-20), registration  
 780 method (ANTS or ANIMAL), and label fusion method (majority vote, cross-correlation weighted majority  
 781 vote, and normalized mutual information weighted majority vote). A total of  $10 \times 69 \times 9 \times 20 \times 2 \times 3 =$   
 782  $7.452 \times 10^5$  validation rounds are conducted. The computed segmentations for a subject are compared to  
 783 the SNT labels provided by ADNI using Dice’s Similarity Coefficient and the score is averaged over the  
 784 validation rounds.

### 785 6.2.2 Experiment 5: Results

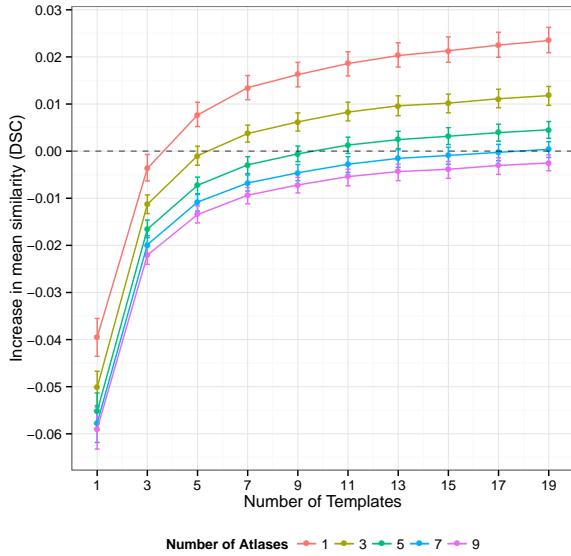
786 As when comparing against manual labels in Experiment 1, we find similar behaviour when comparing  
 787 MAGeT-Brain segmentations to SNT labels: similarity scores increase with increasing numbers of atlases  
 788 and templates, with diminishing increases in improvement trending towards a plateau (Figure 2a). As in

789 Experiment 1, using ANTS registration leads to significantly increased similarity scores, and there is no  
790 significant difference in scores from any of the label fusion methods. Mean DSC score peaks at 0.841 when  
791 using 9 atlases, 19 templates, ANTS registration, and majority vote label fusion. Compared to multi-atlas  
792 segmentations, we find MAGeT-Brain segmentations show increasing improvement with larger atlas and  
793 template libraries when using more than 9 templates and 5 or fewer atlases (Figure 8b). Peak improvement  
794 (+0.023 DSC) is found with a single atlas and template library of 19 images. In addition to a mean increase in  
795 similarity score over multi-atlas-based segmentation, MAGeT-Brain also shows more consistency in similarity  
796 scores across all subjects and validation folds (Figure 8c) with a large enough template library.

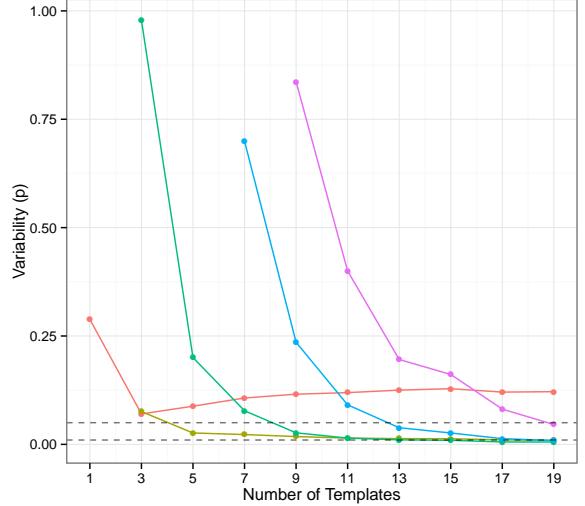
797 S



(a) DSC vs. atlas and template library size



(b) Increase in similarity score over multi-atlas



(c) Difference in variability with multi-atlas

Figure 8: **Whole hippocampus segmentation cross-validation on ADNI subjects with SNT segmentations.** (8a) Average DSC score of MAGeT-Brain with SNT segmentations for 69 ADNI subjects taken over 10 folds of cross-validation at each parameter setting. Error bars indicate standard error. (8b) Increase in DSC of MAGeT-Brain over multi-atlas segmentations. (8c) shows the significance of t-tests comparing the variability in DSC scores of MAGeT-Brain and multi-atlas across validation folds. Only points where MAGeT-Brain mean variability is lower than multi-atlas are shown. Dashed lines indicate p-values of 0.05 and 0.01.

**798 References**

- 799 D. H. Adler, J. Pluta, S. Kadivar, C. Craige, J. C. Gee, B. B. Avants, and P. a. Yushkevich. Histology-derived  
800 volumetric annotation of the human hippocampal subfields in postmortem MRI. *NeuroImage*, 84:505–23,  
801 Jan. 2014. ISSN 1095-9572. doi: 10.1016/j.neuroimage.2013.08.067.
- 802 P. Aljabar, R. a. Heckemann, a. Hammers, J. V. Hajnal, and D. Rueckert. Multi-atlas based segmentation  
803 of brain images: atlas selection and its effect on accuracy. *NeuroImage*, 46(3):726–38, July 2009. ISSN  
804 1095-9572. doi: 10.1016/j.neuroimage.2009.02.018.
- 805 B. B. Avants, C. L. Epstein, M. Grossman, and J. C. Gee. Symmetric diffeomorphic image registration with  
806 cross-correlation: evaluating automated labeling of elderly and neurodegenerative brain. *Medical image  
807 analysis*, 12(1):26–41, Feb. 2008. ISSN 1361-8423. doi: 10.1016/j.media.2007.06.004.
- 808 J. Barnes, J. Foster, R. G. Boyes, T. Pepple, E. K. Moore, J. M. Schott, C. Frost, R. I. Scahill, and N. C. Fox.  
809 A comparison of methods for the automated calculation of volumes and atrophy rates in the hippocampus.  
810 *NeuroImage*, 40(4):1655–71, May 2008. ISSN 1053-8119. doi: 10.1016/j.neuroimage.2008.01.012.
- 811 J. M. Bland and D. Altman. Statistical methods for assessing agreement between two methods of clinical  
812 measurement. *The lancet*, pages 307–310, 1986.
- 813 M. Boccardi, M. Bocchetta, L. G. Apostolova, G. Preboske, N. Robitaille, P. Pasqualetti, L. D. Collins,  
814 S. Duchesne, C. R. Jack, and G. B. Frisoni. Establishing Magnetic Resonance Images Orientation for the  
815 EADC-ADNI Manual Hippocampal Segmentation Protocol. *Journal of neuroimaging : official journal of  
816 the American Society of Neuroimaging*, pages 1–6, Nov. 2013a. ISSN 1552-6569. doi: 10.1111/jon.12065.
- 817 M. Boccardi, M. Bocchetta, R. Ganzola, N. Robitaille, A. Redolfi, S. Duchesne, C. R. Jack, and G. B. Frisoni.  
818 Operationalizing protocol differences for EADC-ADNI manual hippocampal segmentation. *Alzheimer's &  
819 dementia : the journal of the Alzheimer's Association*, pages 1–11, May 2013b. ISSN 1552-5279. doi:  
820 10.1016/j.jalz.2013.03.001.
- 821 M. M. Chakravarty, A. F. Sadikot, S. Mongia, G. Bertrand, and D. L. Collins. Towards a multi-modal  
822 atlas for neurosurgical planning. *Medical image computing and computer-assisted intervention : MICCAI  
823 ... International Conference on Medical Image Computing and Computer-Assisted Intervention*, 9(Pt 2):  
824 389–96, Jan. 2006.
- 825 M. M. Chakravarty, A. F. Sadikot, J. Germann, G. Bertrand, and D. L. Collins. Towards a validation  
826 of atlas warping techniques. *Medical image analysis*, 12(6):713–26, Dec. 2008. ISSN 1361-8423. doi:  
827 10.1016/j.media.2008.04.003.
- 828 M. M. Chakravarty, A. F. Sadikot, J. Germann, P. Hellier, G. Bertrand, and D. L. Collins. Comparison  
829 of piece-wise linear, linear, and nonlinear atlas-to-patient warping techniques: analysis of the labeling of  
830 subcortical nuclei for functional neurosurgical applications. *Human brain mapping*, 30(11):3574–95, Nov.  
831 2009. ISSN 1097-0193. doi: 10.1002/hbm.20780.
- 832 M. M. Chakravarty, P. Steadman, M. C. van Eede, R. D. Calcott, V. Gu, P. Shaw, A. Raznahan, D. L. Collins,  
833 and J. P. Lerch. Performing label-fusion-based segmentation using multiple automatically generated tem-  
834 plates. *Human brain mapping*, 34(10):2635–54, Oct. 2013. ISSN 1097-0193. doi: 10.1002/hbm.22092.

## REFERENCES

---

- 835 G. E. Christensen, S. C. Joshi, and M. I. Miller. Volumetric transformation of brain anatomy. *IEEE*  
836 *transactions on medical imaging*, 16(6):864–77, Dec. 1997. ISSN 0278-0062. doi: 10.1109/42.650882.
- 837 M. Chupin, E. Gérardin, R. Cuingnet, C. Boutet, L. Lemieux, S. Lehéricy, H. Benali, L. Garnero, and  
838 O. Colliot. Fully automatic hippocampus segmentation and classification in Alzheimer’s disease and mild  
839 cognitive impairment applied on data from ADNI. *Hippocampus*, 19(6):579–87, June 2009. ISSN 1098-1063.  
840 doi: 10.1002/hipo.20626.
- 841 D. L. Collins and J. C. Pruessner. Towards accurate, automatic segmentation of the hippocampus and  
842 amygdala from MRI by augmenting ANIMAL with a template library and label fusion. *NeuroImage*, 52  
843 (4):1355–66, Oct. 2010. ISSN 1095-9572. doi: 10.1016/j.neuroimage.2010.04.193.
- 844 D. L. Collins, P. Neelin, T. M. Peters, and A. C. Evans. Automatic 3D intersubject registration of MR  
845 volumetric data in standardized Talairach space. *Journal of computer assisted tomography*, 18(2):192–205,  
846 1994. ISSN 0363-8715.
- 847 D. L. Collins, C. J. Holmes, T. M. Peters, and A. C. Evans. Automatic 3-D model-based neuroanatomical  
848 segmentation. *Human Brain Mapping*, 3(3):190–208, Oct. 1995. ISSN 10659471. doi: 10.1002/hbm.  
849 460030304.
- 850 P. Coupe, V. Fonov, S. Eskildsen, J. Manjón, D. Arnold, and L. Collins. Influence of the training library  
851 composition on a patch-based label fusion method: Application to hippocampus segmentation on the ADNI  
852 dataset. *Alzheimer’s & Dementia*, 7(4):S316, July 2011. ISSN 15525260. doi: 10.1016/j.jalz.2011.05.918.
- 853 P. Coupé, S. F. Eskildsen, J. V. Manjón, V. S. Fonov, and D. L. Collins. Simultaneous segmentation and  
854 grading of anatomical structures for patient’s classification: application to Alzheimer’s disease. *NeuroIm-*  
855 *age*, 59(4):3736–47, Feb. 2012. ISSN 1095-9572. doi: 10.1016/j.neuroimage.2011.10.080.
- 856 J. G. Csernansky, S. Joshi, L. Wang, J. W. Haller, M. Gado, J. P. Miller, U. Grenander, and M. I. Miller. Hip-  
857 pocampal morphometry in schizophrenia by high dimensional brain mapping. *Proceedings of the National*  
858 *Academy of Sciences of the United States of America*, 95(19):11406–11411, 1998.
- 859 T. den Heijer, F. V. der Lijn, M. W. Vernooij, M. de Groot, P. J. Koudstaal, a. V. der Lugt, G. P.  
860 Krestin, a. Hofman, W. J. Niessen, and M. M. B. Breteler. Structural and diffusion MRI measures of  
861 the hippocampus and memory performance. *NeuroImage*, 63(4):1782–9, Dec. 2012. ISSN 1095-9572. doi:  
862 10.1016/j.neuroimage.2012.08.067.
- 863 B. Fischl, D. H. Salat, E. Busa, M. Albert, M. Dieterich, C. Haselgrove, A. van der Kouwe, R. Killiany,  
864 D. Kennedy, S. Klaveness, A. Montillo, N. Makris, B. Rosen, and A. M. Dale. Whole brain segmentation:  
865 automated labeling of neuroanatomical structures in the human brain. *Neuron*, 33(3):341–55, Jan. 2002.  
866 ISSN 0896-6273.
- 867 E. Geuze, E. Vermetten, and J. D. Bremner. MR-based in vivo hippocampal volumetrics: 2. Findings in  
868 neuropsychiatric disorders. *Molecular Psychiatry*, 10(2):160, Sept. 2004. doi: 10.1038/sj.mp.4001579.
- 869 I. S. Gousias, D. Rueckert, R. a. Heckemann, L. E. Dyet, J. P. Boardman, a. D. Edwards, and A. Hammers.  
870 Automatic segmentation of brain MRIs of 2-year-olds into 83 regions of interest. *NeuroImage*, 40(2):  
871 672–84, Apr. 2008. ISSN 1053-8119. doi: 10.1016/j.neuroimage.2007.11.034.

## REFERENCES

---

- 872 J. W. Haller, A. Banerjee, G. E. Christensen, M. Gado, S. Joshi, M. I. Miller, Y. Sheline, M. W. Van-  
873 nier, and J. G. Csernansky. Three-dimensional hippocampal MR morphometry with high-dimensional  
874 transformation of a neuroanatomic atlas. *Radiology*, 202(2):504–510, 1997.
- 875 M. Hartig, D. Truran-sacrey, S. Raptentsetsang, N. Schuff, and M. Weiner. USCF FreeSurfer Overview and  
876 QC Ratings. 2010.
- 877 R. A. Heckemann, J. V. Hajnal, P. Aljabar, D. Rueckert, and A. Hammers. Automatic anatomical brain  
878 MRI segmentation combining label propagation and decision fusion. *NeuroImage*, 46(3):726–38, July  
879 2006a. ISSN 1095-9572. doi: 10.1016/j.neuroimage.2009.02.018.
- 880 R. A. Heckemann, J. V. Hajnal, P. Aljabar, D. Rueckert, and A. Hammers. Automatic anatomical brain MRI  
881 segmentation combining label propagation and decision fusion. *NeuroImage*, 33(1):115–26, Oct. 2006b.  
882 ISSN 1053-8119. doi: 10.1016/j.neuroimage.2006.05.061.
- 883 R. A. Heckemann, S. Keihaninejad, P. Aljabar, K. R. Gray, C. Nielsen, D. Rueckert, J. V. Hajnal, and  
884 A. Hammers. Automatic morphometry in Alzheimer’s disease and mild cognitive impairment. *NeuroImage*,  
885 56(4):2024–37, July 2011. ISSN 1095-9572. doi: 10.1016/j.neuroimage.2011.03.014.
- 886 Y.-Y. Hsu, N. Schuff, A.-T. Du, K. Mark, X. Zhu, D. Hardin, and M. W. Weiner. Comparison of automated  
887 and manual MRI volumetry of hippocampus in normal aging and dementia. *Journal of magnetic resonance*  
888 *imaging : JMRI*, 16(3):305–10, Sept. 2002. ISSN 1053-1807. doi: 10.1002/jmri.10163.
- 889 C. R. Jack, M. A. Bernstein, N. C. Fox, P. Thompson, G. Alexander, D. Harvey, B. Borowski, P. J. Britson,  
890 J. L Whitwell, C. Ward, A. M. Dale, J. P. Felmlee, J. L. Gunter, D. L. G. Hill, R. Killiany, N. Schuff,  
891 S. Fox-Bosetti, C. Lin, C. Studholme, C. S. DeCarli, G. Krueger, H. A. Ward, G. J. Metzger, K. T.  
892 Scott, R. Mallozzi, D. Blezek, J. Levy, J. P. Debbins, A. S. Fleisher, M. Albert, R. Green, G. Bartzokis,  
893 G. Glover, J. Mugler, and M. W. Weiner. The Alzheimer’s Disease Neuroimaging Initiative (ADNI): MRI  
894 methods. *Journal of magnetic resonance imaging : JMRI*, 27(4):685–91, Apr. 2008. ISSN 1053-1807. doi:  
895 10.1002/jmri.21049.
- 896 C. R. Jack, F. Barkhof, M. A. Bernstein, M. Cantillon, P. E. Cole, C. Decarli, B. Dubois, S. Duchesne,  
897 N. C. Fox, G. B. Frisoni, H. Hampel, D. L. G. Hill, K. Johnson, J.-F. Mangin, P. Scheltens, A. J. Schwarz,  
898 R. Sperling, J. Suhy, P. M. Thompson, M. Weiner, and N. L. Foster. Steps to standardization and validation  
899 of hippocampal volumetry as a biomarker in clinical trials and diagnostic criterion for Alzheimer’s disease.  
900 *Alzheimer’s & dementia : the journal of the Alzheimer’s Association*, 7(4):474–485.e4, July 2011. ISSN  
901 1552-5279. doi: 10.1016/j.jalz.2011.04.007.
- 902 A. Jeneson and L. Squire. Working memory, long-term memory, and medial temporal lobe function. *Learning*  
903 *& Memory*, 19(1):15–25, 2012. doi: 10.1101/lm.024018.111.
- 904 M. S. Karnik-Henry, L. Wang, D. M. Barch, M. P. Harms, C. Campanella, and J. G. Csernansky. Medial  
905 temporal lobe structure and cognition in individuals with schizophrenia and in their non-psychotic siblings.  
906 *Schizophrenia research*, 138(2-3):128–35, July 2012. ISSN 1573-2509. doi: 10.1016/j.schres.2012.03.015.
- 907 K. K. Leung, J. Barnes, G. R. Ridgway, J. W. Bartlett, M. J. Clarkson, K. Macdonald, N. Schuff, N. C. Fox,  
908 and S. Ourselin. Automated cross-sectional and longitudinal hippocampal volume measurement in mild  
909 cognitive impairment and Alzheimer’s disease. *NeuroImage*, 51(4):1345–59, July 2010. ISSN 1095-9572.  
910 doi: 10.1016/j.neuroimage.2010.03.018.

- 911 C. Loken, D. Gruner, L. Groer, R. Peltier, N. Bunn, M. Craig, T. Henriques, J. Dempsey, C.-H. Yu, J. Chen,  
912 L. J. Dursi, J. Chong, S. Northrup, J. Pinto, N. Knecht, and R. V. Zon. SciNet: Lessons Learned from  
913 Building a Power-efficient Top-20 System and Data Centre. *Journal of Physics: Conference Series*, 256:  
914 012026, Nov. 2010. ISSN 1742-6596. doi: 10.1088/1742-6596/256/1/012026.
- 915 J. M. Lötjönen, R. Wolz, J. R. Koikkalainen, L. Thurfjell, G. Waldemar, H. Soininen, and D. Rueckert. Fast  
916 and robust multi-atlas segmentation of brain magnetic resonance images. *NeuroImage*, 49(3):2352–65,  
917 Mar. 2010. ISSN 1095-9572. doi: 10.1016/j.neuroimage.2009.10.026.
- 918 A. Malla, R. Norman, T. McLean, D. Scholten, and L. Townsend. A Canadian programme for early inter-  
919 vention in non-affective psychotic disorders. *The Australian and New Zealand journal of psychiatry*, 37  
920 (4):407–13, Aug. 2003. ISSN 0004-8674.
- 921 J. Mazziotta, A. Toga, A. Evans, P. Fox, J. Lancaster, K. Zilles, R. Woods, T. Paus, G. Simpson, B. Pike,  
922 C. Holmes, L. Collins, P. Thompson, D. MacDonald, M. Iacoboni, T. Schormann, K. Amunts, N. Palomero-  
923 Gallagher, S. Geyer, L. Parsons, K. Narr, N. Kabani, G. Le Goualher, J. Feidler, K. Smith, D. Boomsma,  
924 H. Hulshoff Pol, T. Cannon, R. Kawashima, and B. Mazoyer. A four-dimensional probabilistic atlas of  
925 the human brain. *Journal of the American Medical Informatics Association : JAMIA*, 8(5):401–30. ISSN  
926 1067-5027.
- 927 J. Mazziotta, A. Toga, A. Evans, P. Fox, J. Lancaster, K. Zilles, R. Woods, T. Paus, G. Simpson,  
928 B. Pike, C. Holmes, L. Collins, P. Thompson, D. MacDonald, M. Iacoboni, T. Schormann, K. Amunts,  
929 N. Palomero-Gallagher, S. Geyer, L. Parsons, K. Narr, N. Kabani, G. Le Goualher, D. Boomsma,  
930 T. Cannon, R. Kawashima, and B. Mazoyer. A probabilistic atlas and reference system for the hu-  
931 man brain: International Consortium for Brain Mapping (ICBM). *Philosophical transactions of the Royal  
932 Society of London. Series B, Biological sciences*, 356(1412):1293–322, Aug. 2001. ISSN 0962-8436. doi:  
933 10.1098/rstb.2001.0915.
- 934 J. C. Mazziotta, A. W. Toga, A. Evans, P. Fox, and J. Lancaster. A probabilistic atlas of the human brain:  
935 theory and rationale for its development. The International Consortium for Brain Mapping (ICBM).  
936 *NeuroImage*, 2(2):89–101, June 1995. ISSN 1053-8119.
- 937 J. H. Morra, Z. Tu, L. G. Apostolova, A. E. Green, C. Avedissian, S. K. Madsen, N. Parikshak, X. Hua, A. W.  
938 Toga, C. R. Jack, M. W. Weiner, and P. M. Thompson. Validation of a fully automated 3D hippocampal  
939 segmentation method using subjects with Alzheimer's disease mild cognitive impairment, and elderly  
940 controls. *NeuroImage*, 43(1):59–68, Oct. 2008. ISSN 1095-9572. doi: 10.1016/j.neuroimage.2008.07.003.
- 941 S. Mueller, L. Stables, A. Du, and N. Schuff. Measurement of hippocampal subfields and age-related changes  
942 with high resolution MRI at 4T. *Neurobiology of ...*, 1(5):719–726, 2007. doi: 10.1016/j.neurobiolaging.  
943 2006.03.007.
- 944 S. G. Mueller and M. W. Weiner. Selective effect of age, Apo e4, and Alzheimer's disease on hippocampal  
945 subfields. *Hippocampus*, 19(6):558–64, June 2009. ISSN 1098-1063. doi: 10.1002/hipo.20614.
- 946 K. L. Narr, P. M. Thompson, P. Szeszko, D. Robinson, S. Jang, R. P. Woods, S. Kim, K. M. Hayashi,  
947 D. Asunction, A. W. Toga, and R. M. Bilder. Regional specificity of hippocampal volume reductions  
948 in first-episode schizophrenia. *NeuroImage*, 21(4):1563–75, Apr. 2004. ISSN 1053-8119. doi: 10.1016/j.  
949 neuroimage.2003.11.011.

## REFERENCES

---

- 950 S. M. Nestor, E. Gibson, F.-Q. Gao, A. Kiss, and S. E. Black. A Direct Morphometric Comparison of Five  
951 Labeling Protocols for Multi-Atlas Driven Automatic Segmentation of the Hippocampus in Alzheimer's  
952 Disease. *NeuroImage*, Nov. 2012. ISSN 1095-9572. doi: 10.1016/j.neuroimage.2012.10.081.
- 953 Z. Pausova, T. Paus, M. Abrahamowicz, J. Almerigi, N. Arbour, M. Bernard, D. Gaudet, P. Hanzalek,  
954 P. Hamet, A. C. Evans, M. Kramer, L. Laberge, S. M. Leal, G. Leonard, J. Lerner, R. M. Lerner, J. Math-  
955 ieu, M. Perron, B. Pike, A. Pitiot, L. Richer, J. R. Séguin, C. Syme, R. Toro, R. E. Tremblay, S. Veillette,  
956 and K. Watkins. Genes, maternal smoking, and the offspring brain and body during adolescence: design  
957 of the Saguenay Youth Study. *Human brain mapping*, 28(6):502–18, June 2007. ISSN 1065-9471. doi:  
958 10.1002/hbm.20402.
- 959 K. M. Pohl, S. Bouix, M. Nakamura, T. Rohlfing, R. W. McCarley, R. Kikinis, W. E. L. Grimson, M. E.  
960 Shenton, and W. M. Wells. A hierarchical algorithm for MR brain image parcellation. *IEEE transactions  
961 on medical imaging*, 26(9):1201–12, Sept. 2007. ISSN 0278-0062. doi: 10.1109/TMI.2007.901433.
- 962 J. Poppenk and M. Moscovitch. A Hippocampal Marker of Recollection Memory Ability among Healthy  
963 Young Adults: Contributions of Posterior and Anterior Segments. *Neuron*, 72(6):931–937, Dec. 2011.  
964 ISSN 0896-6273. doi: 10.1016/j.neuron.2011.10.014.
- 965 S. Powell, V. A. Magnotta, H. Johnson, V. K. Jammalamadaka, R. Pierson, and N. C. Andreasen. Registration  
966 and machine learning-based automated segmentation of subcortical and cerebellar brain structures.  
967 *NeuroImage*, 39(1):238–47, Jan. 2008. ISSN 1053-8119. doi: 10.1016/j.neuroimage.2007.05.063.
- 968 J. C. Pruessner, L. M. Li, W. Serles, M. Pruessner, D. L. Collins, N. Kabani, S. Lupien, and A. C. Evans.  
969 Volumetry of hippocampus and amygdala with high-resolution MRI and three-dimensional analysis soft-  
970 ware: minimizing the discrepancies between laboratories. *Cerebral cortex (New York, N.Y. : 1991)*, 10  
971 (4):433–42, Apr. 2000. ISSN 1047-3211.
- 972 S. Robbins, A. C. Evans, D. L. Collins, and S. Whitesides. Tuning and comparing spatial normalization  
973 methods. *Medical image analysis*, 8(3):311–23, Sept. 2004. ISSN 1361-8415. doi: 10.1016/j.media.2004.  
974 06.009.
- 975 N. Robitaille and S. Duchesne. Label fusion strategy selection. *International journal of biomedical imaging*,  
976 2012:431095, Jan. 2012. ISSN 1687-4196. doi: 10.1155/2012/431095.
- 977 M. R. Sabuncu, R. S. Desikan, J. Sepulcre, B. T. T. Yeo, H. Liu, N. J. Schmansky, M. Reuter, M. W.  
978 Weiner, R. L. Buckner, R. a. Sperling, and B. Fischl. The dynamics of cortical and hippocampal atrophy  
979 in Alzheimer disease. *Archives of neurology*, 68(8):1040–8, Aug. 2011. ISSN 1538-3687. doi: 10.1001/  
980 archneurol.2011.167.
- 981 W. B. Scoville and B. Milner. Loss of recent memory after bilateral hippocampal lesions. 1957. *The Journal  
982 of neuropsychiatry and clinical neurosciences*, 12(1):103–113, 2000.
- 983 J. Shao. Linear Model Selection by Cross-validation. *Journal of the American Statistical Association*, 88  
984 (422):486–494, June 1993. ISSN 0162-1459. doi: 10.1080/01621459.1993.10476299.
- 985 J. G. Sled, a. P. Zijdenbos, and a. C. Evans. A nonparametric method for automatic correction of intensity  
986 nonuniformity in MRI data. *IEEE transactions on medical imaging*, 17(1):87–97, Feb. 1998. ISSN 0278-  
987 0062. doi: 10.1109/42.668698.

## REFERENCES

---

- 988 C. Studholme, D. Hill, and D. Hawkes. An overlap invariant entropy measure of 3D medical image alignment.  
989 *Pattern Recognition*, 32(1):71–86, Jan. 1999. ISSN 00313203. doi: 10.1016/S0031-3203(98)00091-0.
- 990 C. Studholme, E. Novotny, I. G. Zubal, and J. S. Duncan. Estimating tissue deformation between functional  
991 images induced by intracranial electrode implantation using anatomical MRI. *NeuroImage*, 13(4):561–76,  
992 Apr. 2001. ISSN 1053-8119. doi: 10.1006/nimg.2000.0692.
- 993 F. van der Lijn, T. den Heijer, M. M. B. Breteler, and W. J. Niessen. Hippocampus segmentation in MR  
994 images using atlas registration, voxel classification, and graph cuts. *NeuroImage*, 43(4):708–20, Dec. 2008.  
995 ISSN 1095-9572. doi: 10.1016/j.neuroimage.2008.07.058.
- 996 K. Van Leemput, A. Bakkour, T. Benner, G. Wiggins, L. L. Wald, J. Augustinack, B. C. Dickerson, P. Gol-  
997 land, and B. Fischl. Automated segmentation of hippocampal subfields from ultra-high resolution in vivo  
998 MRI. *Hippocampus*, 19(6):549–57, June 2009. ISSN 1098-1063. doi: 10.1002/hipo.20615.
- 999 H. Wang, J. W. Suh, J. Pluta, M. Altinay, and P. Yushkevich. Optimal weights for multi-atlas label fusion.  
1000 *Information processing in medical imaging : proceedings of the ... conference*, 22:73–84, Jan. 2011. ISSN  
1001 1011-2499.
- 1002 S. K. Warfield, K. H. Zou, and W. M. Wells. Simultaneous truth and performance level estimation (STAPLE):  
1003 an algorithm for the validation of image segmentation. *IEEE transactions on medical imaging*, 23(7):903–  
1004 21, July 2004. ISSN 0278-0062. doi: 10.1109/TMI.2004.828354.
- 1005 J. L. Winterburn, J. C. Pruessner, S. Chavez, M. M. Schira, N. J. Lobaugh, A. N. Voineskos, and M. M.  
1006 Chakravarty. A novel in vivo atlas of human hippocampal subfields using high-resolution 3 T magnetic  
1007 resonance imaging. *NeuroImage*, 74:254–65, July 2013. ISSN 1095-9572. doi: 10.1016/j.neuroimage.2013.  
1008 02.003.
- 1009 L. E. M. Wisse, L. Gerritsen, J. J. M. Zwanenburg, H. J. Kuijf, P. R. Luijten, G. J. Biessels, and M. I.  
1010 Geerlings. Subfields of the hippocampal formation at 7 T MRI: in vivo volumetric assessment. *NeuroImage*,  
1011 61(4):1043–9, July 2012. ISSN 1095-9572. doi: 10.1016/j.neuroimage.2012.03.023.
- 1012 J. Wixted and L. Squire. The medial temporal lobe and the attributes of memory. *Trends in cognitive  
1013 sciences*, 15(5):210–217, 2011. doi: 10.1016/j.tics.2011.03.005.
- 1014 R. Wolz, P. Aljabar, J. V. Hajnal, A. Hammers, and D. Rueckert. LEAP: learning embeddings for atlas  
1015 propagation. *NeuroImage*, 49(2):1316–25, Jan. 2010. ISSN 1095-9572. doi: 10.1016/j.neuroimage.2009.09.  
1016 069.
- 1017 B. T. Wyman, D. J. Harvey, K. Crawford, M. A. Bernstein, O. Carmichael, P. E. Cole, P. K. Crane,  
1018 C. Decarli, N. C. Fox, J. L. Gunter, D. Hill, R. J. Killiany, C. Pachai, A. J. Schwarz, N. Schuff, M. L.  
1019 Senjem, J. Suhy, P. M. Thompson, M. Weiner, and C. R. Jack. Standardization of analysis sets for reporting  
1020 results from ADNI MRI data. *Alzheimer's & dementia : the journal of the Alzheimer's Association*, Oct.  
1021 2012. ISSN 1552-5279. doi: 10.1016/j.jalz.2012.06.004.
- 1022 J. Yelnik, E. Bardinet, D. Dormont, G. Malandain, S. Ourselin, D. Tandé, C. Karachi, N. Ayache, P. Cornu,  
1023 and Y. Agid. A three-dimensional, histological and deformable atlas of the human basal ganglia. I. Atlas  
1024 construction based on immunohistochemical and MRI data. *NeuroImage*, 34(2):618–38, Jan. 2007. ISSN  
1025 1053-8119. doi: 10.1016/j.neuroimage.2006.09.026.

## *REFERENCES*

---

- 1026 P. A. Yushkevich, B. B. Avants, J. Pluta, S. Das, D. Minkoff, D. Mechanic-Hamilton, S. Glynn, S. Pickup,  
1027 W. Liu, J. C. Gee, M. Grossman, and J. A. Detre. A high-resolution computational atlas of the human  
1028 hippocampus from postmortem magnetic resonance imaging at 9.4 T. *NeuroImage*, 44(2):385–98, Jan.  
1029 2009. ISSN 1095-9572. doi: 10.1016/j.neuroimage.2008.08.042.
- 1030 P. A. Yushkevich, H. Wang, J. Pluta, S. R. Das, C. Craige, B. B. Avants, M. W. Weiner, and S. Mueller.  
1031 Nearly automatic segmentation of hippocampal subfields in in vivo focal T2-weighted MRI. *NeuroImage*,  
1032 53(4):1208–24, Dec. 2010. ISSN 1095-9572. doi: 10.1016/j.neuroimage.2010.06.040.

**LaTeX Source Files**

[Click here to download LaTeX Source Files: dist.zip](#)