

# Predictive Markers for AD in a Multi-Modality Framework: An Analysis of MCI Progression in the ADNI Population

Chris Hinrichs<sup>a,b,\*†</sup> Vikas Singh<sup>b,a</sup> Guofan Xu<sup>c,d</sup> Sterling C. Johnson<sup>c,d</sup>  
and the Alzheimers Disease Neuroimaging Initiative<sup>‡</sup>

## Abstract

Alzheimer’s Disease (AD) and other neurodegenerative diseases affect over 20 million people worldwide, and this number is projected to significantly increase in the coming decades. Proposed imaging-based markers have shown steadily improving levels of sensitivity/specificity in classifying individual subjects as AD or normal. Several of these efforts have utilized statistical machine learning techniques, using brain images as input, as means of deriving such AD-related markers. A common characteristic of this line of research is a focus on either (1) using a single imaging modality for classification, or (2) incorporating several modalities, but reporting *separate* results for each. One strategy to improve on the success of these methods is to leverage *all* available imaging modalities *together* in a single automated learning framework. The rationale is that some subjects may show signs of pathology in one modality but not in another – by combining all available images a clearer view of the progression of disease pathology will emerge. Our method is based on the Multi-Kernel Learning (MKL) framework, which allows the inclusion of an arbitrary number of views of the data in a maximum margin, kernel learning framework. The principal innovation behind MKL is that it learns an optimal combination of kernel (similarity) matrices while simultaneously training a classifier. **In classification experiments MKL outperformed an SVM trained on all available features by 3% – 4%.** We are especially interested in whether such markers are capable of identifying *early* signs of the disease. To address this question, we have examined whether our multi-modal disease marker (MMDM) can predict conversion from Mild Cognitive Impairment (MCI) to AD. Our experiments reveal that this measure shows significant group differences between MCI subjects who progressed to AD, and those who remained stable for 3 years. **These differences were most significant in MMDMs based on imaging data.** We also discuss the relationship between our MMDM and an individual’s conversion from MCI to AD.

## 1 Introduction

A significant body of existing literature (Johnson et al., 2006; Whitwell et al., 2007; Reiman et al., 1996; Canu et al., 2010; Thompson and Apostolova, 2007) suggests that pathological manifestations of Alzheimer’s disease begin many years before the patient becomes *symptomatic* – which is typically when cognitive tests can be used to make a diagnosis (Albert et al., 2001). Unfortunately, by this time significant neurodegeneration has already occurred. In an effort to identify AD-related changes early, a promising direction of ongoing research is focused on exploiting advanced imaging-based techniques to characterize prominent neurodegenerative patterns during the prodromal stages of the disease, when only mild symptoms of the disease are evident. A set of recent papers (Davatzikos et al., 2008a,b; Fan et al., 2008b; Vemuri et al., 2008) including work from our

\*Corresponding author. 5765 Medical Science Center, Madison, WI 53706, USA

†<sup>a</sup>Department of Computer Sciences, University of Wisconsin-Madison, Madison, WI 53706.

<sup>b</sup>Department of Biostatistics and Medical Informatics, University of Wisconsin-Madison Madison, WI 53705.

<sup>c</sup>William S. Middleton VA Medical Center, Madison, WI 53792.

<sup>d</sup>Department of Medicine, University of Wisconsin-Madison Madison, WI 53792.

Email addresses: hinrichs@cs.wisc.edu (Chris Hinrichs), vsingh@biostat.wisc.edu (Vikas Singh) gxu@medicine.wisc.edu (Guofan Xu), scj@medicine.wisc.edu (Sterling Johnson)

<sup>‡</sup>Data used in the preparation of this article were obtained from the Alzheimers Disease Neuroimaging Initiative (ADNI) database <http://www.loni.ucla.edu/ADNI>. 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. ADNI investigators include (complete listing available at [http://www.loni.ucla.edu/ADNI/Collaboration/ADNI\\_Manuscript\\_Citations.pdf](http://www.loni.ucla.edu/ADNI/Collaboration/ADNI_Manuscript_Citations.pdf))

35 group (Hinrichs et al., 2009a,b) have demonstrated that this is indeed feasible by leveraging and extending  
36 state-of-the-art methods from Statistical Machine Learning and Computer Vision. Good discrimination (in  
37 identifying whether an image corresponds to a control or AD subject) has been obtained on classification  
38 tasks making use of MR *or* FDG-PET images (*i.e.*, *one* type of image data) (Davatzikos et al., 2008a,b; Fan  
39 et al., 2008b; Vemuri et al., 2008; Hinrichs et al., 2009a). A natural question then is whether we can exploit  
40 data from multiple modalities and biological measures (if available) *in conjunction* to (1) obtain improved  
41 accuracy, and (2) identify more subtle class differences (*e.g.*, sub-groups within MCI). This paper considers  
42 exactly this problem – *i.e.*, methods for systematic combination of multiple imaging modalities and clinical  
43 data for classification (*i.e.*, class prediction) at the level of individual subjects.

44 Recently, we have seen evidence that various aspects of AD-related neurodegeneration such as structural  
45 atrophy (Jack Jr. et al., 2005; deToledo-Morrell et al., 2004; Thompson et al., 2001), decreased blood  
46 perfusion (Ramírez et al., 2009), and decreased glucose metabolism (Hoffman et al., 2000; Matsuda, 2001;  
47 Minoshima et al., 1994) can be identified (in structural and functional images) in Mild Cognitive Impaired  
48 (MCI) and AD subjects, as well as at-risk individuals (Small et al., 2000; Querbes et al., 2009; Davatzikos  
49 et al., 2009). A number of groups have made significant progress by adapting well-known machine learning  
50 tools to the problem – this includes Support Vector Machines (SVMs), logistic regression, boosting, and  
51 other classification mechanisms. In the usual classification setting, a number of image acquisitions (training  
52 examples) are provided for which the subjects’ clinical diagnosis is as certain as diagnostically possible.  
53 The objective is to choose a discriminating function which optimizes a statistical measure of the likelihood  
54 of correctly labeling ‘future’ examples. Such measures may be based on certain brain regions, (*e.g.*, the  
55 hippocampus or posterior cingulate cortex) for example. The function’s output can then be used as a targeted  
56 disease marker in individuals that are not part of the training cohort. In the remainder of this section,  
57 we briefly review several interesting AD classification-focused research efforts, and lay the groundwork for  
58 introducing our contributions (*i.e.*, truly multi-modal analysis).

59 The machine learning, or classification approach has been used to provide markers for various neurological  
60 disorders including Alzheimer’s disease (Davatzikos et al., 2008b; Klöppel et al., 2008; Vemuri et al., 2008;  
61 Duchesne et al., 2008; Arimura et al., 2008; Soriano-Mas et al., 2007; Shen et al., 2003; Demirci et al., 2008).  
62 These efforts have primarily utilized brain *images*, though some have also used other available biological  
63 measures. In (Fan et al., 2008b,a; Davatzikos et al., 2008a,b), the authors implemented a classification /  
64 pattern recognition technique using structural (sMR) images provided by the Baltimore Longitudinal Study  
65 of Aging (BLSA) dataset (Shock et al., 1984). The proposed methodology was to first segment the images  
66 into different tissue types, and then perform a non-linear warp to a common template space to allow voxel-  
67 wise comparisons. Next, voxels were selected to serve as “features” (using statistical measures of (clinical)  
68 group differences), used to train a linear Support Vector Machine (SVM) (Bishop, 2006). The reported  
69 accuracy was quite encouraging. The authors of (Klöppel et al., 2008) also used linear SVMs to classify AD  
70 subjects from controls using whole-brain MR images. An additional focus of their research was to separate  
71 AD cases from Frontal Temporal Lobar Degeneration (FTLD). The authors reported high accuracy (> 90%)  
72 on confirmed AD patients, and less where post-mortem diagnosis was unavailable. In related work, Vemuri  
73 *et. al.* (Vemuri et al., 2008) demonstrated a slightly different method of applying linear SVMs on another  
74 dataset obtaining 88 – 90% classification accuracy. More recently, the methods in (Fan et al., 2008a; Misra  
75 et al., 2008; Hinrichs et al., 2009a) have been applied to the Alzheimer’s Disease Neuroimaging Initiative  
76 (ADNI) dataset, (<http://www.loni.ucla.edu/ADNI/Data/>) (Mueller et al., 2005) consisting of a large set of  
77 Magnetic Resonance (MR) and (18-fluorodeoxyglucose Positron Emission Tomography) FDG-PET images,  
78 giving accuracy measures similar to those reported in (Fan et al., 2008b,a; Davatzikos et al., 2008a,b). **In  
79 (Hinrichs et al., 2009a), we proposed a combination of  $\ell_1$  sparsity and spatial smoothness bias, implemented  
80 via augmentation of the linear program used in training. The spatial bias lead to an increase in accuracy, and  
81 made the resulting images more interpretable.** Steady increases in the levels of accuracy on this problem,  
82 *i.e.*, separating AD subjects from controls, have lead some researchers in the field to move towards the more  
83 challenging problem of making similar classifications on MCI subjects, with the expectation of extending  
84 such methods for identifying signs of the disease in its earlier stages. We provide a brief review of some  
85 preliminary efforts in this direction next.

86 Several recent studies (Schroeter et al., 2009; deToledo-Morrell et al., 2004; Dickerson et al., 2001; Hua  
87 et al., 2008) have shown that certain markers are significantly associated with conversion from MCI to  
88 AD. In (deToledo-Morrell et al., 2004; Dickerson et al., 2001), the authors show that traced volumes of

89 the hippocampus and entorhinal cortex show significant group-level differences between converting and non-  
90 converting MCI subjects. We note that these studies show (in a *post-hoc* manner) that certain brain regions  
91 are correlated with AD histopathology; what we seek to do instead is to evaluate such markers in terms of  
92 their ability to classify novel examples. In (Hua et al., 2008) a large number of ADNI subjects were tracked  
93 longitudinally using Tensor-Based Morphometry (TBM). The authors compared conversion from MCI to AD  
94 over 1 year with atrophy in various regions, but a discussion of the predictive accuracy results was relatively  
95 limited (*i.e.*, included  $p$ -values of 0.02 between converters and non-converters). In (Davatzikos et al., 2009),  
96 the authors applied statistical techniques to both ADNI and BLSA subjects (Shock et al., 1984). A classifier  
97 was trained using ADNI subjects, and applied to MCI and control subjects (in the BLSA cohort) to provide a  
98 SPARE-AD disease marker. This procedure could successfully separate MCI and control subjects with high  
99 confidence (AUC of 0.885), and it was demonstrated that the MCI group had a larger increase in SPARE-AD  
100 scores longitudinally. However, the main focus in (Davatzikos et al., 2009) was *not* on predicting which MCI  
101 subjects would progress to AD, but rather on finding a marker for MCI itself. In (Querbes et al., 2009),  
102 cortical thickness measures were used on a large set of ADNI subjects to characterize disease progression in  
103 AD and MCI subjects. Freely available tools (FreeSurfer) were used to calculate cortical thickness values at  
104 points on the surface of each subject’s brain (after warping to MNI template space) and then the thickness  
105 measures were agglomerated into 22 Regions of Interest (ROI), which the authors used as features (*i.e.*,  
106 covariates) in a logistic regression framework. Using age as a covariate, a set of AD and control subjects  
107 were used to train a logistic regression classifier for each subject, yielding a Normalized Thickness Index  
108 (NTI). It was found that this NTI was able to give 85% accuracy in separating AD subjects vs. controls,  
109 and had 73% accuracy (0.76 AUC) in predicting which MCI subjects would progress to full AD within 3  
110 years. The latter objective is of special interest in the context of the techniques presented in this paper.

111 A common trend in the studies mentioned above is their focus on using a single scanning modality and  
112 processing pipeline. For instance, in a recent study (Schroeter et al., 2009), the authors surveyed 62 original  
113 research papers in a meta-analysis aimed at identifying which brain regions might make the most useful  
114 markers of AD-related atrophy, in a variety of different scanning modalities. A fundamental assumption is  
115 that the studies use only one scanning modality and analysis method in isolation, rather than combining the  
116 several available modalities into a single disease marker. However, each scanning modality and processing  
117 method can reveal information about different aspects of the underlying pathology. For instance, structural  
118 MR images may reveal patterns of gray matter atrophy, while FDG-PET images may reveal reduced glucose  
119 metabolism (Ishii et al., 2005), PIB imaging highlights the level of amyloid burden in brain tissue (Klunk  
120 et al., 2004), and SPECT imaging can allow an examination of cerebral blood flow (Ramírez et al., 2009);  
121 similarly, Voxel-Based Morphometry (VBM) shows gray matter density at baseline, while Tensor-Based  
122 Morphometry (TBM) shows longitudinal patterns of change (Hua et al., 2008). Another important issue  
123 one must consider is that as new types of biologically relevant imaging modalities become available, (*e.g.*,  
124 new tracers for use in PET scanners, or new pulse sequences in MRI scanners), it is desirable for the  
125 diagnostic process to incorporate such advances seamlessly. Further, since AD pathology is known to be  
126 heterogeneous, (Thompson et al., 2001) it may be advantageous to include multiple scanning modalities in  
127 a single classification framework. Indeed, a wide variety of markers may be available, and it is desirable to  
128 make the best use of *all* such information in a predictive setting. The main difficulty is that as the number  
129 of available input features grows, many machine learning algorithms may lose their ability to generalize  
130 to unseen examples, due to the disparity between the sample size and the increased dimensionality. To  
131 address this problem, we propose to employ a recent development in the machine learning literature, called  
132 Multi-Kernel Learning (MKL), which is designed to deal with multiple data sources while controlling model  
133 complexity. We have evaluated this method’s performance on subjects from the ADNI data set, and report  
134 these results below. We have also applied the multi-modal classifier to MCI subjects, showing a promising  
135 ability to predict which subjects will convert from MCI to full AD in the ADNI sample.

136 The principal **contributions** of this paper are: **(1)** We propose a new application of Multi-Kernel Learn-  
137 ing (MKL) to the task of classifying AD, MCI, and control subjects, which permits seamless incorporation  
138 of tens of imaging modalities, clinical measures, and cognitive status markers into a single predictive frame-  
139 work. The main ideas behind MKL are presented in Section 2.2; **(2)** We have conducted an extensive set  
140 of experiments using ADNI subjects, aimed at providing a rigorous evaluation of the method’s ability to  
141 predict disease progression under conditions designed to match a clinical setting. We present these results  
142 in Section 4; **(3)** We employ our method to produce a Multi-Modality Disease Marker (MMDM) for MCI

143 subjects, and present an analysis of its predictive value on rates of conversion from MCI to AD in Section  
 144 4.3. A discussion of our results is given in Section 5. <sup>1</sup>

## 145 2 Algorithm

### 146 2.1 Support Vector Classification

147 In the following section, we present a brief overview of Support Vector Machines, (Cortes and Vapnik,  
 148 1995) illustrate the connection to Multi-Kernel Learning, and how this relates to the problem of disease  
 149 classification from multiple modalities.

150 Machine learning methods are designed to find a classifier (*i.e.*, function) that correctly (or maximally)  
 151 classifies a set of  $n$  training examples (*i.e.*, where class labels are known), while simultaneously satisfying  
 152 some other form of *inductive bias* which will allow the algorithm to generalize, *i.e.*, correctly label future  
 153 examples. Given a collection of points in a high dimensional space, SVM frameworks output a decision  
 154 function separating classes (in a maximum margin sense) in that space; the ‘bias’ here is toward selecting  
 155 functions with large margins. A linear decision boundary describes a *separating hyper-plane* – parameterized  
 156 by a weight vector  $\mathbf{w}$ , and an offset  $b$ . Classifying a new example  $\mathbf{x}$  involves taking the inner product between  
 157  $\mathbf{x}$  and  $\mathbf{w}$  plus the offset  $b$ ; the sign of this quantity indicates which side of the hyperplane  $\mathbf{x}$  falls on (*i.e.*, its  
 158 predicted class). In order to find the classifier, SVMs try not only to assign correct labels to each training  
 159 example by placing them on the correct side of the hyperplane, but also attempt to place them some distance  
 160 away. The measure of this distance is controlled by  $\|\mathbf{w}\|_2$ , or  $\ell_2$ -norm of  $\mathbf{w}$ . Thus, by rewarding the algorithm  
 161 for reducing the magnitude of  $\mathbf{w}$ , classifiers that correctly label the data (*and* have the widest margin) are  
 162 selected, see (Schoelkopf and Smola, 2002) for details. SVMs choose an optimal classifier by optimizing the  
 163 following primal/dual problem, whose solution  $\mathbf{w}$  gives the separating hyperplane:

$$\begin{array}{ll}
 \text{(primal)} & \text{(dual)} \\
 \min_{\mathbf{w}, \xi} \frac{\|\mathbf{w}\|_2}{2} + C \sum_i \xi_i & \max_{\alpha} \sum_i \alpha_i - \sum_{i,j} \alpha_i \alpha_j y_i y_j \underbrace{x_i^T x_j}_{\text{kernel}} \quad (2) \\
 \text{s.t. } y_i (\mathbf{w}^T x_i + b) \geq 1 - \xi_i \quad \forall i & \text{s.t. } 0 \leq \alpha_i \leq C \quad \forall i \\
 \xi_i \geq 0 \quad \forall i & \sum_i y_i \alpha_i = 0 \quad \forall i
 \end{array}$$

165 In the primal problem (1), the *slack variables*  $\xi$  implement a *soft margin* objective. That is, for each  
 166 example  $i$  that is not placed more than unit distance from the separating hyperplane, the slack variable  
 167  $\xi_i$  takes the value of the remaining distance from example  $i$  to the margin, which is then penalized in the  
 168 objective.  $C$  is a constant parameter controlling the amount of emphasis on separating the data (if  $C$   
 169 is large,) vs. widening the margin (if  $C$  is small). Thus, the soft-margin objective allows for a trade-off  
 170 between perfectly classifying every example, and widening the margin. The bias term  $b$  allows for separating  
 171 hyperplanes ( $\mathbf{w}^T x + b$ ) which do not pass through the origin. Class labels for each example are given as  
 172  $y_i = \pm 1$ , so that  $y_i(\mathbf{w}^T x_i + b)$  will be positive iff  $\mathbf{w}^T x + b$  gives  $x_i$  the correct sign specified by  $y_i$ .

173 Note that the hyperplane parameters  $\mathbf{w}$  can be given as a linear combination of examples. It is a special  
 174 property of the SVM formulation that the dual variables <sup>2</sup>  $\alpha$  are exactly the coefficients of such a linear  
 175 combination, *i.e.*,  $\mathbf{w} = \sum_i \alpha_i y_i x_i$ . For typical settings of  $C$ , the support of  $\alpha$  will be sparse, giving rise to  
 176 the term ‘Support Vector Machine’.

177 Note that in the dual problem (2), the examples only occur as inner products  $\langle x_i, x_j \rangle$ . These inner  
 178 products can be captured in a single  $n \times n$  matrix called a Gram matrix or kernel matrix,  $\mathcal{K}$ ; see (Bishop,  
 179 2006). In practice,  $\mathcal{K}$  is specified by the user and expresses some notion of similarity between the examples –  
 180 that is, the magnitude of a kernel function of two examples expresses an inner product between corresponding

<sup>1</sup>A preliminary conference version of this paper appeared as (Hinrichs et al., 2009b).

<sup>2</sup>In linear and quadratic optimization, every primal problem has an associated dual problem; the optimal solution to one can be used to recover the optimal solution to the other.

181 points in an implicit Reproducing Kernel Hilbert Space  $\mathcal{H}$ . The translation from the original data space to  
 182  $\mathcal{H}$  is commonly denoted as  $\phi(x)$ ; when the kernel function is modified,<sup>3</sup> the kernel space  $\mathcal{H}$  and translation  
 183 function  $\phi(x)$  are correspondingly modified. The kernel function can also be calculated analytically – among  
 184 those commonly used are Linear, Polynomial, and Gaussian kernels. Briefly, a linear kernel function is simply  
 185 the inner product of two examples in the original data space; thus, unmodified SVMs use a linear kernel. A  
 186 polynomial kernel function is one in which each inner product is squared (or cubed etc.). Such kernels allow  
 187 for polynomial decision boundaries, rather than simple hyperplanes. Finally, Gaussian kernels are based on  
 188 the Euclidean distance between examples, by the formula

$$\exp\left(\frac{-\|x_i - x_j\|}{2\sigma}\right)$$

189 where  $\sigma$  is a bandwidth parameter and  $x_i$  and  $x_j$  may denote examples  $i$  and  $j$ . Gaussian kernel-based  
 190 SVMs can be thought of as training a Gaussian mixture model as the pattern classifier. If a modified kernel  
 191 function is used, corresponding to a non-linear transformation of the data, then the learned classifier is a  
 192 linear function (*i.e.*, hyperplane) in the kernel space  $\mathcal{H}$ . Such a function typically maps back to a non-linear  
 193 decision function in the original data space. A thorough treatment is given in (Bishop, 2006).

## 194 2.2 Multi-Kernel Pattern Classification

195 An extension of this idea is to combine many such functions of the data (*i.e.*, multiple kernels, each pertaining  
 196 to one modality for example, or to different parameterizations of the kernel function, or to different sets of  
 197 selected features), to create a single kernel matrix from which a better classifier can be learnt. Multi-kernel  
 198 learning (MKL) (Lanckriet et al., 2004; Sonnenburg et al., 2006; Rakotomamonjy et al., 2008; Gehler and  
 199 Nowozin, 2009; Mukherjee et al., 2010) formalizes this idea. This is achieved by adding a set of optimization  
 200 variables called *subkernel weights* which are coefficients in a linear combination of kernels. The subkernel  
 201 weights are chosen so that the resulting linear combination of kernel matrices (another kernel matrix) yields  
 202 the best margin and separation on the training set, with additional regularization to reduce the chances of  
 203 overfitting the data due to the increase in the degrees of freedom of the model.

$$\begin{aligned} \min_{\mathbf{w}_k, \xi, \beta, b} & \left( \sum_k \frac{\|\mathbf{w}_k\|_2}{\beta} \right)^2 + C \sum_i \xi_i + \|\beta_k\|_2^2 & (3) \\ \text{s.t. } & y_i \left( \sum_k \mathbf{w}_k^T \phi_k(x_i) + b \right) \geq 1 - \xi_i \quad \forall i \end{aligned}$$

204 Here,  $\beta_k$  is the subkernel weight of the  $k$ -th kernel, and  $\mathbf{w}_k$  is the set of weights for the  $k$ -th feature space,  
 205 while  $\xi_i$  is a *slack variable* as described above. Regularization of the subkernel weights is accomplished by  
 206 penalizing the squared 2-norm of  $\beta$  in the objective. Thus, in addition to minimizing the magnitude of each  
 207 set of weights, the MKL algorithm also tries to minimize the magnitude of the subkernel weight vector. Thus  
 208 as  $\beta_k$  grows larger, the corresponding  $\mathbf{w}_k$  is penalized less, and therefore tends to have a larger contribution  
 209 to the final classifier. The combined classifier is defined as  $f(x) = \sum_k \mathbf{w}_k^T \phi_k(x) + b$ . Thus, the implicit  
 210 kernel function is equal to  $\sum_k \beta_k \phi_k(x_i)^T \phi_k(x_j)$ . In the context of our application, it is helpful to think  
 211 of the various kernel matrices as being derived from different sources of data (e.g., different modalities),  
 212 different choice of kernel function or parameters, (*e.g.*, bandwidth parameter in a Gaussian kernel function,)   
 213 or a different set of features. Their assigned weights can then be interpreted as their relative influence in  
 214 learning a good classifier (*i.e.*, discriminative ability). Because there is a natural mechanism to control the  
 215 greater complexity resulting from the increased dimensionality of multi-modality data, we believe that MKL  
 216 is a preferable option rather than simply ‘concatenating’ all features together and using a regular SVM. Our  
 217 proposed method then, is to calculate various kernel matrices from each available input modality – including  
 218 brain images, cognitive scores and other characteristics, such as CSF assays or APOE genotype, and use  
 219 MKL to train an optimal combined kernel and classifier.

<sup>3</sup>Any such modification must preserve the positive-definite property of the original kernel function.

220 Note that in the term  $\|\beta_k\|_2^2$  the subkernel weights are penalized according to the Euclidean, or 2-norm. <sup>4</sup>  
221 A recent focus in MKL research has been to generalize this formulation to include other norms (Kloft et al.,  
222 2010), having different effects on the sparsity of the resulting vector of subkernel weights. For instance, the  
223 1-norm is a sparsity inducing norm, while the 2-norm is not; norms between 1 and 2 allow a trade-off of  
224 emphasis between sparse and non-sparse solutions. When combining multiple imaging modalities for AD  
225 classification, it is preferable not to encourage sparsity, as the algorithm will be very likely to completely  
226 ignore some modalities.

## 227 3 Experimental Setup

### 228 3.1 Data

229 Data used in the evaluations of our algorithm were taken from the Alzheimer’s Disease Neuroimaging Initia-  
230 tive (ADNI) database ([www.loni.ucla.edu/ADNI](http://www.loni.ucla.edu/ADNI)). The ADNI was launched in 2003 by the National Institute  
231 on Aging (NIA), the National Institute of Biomedical Imaging and Bioengineering (NIBIB), the Food and  
232 Drug Administration (FDA), private pharmaceutical companies and non-profit organizations, as a \$60 mil-  
233 lion, 5-year public-private partnership. The primary goal of ADNI has been to test whether serial magnetic  
234 resonance imaging (MRI), positron emission tomography (PET), other biological markers, and clinical and  
235 neuropsychological assessment can be combined to measure the progression of mild cognitive impairment  
236 (MCI) and early Alzheimers disease (AD). Determination of sensitive and specific markers of very early AD  
237 progression is intended to aid researchers and clinicians to develop new treatments and monitor their effec-  
238 tiveness, as well as lessen the time and cost of clinical trials. The Principal Investigator of this initiative is  
239 Michael W. Weiner, M.D., VA Medical Center and University of California San Francisco. ADNI is the result  
240 of efforts of many co-investigators from a broad range of academic institutions and private corporations, and  
241 subjects have been recruited from over 50 sites across the U.S. and Canada. The initial goal of ADNI was to  
242 recruit 800 adults, ages 55 to 90, to participate in the research approximately 200 cognitively normal older  
243 individuals to be followed for 3 years, 400 people with MCI to be followed for 3 years, and 200 people with  
244 early AD to be followed for 2 years.

245 Our data consisted of ADNI subjects for whom both MR and FDG-PET scans roughly 24 months apart  
246 were available (as of October 2009). For quality control purposes, several (16) subjects were removed due  
247 to motion artifacts (MR), reconstruction artifacts (FDG-PET) or other problems visible to an expert. All  
248 such evaluations were made *before* any classification experiments were conducted, so as not to unfairly bias  
249 the experimental results. Finally, we had data for 233 subjects (48 AD, 66 healthy controls, and 119 MCI  
250 subjects). Demographic data are shown in Table 1.

### 251 3.2 Preliminary Image-processing

252 In order to apply SVM and MKL methods to imaging data, it is necessary to extract features which are  
253 common to all subjects. Using standard voxel-based morphometry methods, as described below, we warped  
254 the scans into a common template space, and used voxel intensities as features. That is, after extracting  
255 foreground voxels, (*i.e.*, those corresponding to brain tissue,) each subject can then be treated as a vector  
256 of fixed length.

257 **T1-weighted MR images.** Cross-sectional image processing of the baseline T1-weighted images was  
258 first performed using Voxel-Based Morphometry (VBM) toolbox in Statistical Parametric Mapping software  
259 (SPM, <http://www.fil.ion.ucl.ac.uk/spm>). **The ADNI study provides repeated acquisitions of the MR scans,**  
260 **which we utilized by first performing an affine warp between duplicates, and then averaging them in order**  
261 **to boost the signal/noise ratio.** We then segmented the original anatomical MR images into gray matter  
262 (GM), white matter (WM), and cerebrospinal fluid (CSF) segments. Then by using the “DARTEL Tools”  
263 facility in SPM5, a study-cohort customized template was calculated based on all subjects’ baseline MR  
264 images with the registration results as well as all relevant flow fields (representing the transformations). All  
265 individual MR scans were subsequently warped to this new template. Modulated GM and WM segments  
266 were produced in the DARTEL template space, using both the original scans (Ashburner, 2007). Finally,  
267 the normalized maps were smoothed using an 8 mm isotropic Gaussian kernel to optimize signal to noise and

<sup>4</sup>In general, the p-norm of a space  $\mathcal{X}$  is given as  $\|(\mathbf{x})\|_p = (\sum_i |x_i|^p)^{1/p}$ , for  $x \in \mathcal{X}$ .

268 facilitate comparison across participants. Analysis of gray matter volume employed an absolute threshold  
269 masking of 0.1 to minimize the inclusion of the white matter in analysis. Longitudinal MR image processing  
270 of baseline and 24-Month MR scans was performed with a tensor-based morphometry (TBM) approach in  
271 SPM5. We first co-registered the baseline and follow-up scans with rigid body affine transformation, and  
272 applied bias correction and intensity normalization to make both images comparable. Pre-processing TBM  
273 procedures are described in detail in a previous article (Kipps et al., 2005). Briefly, a deformation field was  
274 used to warp the corrected late image to match the early one within subject (Ashburner and Friston, 2000).  
275 The amount of volume change was quantified by taking the determinant of the gradient of deformation at a  
276 single-voxel level (*i.e.*, Jacobian determinant). Each subject’s Jacobian determinant map was normalized to  
277 the cohort-specific DARTEL template and smoothed using a 12 mm isotropic Gaussian kernel.

278 **FDG-PET images.** All FDG-PET images were first co-registered to each individual’s baseline MR-T1  
279 images and subsequently warped to the cohort-specific DARTEL template (see above). A mask of the Pons  
280 was manually drawn in the DARTEL template as the reference region. All of the normalized FDG-PET  
281 images were scaled to each individual’s Pons average FDG uptake value and smoothed with a 12 mm isotropic  
282 Gaussian kernel.

283 **Other biological and neurological data.** In addition to MR and FDG-PET images, other biological  
284 measures and cognitive status measures are provided by ADNI for some subjects. These include CSF assays  
285 for certain compounds thought to be involved in neurodegeneration, such as AB1-42, Total Tau, and P-tau  
286 181; NeuroPsychological Status Exam scores (NPSEs); and APOE genotype data. The complete list of  
287 biological measures, and their availability in the study population is shown in Tables 2 and 3.

### 288 3.3 Experimental Methodology

289 We performed two sets of classification experiments: **(1)** We first performed *multi-modal* classification ex-  
290 periments for separating AD and control subjects using baseline and longitudinal imaging data, (MR and  
291 FDG-PET), and other available cognitive / biological measures (CSF assays, NeuroPsychological Status  
292 Exams (NPSE), and APOE genotype). For comparison, we also present single-kernel experiments for each  
293 data modality (except APOE, since APOE genotype alone is not sufficient to diagnose AD), and on an SVM  
294 trained on the sum of all kernels, (or equivalently, the concatenation of all feature vectors). **(2)** Finally,  
295 we trained a classifier on the entire set of AD and control subjects and then applied it to the MCI popu-  
296 lation, giving a *Multi-Modality Disease Marker (MMDM)*. We compared this marker with NPSEs taken at  
297 24 months, and examined its utility in predicting which MCI subjects would progress to AD, as opposed to  
298 remaining stable as MCI. Note that this is different from separating MCI subjects from AD/controls.

299 **Kernel matrices** Kernel matrices used in our experiments were computed using a varying number  
300 of voxel-wise features, (*i.e.*, intensity values at each voxel,) and kernel functions *i.e.*, linear, quadratic and  
301 Gaussian, for each imaging modality. For each fold, voxels were ranked by *t*-statistic between AD and control  
302 training subjects. That is, each voxel’s intensity value can be thought of as a random variable, upon which  
303 we performed a *t*-test, and ranked the features by the resulting p-values. Separate kernels were computed  
304 using the top 250,000, 150,000, 100,000, 65,000, 25,000, 10,000, 5000 and 2000 features, respectively. These  
305 sets of features were chosen beforehand so as to give a reasonable coverage of the range of features available,  
306 while allowing the algorithm to choose a linear combination that leads to a discriminative kernel. In addition  
307 to performing an implicit feature selection step, this allows us to evaluate the MKL algorithm’s ability to  
308 integrate tens to hundreds of kernels, as in the case when many more modalities are available. For each set  
309 of features, we constructed linear, quadratic, and Gaussian kernels, using a bandwidth parameter of 2 times  
310 the number of features for the Gaussian kernel. The Gaussian kernel bandwidth parameter should be chosen  
311 to be within the same order of magnitude as the majority of pairwise distances. Thus, when voxel-wise  
312 intensity values fall in the range  $[0, 1]$ , a common choice for the bandwidth parameter is a small number  
313 times the number of features. By this process, we obtained 24 separate kernel matrices for each imaging  
314 modality. For non-imaging modalities, *i.e.*, CSF assays, NPSEs, and APOE genotype, all features were used,  
315 giving three kernels per modality. The biological measures used are shown in Table 2. Because only a subset  
316 of subjects had such measures available, we used zero values for those who did not. **This means that kernel**  
317 **matrices had zero values where such data were missing, and therefore added nothing to the classification on**  
318 **those subjects. We chose a conservative approach to this problem, meaning that results can only improve if**  
319 **a statistical interpolation method were to be introduced.** For computing the MMDM for MCI subjects, all  
320 AD and CN subjects were used both in feature selection and training.

321 Before training a classifier using the kernels constructed as described above, it is necessary to perform some  
322 normalization; consider that the vector  $\mathbf{w}$  which defines the separating hyperplane is a linear combination  
323 of examples. If the average magnitude of examples as implicitly represented by one kernel is orders of  
324 magnitude larger than that of another kernel, then for the same subkernel weights, one kernel will have a far  
325 greater contribution to  $\mathbf{w}$ . In order to ensure that this is not the case, we adopted a standard approach to  
326 kernel normalization. The first step is to divide each kernel by the largest entry, so that all entries are in the  
327 range  $[0, 1]$ . Second, we re-centered the points in each kernel space by subtracting row and column mean  
328 values, and then dividing by the trace. See Bakir et al. (2007) for details. As a consequence of normalizing  
329 the kernels, the  $C$  parameter which controls the regularization trade-off can be set to a small integer. We  
330 therefore set  $C = 10$ ; no fine tuning or model selection was necessary.

331 Recall that when longitudinal data are available, there is more than one way to perform spatial normal-  
332 ization of scans, and we treat them as different imaging modalities, because we expect different types of  
333 information to be revealed by each. From MR images, we have both baseline VBM, and TBM modalities;  
334 in FDG-PET we have baseline and 24 month scans, as well as the voxel-wise difference and ratio between  
335 scans at different time points. Kernels based on the longitudinal voxel-wise difference and ratio in FDG-PET  
336 images were found to have poor performance relative to the raw FDG-PET values (60% – 70% accuracy),  
337 and we did not make further use of them in our experiments.

338 **ROC curves** We also computed Receiver Operator Characteristic curves (ROCs) for each set of ex-  
339 periments. Briefly, while a classification algorithm must output a  $\pm 1$  group label, our algorithm can also  
340 output a ‘confidence’ level for each test subject which in this case is the signed output of the classifier . By  
341 ordering the confidence levels of the entire study population, and calculating a True Positive Rate (TPR or  
342 sensitivity) and False Positive Rate (FPR or 1 - specificity) for each level, an ROC curve qualitatively shows  
343 not only how many examples are misclassified, but provides a sense of how the classifier’s confidence relates  
344 to its correctness.

345 **Cross-validated classification** For the first set of experiments, we performed AD vs. control classi-  
346 fication experiments using 30 realizations of 10-fold cross-validation. That is, in each realization the study  
347 population was randomly divided into ten separate groups, or folds. Each fold was used as a “test” set,  
348 while the remaining data was used as a “training” set. Therefore, the algorithm was evaluated on AD and  
349 control examples which were unseen during the training process, while permitting us to use the entire dataset  
350 effectively. Various accuracy measures, such as test-set accuracy (% of test examples properly labeled as AD  
351 or control,) sensitivity, (% of AD cases labeled as such) and specificity (% of controls labeled as such), and  
352 area under ROC curves were computed by averaging over all 30 realizations. Using this methodology, we first  
353 evaluated each kernel function on its own, in an SVM framework. We then evaluated each modality in an  
354 MKL framework, by combining different kernel functions, all derived from the same modality and features.  
355 Finally, we combined all imaging modalities into a multi-modality MKL classification framework. We did  
356 the same for cognitive scores and biological measures, allowing for a comparison between different types of  
357 subject data in terms of their ability to identify signs of AD.

358 **Comparison of subkernel weight vector regularization norms** Another interesting area of investi-  
359 gation is on the effect of different MKL norm regularizers, especially with regard to sparsity of the resulting  
360 classifier. Sparsity is often advantageous in the presence of non-informative or error-prone kernels, however  
361 an overly sparse combination can discard useful information, leading to a sub-optimal classifier. Thus, it  
362 is important to understand this trade-off. Using the cross-validation setup described above, we compared  
363 different subkernel norm regularizers, (1, 1.25, 1.5, 1.75, and 2), using all available kernel types, as shown  
364 in Tables 2 and 3. In order to demonstrate MKL’s ability to combine fundamentally different sources of  
365 information, we also constructed additional kernels using subject age, APOE genotype, years of education,  
366 and geriatric depression scale as features. We expect that some of these additional kernels may or may not  
367 be as useful to the learning algorithm, so as to allow a meaningful assessment of the usefulness of applying  
368 sparsity in the kernel norm. For baseline comparison we trained an SVM on the sum of all kernels, which is  
369 equivalent to simply concatenating all feature vectors, by definition of the inner product of vectors.

370 **MMDMs** Our next set of experiments were conducted to evaluate the ability of imaging-based markers  
371 to predict which subjects would convert from MCI to AD. In order to do this, we first trained an MKL  
372 classifier using all 114 AD and CN subjects, and then applied it to all 119 MCI subjects, giving an MMDM  
373 measure. This procedure was repeatedly performed using (a) imaging-based, (b) cognitive marker-based,  
374 and (c) biological measure-based kernels, so as to evaluate each type of data separately, and facilitated a

375 better comparison among them. We also differentiated between baseline and longitudinal data.

376 To quantify the predictive value of the MMDMs, we separated the MCI subjects into three groups –  
 377 those who had progressed to AD after three years, those who remained stable, and those who reverted to  
 378 normal status – and calculated p-values of group differences using a *t*-test. We also computed ROC curves  
 379 to quantitatively measure the degree of differentiation between the MCI groups as given by different types  
 380 of biological measures. There are two ways to compute such ROCs: based on the differentiation between  
 381 progressing and reverting MCI subjects, ignoring the stable MCI subjects; and based on the differentiation  
 382 between progressing and non-progressing MCI subjects. In the former case, we treat stable MCI subjects  
 383 as though their final status is not yet known, and thus the task is to predict whether a given subject will  
 384 eventually revert, or progress. For our analysis, we calculated both kinds of ROC curves, and present results  
 385 below.

386 **Implementation** Our validation experiments and analysis framework were implemented in Matlab using  
 387 an interface to the Shogun toolbox (Sonnenburg et al., 2006) (<http://www.shogun-toolbox.org>). The  
 388 source code for this project and supplemental information will be made available at [http://pages.cs.wisc.edu/~hinrichs/MKL\\_ADNI](http://pages.cs.wisc.edu/~hinrichs/MKL_ADNI) [upon publication].  
 389

**TABLE 1** Study population demographics

	controls (mean)	controls (s.d.)	MCI (mean)	MCI (s.d.)	AD (mean)	AD (s.d.)
Age at baseline	76.2	4.59	75.1	7.44	76.6	6.28
Gender(M/F)	40/26	–	79/40	–	25/23	–
APOE carriers	17	–	63	–	37	–
MMSE at Baseline	29.17	0.85	27.18	1.64	23.50	1.92
MMSE at 24 months	28.67	3.73	25.54	4.84	18.98	6.60
ADAS at baseline	9.94	4.27	17.26	6.13	28.27	9.80
Years of Education	16.15	3.02	15.73	2.82	14.60	3.17
Geriatric Depression	0.97	1.35	1.40	1.28	1.71	1.47

Table 1: Demographic and neuropsychological characteristics of the study population.

**TABLE 2** Biological measures data used in kernel functions

Type	Subjects available
Tau	130
Amyloid-Beta 142	130
P-Tau 181P	130
T-Tau	130
APOE Genotype	233

Table 2: Non-imaging biological measures used to construct kernels for experiments. Cerebro-Spinal Fluid (CSF) assays and APOE genotype data were utilized.

## 390 4 Results and Analysis

391 We present here the results of our experiments on the ADNI data described in Section 3, and an analysis of  
 392 the MKL algorithm in the context of MCI progression.

### 393 4.1 Separating AD subjects and Controls

394 As a first step, we separately evaluated the kernels produced by each modality by comparing their perfor-  
 395 mance at classifying AD vs. control subjects using an MKL norm of 2.0, so as not to discard any useful

**TABLE 3** Cognitive markers used in kernel functions

Cognitive measure	Subjects available
Rey auditory / verbal 1-5 scores	233
Rey auditory delayed recall scores	233
Category Fluency scores	233
Trail-making A & B	233
Digit-span scores	233
Boston Naming scores	233
ANART errors	233

Table 3: Non-imaging cognitive markers used to construct kernels for experiments.

396 information. Results of these experiments are shown in Figure 1. Note that the color scale is the same  
 397 between all figures.

398 Our first set of multi-kernel experiments also focused on whether the algorithm could learn to separate  
 399 AD subjects from controls. Our experimental method was to use 10-fold cross-validation repeated 30 times,  
 400 using kernel matrices computed as described in 3.3. Accuracy, sensitivity, and specificity results are shown  
 401 in Table 4. In order to compare the efficacy of imaging-based disease markers with other biological measures,  
 402 we performed experiments **(1)** using only image-derived data, **(2)** using other biological measures, **(3)** using  
 403 only NPSEs, and finally using all available data modalities.

404 Note that the accuracy achieved using imaging-based MMDMs is nearly as good as that achieved using  
 405 NPSEs. **We believe this is promising, because NPSEs should be expected to perform better than imaging**  
 406 **modalities when AD-related cognitive decline is present, even if the NPSEs were not used in making the**  
 407 **diagnosis. This is because AD is currently diagnosed according to the patient’s cognitive status, and while**  
 408 **the NPSEs we utilized are *not* the same as those used in making a clinical diagnosis, they are nonetheless**  
 409 **markers of detectable decline in cognition, and as such are not directly comparable to imaging-based markers.**  
 410 **Rather, we include these experiments only to facilitate indirect comparison.** Thus, for the imaging-based  
 411 markers to be nearly as effective is quite promising.

412 The areas under each ROC curve (another measure of classification performance) are provided in Table  
 413 4. In terms of area under ROC curve, all modalities performed about as well as other accuracy measures  
 414 would suggest. Again, we note that imaging modalities and cognitive scores performed very similarly under  
 415 this measure.

416

**TABLE 4** Accuracy results of validation experiments using 2-norm MKL

Modalities used	Accuracy	Sensitivity	Specificity	Area under ROC
Imaging modalities	0.876	0.789	<b>0.938</b>	0.944
Biological measures	0.704	0.581	0.794	0.767
Cognitive scores	<b>0.912</b>	<b>0.892</b>	0.926	<b>0.983</b>
All modalities	0.924	0.867	0.966	0.977

Table 4: Comparison of 2-norm MKL with different types of input data modalities.

417 In order to compare the effect of subkernel weight norms, we repeated the above experiments using all  
 418 kernels and modalities available and MKL norms in the range of (1, 1.25, 1.5, 1.75, 2). These results are  
 419 shown in Table 5. Note that among the MKL norms, accuracy increases slightly with MKL norm up to the  
 420 point where sparsity is no longer strongly encouraged (at about 1.5), suggesting that overly sparse MKL norm  
 421 regularizers do indeed lose information. We also note that the SVM’s performance suffered significantly.

422 **When using a 1-norm, out of the 72 available kernels, only 4 had non-zero weights: one TBM Gaussian**  
 423 **kernel using 10,000 features, two VBM kernels, (one linear with 10,000 features, one quadratic with 25,000),**

FIGURE 1

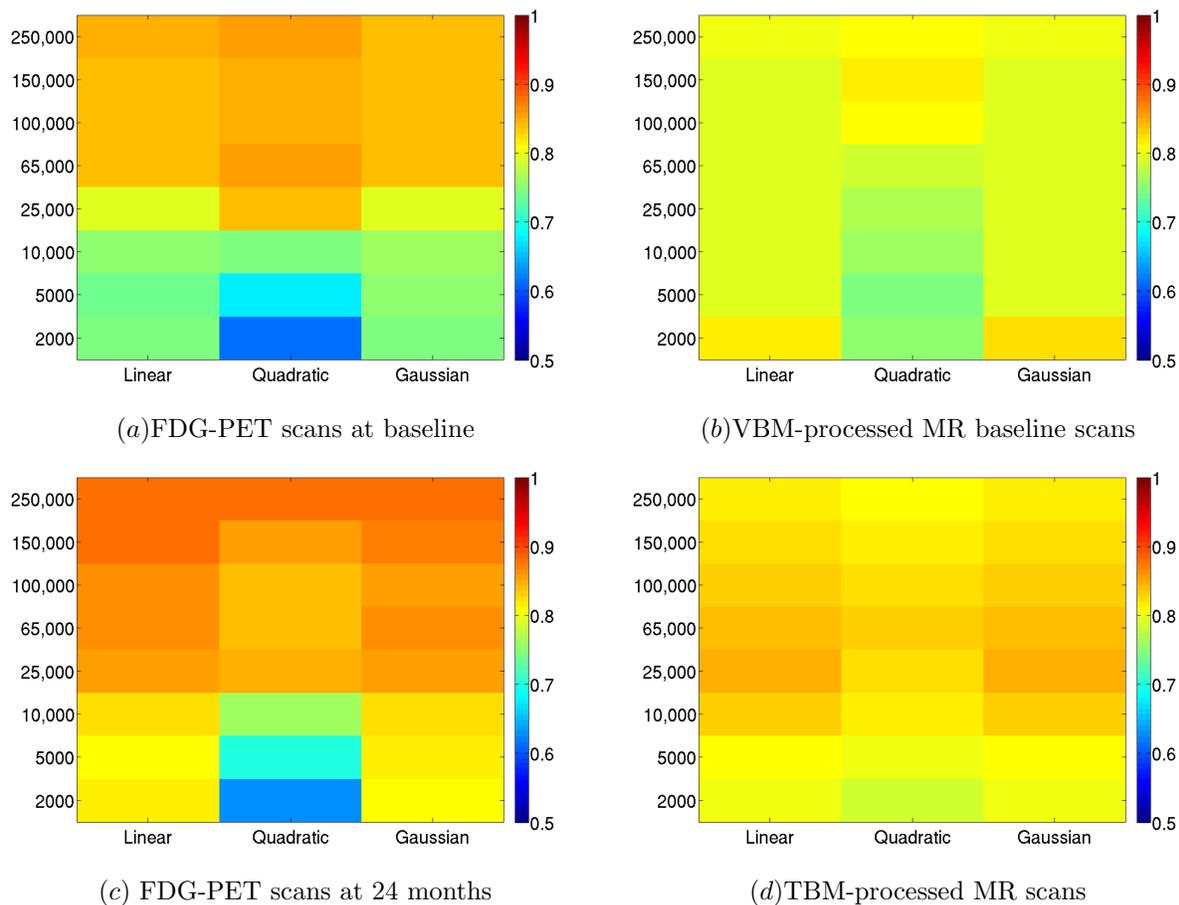


Figure 1: Accuracies of single-kernel, single-modality methods. Color represents classification accuracy on unseen test data, ranging from blue (lowest, 50% accuracy,) to red (highest, 100% accuracy). The modalities used are, (a) FDG-PET scans at baseline, (b) VBM-processed MR baseline scans, (c) FDG-PET scans at 24 months, and (d) TBM-processed MR scans.

TABLE 5 Comparison of different MKL norms with the SVM trained on concatenated-features

MKL norm used	Accuracy	Sensitivity	Specificity	Area under ROC
1.0	0.914	0.867	0.949	0.977
1.25	0.916	0.865	0.954	0.980
1.5	0.921	0.874	0.956	0.982
1.75	0.923	0.872	0.961	0.982
2.0	0.922	0.870	0.959	0.981
SVM (concatenated features)	0.882	0.844	0.910	0.970

Table 5: Comparison of different MKL norms in the presence of uninformative kernels, and an SVM trained on a concatenation of all features for comparison.

424 none from the baseline FDG-PET scans, and one linear kernel with 2,000 features. In contrast, the subkernel  
 425 weights chosen when using an MKL norm of 2 were *all* non-zero, and are shown in Figure 2. This means  
 426 that in the context of AD classification, different modalities (and different representations of information

427 from those modalities) contributed to in varying proportions to yield a discriminative classifier. It is perhaps  
428 interesting to note that most of the weight was placed on the VBM kernels, followed by the TBM and  
429 FDG-PET kernels.

## 430 4.2 Classifier brain regions

431 An important component of the evaluation of our method is an analysis of the brain regions selected by  
432 the algorithm. That is, if the algorithm is only given linear kernels from brain images, then the decision  
433 boundary itself can be interpreted as a set of voxel weights, using the formula  $\mathbf{w}_m = \beta_m \sum_i \alpha_i \phi_m(\mathbf{x}_i)$  where  
434  $\phi_m(\mathbf{x})$  is the implicit (possibly non-linear) transform from the original data space to the kernel Hilbert  
435 space. An examination of these weights can reveal which brain regions were found to be most useful or  
436 discriminative (by the algorithm) in its predictions. Thus, the images of brain regions below are taken from  
437 the multi-modality classifier trained on all four imaging modalities used in our experiments, using *only* linear  
438 kernels. Note that from Figure 1, we can see that among the kernels derived from FDG-PET images, the  
439 most informative kernel used more than 65000 voxels, which implies that *classification strategies can benefit*  
440 *from using whole-brain images rather than examining small, localized brain regions, or ROIs* in FDG-PET  
441 imaging. The results are shown in Figures 3 – 6. Note that these weights were all calculated simultaneously  
442 in the MKL setting. These images can be interpreted as follows: image intensity in voxels showing a stronger  
443 red color contributes to a subject’s healthy (positive) diagnosis, while intensity in voxels showing a stronger  
444 blue color contributes to a subject’s diseased (negative) diagnosis, and intensity in yellow-, green- or cyan-  
445 colored voxels is essentially ignored. Note that these weights are purely relative, and thus have no applicable  
446 units. Each subject’s final score is thus the difference between the weighted average intensity in the red  
447 and orange regions and the blue and cyan regions. We interpret this as meaning that red-orange (positive  
448 weighted) regions are those in which image intensity is a prerequisite of healthy status. For blue-cyan  
449 (negative weighted) regions, the literal interpretation is that the algorithm found higher intensity among the  
450 AD group than in the controls.

451 In some cases, we observe that negative weights are assigned in regions where higher image intensity  
452 is usually associated with positive status. There are several possible explanations for this, such as image  
453 normalization artifacts which artificially boost the intensity of these regions in some AD subjects. For  
454 instance in FDG-PET images, image intensity was normalized using a map of the Pons, and thus irregularities  
455 in this region could produce artificially inflated intensities in the rest of the image. Another possibility is  
456 brought up by (Davatzikos et al., 2009), which is that in MR images of gray matter, periventricular white  
457 matter may be mis-segmented as gray-matter, due to certain types of vascular pathology. A third possibility  
458 is that there is a small set of subjects whose characteristics is heterotypical of their group, and thus induce  
459 negative weights in regions which would otherwise have positive weights. Evidence of such a group was  
460 found in (Hinrichs et al., 2009a). In order to examine this possibility we found a set of subjects (5 subjects  
461 based on baseline FDG-PET scans, and 4 subjects based on baseline MR scans) who had unusually strong  
462 intensity in regions which had been assigned negative weights, and re-trained the MKL classifier without  
463 them. The resulting classifier was nearly free of such anomalous negative weights, which strongly suggests  
464 that these negative weights are entirely the result of the influence of a small group of outlier subjects, (9 out  
465 of 114). We have investigated this issue briefly in our previous work. (Hinrichs et al., 2009a) The weights  
466 assigned by this classifier can be seen in Figure 7. It is important to note that these subjects were removed  
467 for visualization purposes only, and were still used in computing accuracy and other performance estimates,  
468 and in the MCI analyses described below.

469 In Fig. 3, we can see that heteromodal, frontal, parietal regions and temporal lobes are given negative  
470 weights. The posterior cingulate cortex, lateral parietal lobules (bilaterally) and pre-frontal midline struc-  
471 tures prerequisite of an indication of healthy status. The weights assigned to the FDG-PET scans taken at  
472 24 months show a similar pattern, and are shown in Figure 4.

473 Among the MR-based kernels, the most informative kernels (as measured in a single-kernel setting,)   
474 used 5000 to 25000 voxels, implying that smaller regions, can be used to identify signs of AD-related gray  
475 matter atrophy. Thus, we expect to see a similar pattern in the multi-modality setting. Using the same  
476 interpretation of color as above, we can see that in the baseline GM density images, (VBM) hippocampal  
477 and parahippocampal regions are highlighted more clearly, consistent with the single-modality results which  
478 indicated that a small number of voxels are most informative in this modality. In the TBM-based images,  
479 we see that the hippocampal regions and parahippocampal gyri are highlighted, as well as middle temporal

480 lobar structures bilaterally, indicating that longitudinal atrophy is concentrated in these regions, which is  
 481 again consistent with the single kernel results, (and prior literature), (Braak et al., 1999) in which the top  
 482 25000 voxels produced the most informative classifier.

**FIGURE 2**

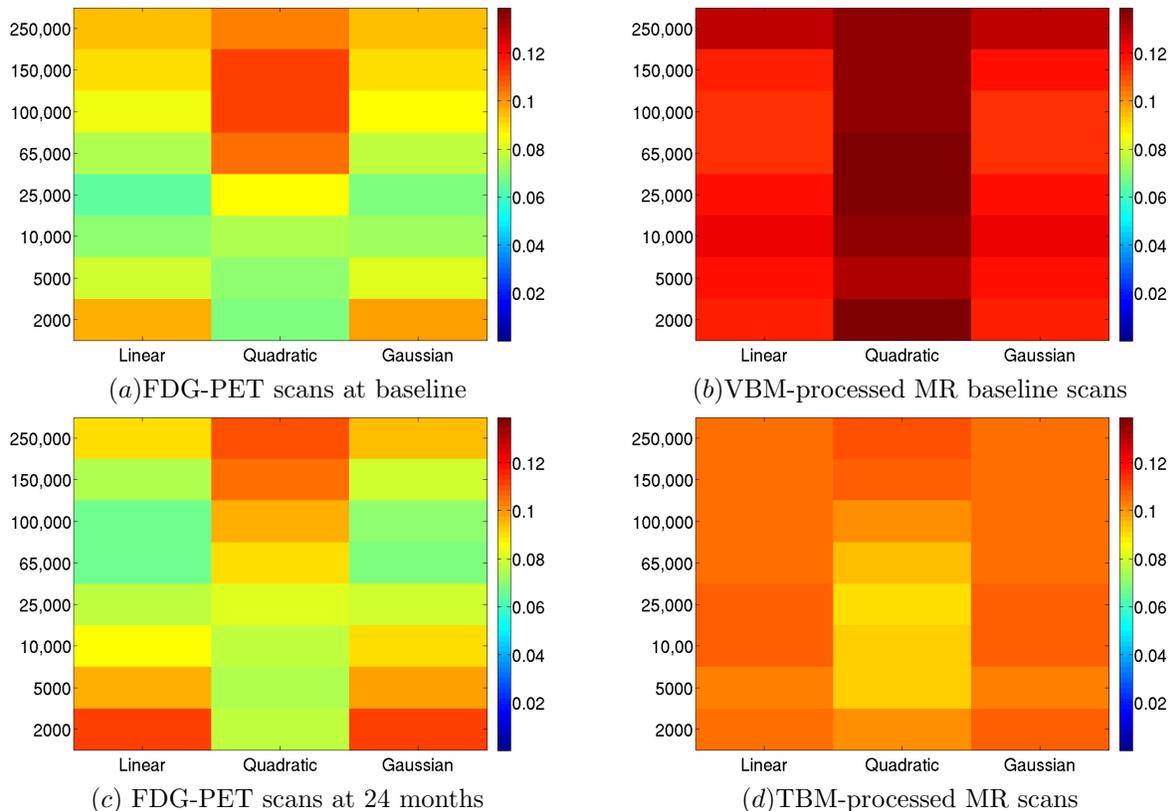


Figure 2: Subkernel weights ( $\beta$ ) chosen by the MKL algorithm with 2-norm regularization. Weights are relative, and have no applicable units. The modalities used are, (a) FDG-PET scans at baseline, (b) VBM-processed MR baseline scans, (c) FDG-PET scans at 24 months, and (d) TBM-processed MR scans.

483  
 484  
 485  
 486  
 487  
 488

### 489 4.3 Correlations and predictions on the MCI population

490 For the second set of experiments, which involved MCI subjects, we trained a classifier on the entire AD  
 491 and control population using MKL. This classifier was then applied to the MCI population, giving a Multi-  
 492 Modality Disease Marker (MMDM). Using this methodology, only AD and control subjects were used to  
 493 train the model, while MCI subjects were only used for evaluation, rather than other methodologies in which

FIGURE 3

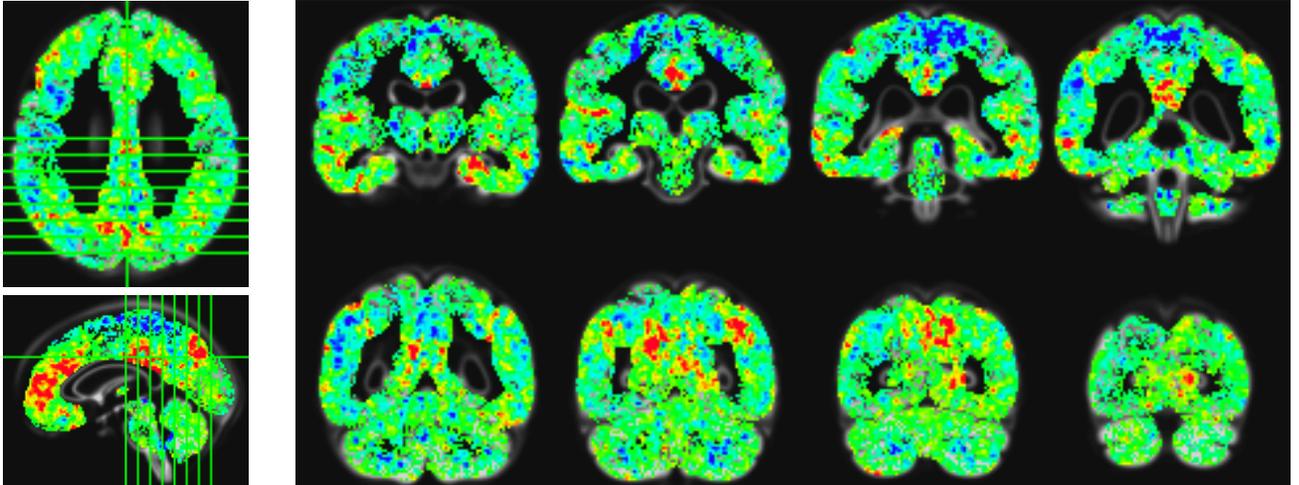


Figure 3: Voxels used in the classifier for FDG-PET baseline images. Weights are relative, and have no applicable units. Blue indicates negative weights, associated with AD, while green indicates zero or neutral weight, while red indicates positively weighted regions associated with healthy status. Green bars in the axial and sagittal views correspond to coronal slices.

FIGURE 4

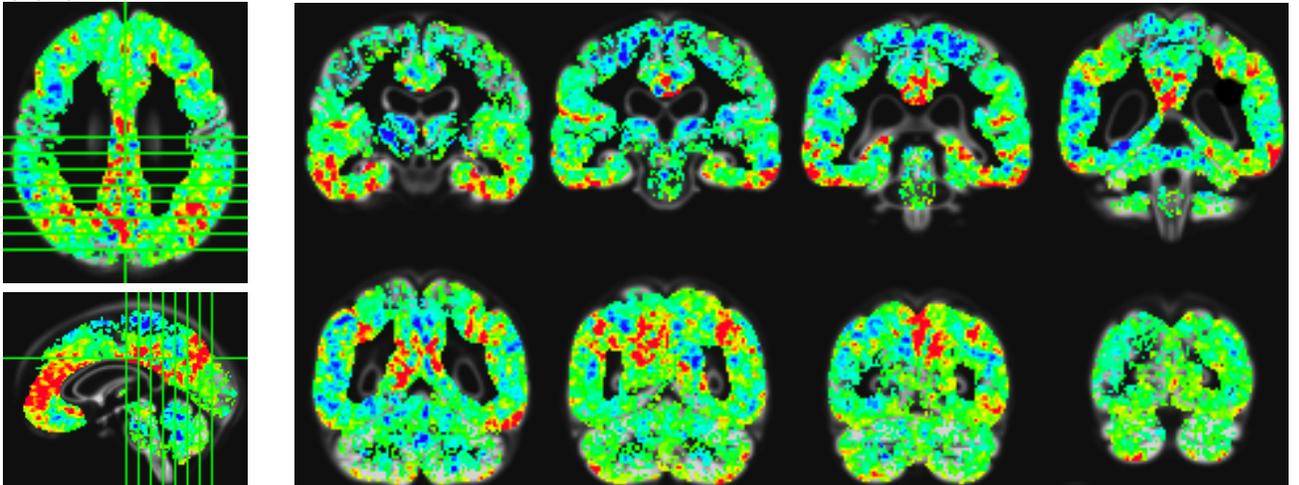


Figure 4: Voxels used in the classifier for FDG-PET images at 24 months. Weights are relative, and have no applicable units. Blue indicates negative weights, associated with AD, while green indicates zero or neutral weight, while red indicates positively weighted regions associated with healthy status. Green bars in the axial and sagittal views correspond to coronal slices.

494 MCI subjects are used for training purposes. (Hua et al., 2008, 2009; Davatzikos et al., 2009) This process  
495 was repeated for each modality separately, as well as in groups of modalities. That is, all imaging modalities  
496 were combined, as were all NPSEs and biological measures. The outputs for each subject are shown in Figure  
497 8. Subjects who remained stable are shown in blue; subjects who progressed to AD after 3 years or less  
498 are shown in red; subjects who reverted to normal cognitive status are shown in green. The four plots are  
499 divided between baseline (left) and longitudinal (right), and imaging-based (top) and NPSE-based (bottom)  
500 MMDMs. In each plot, a maximum accuracy cut-point is plotted as a solid black line. On the left we can  
501 see that neither of the baseline scans shows much differentiation between the groups, and the maximum  
502 accuracy separating line is essentially choosing the majority class. On the right, both the imaging-based and  
503 NPSE-based MMDMs provide better separation of the 2 groups. We also computed a set of MMDM scores

FIGURE 5

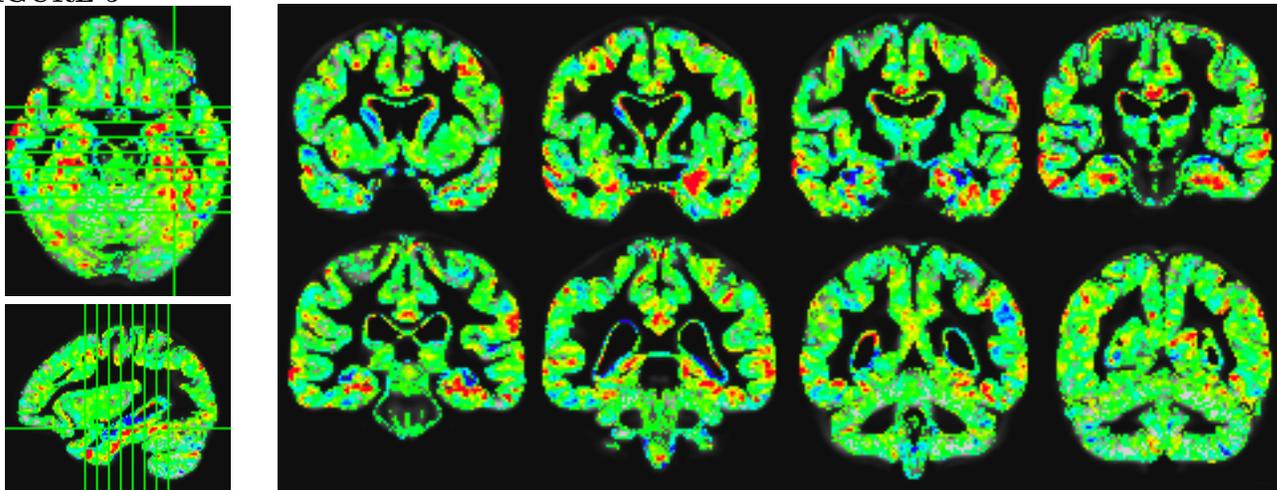


Figure 5: Voxels used in the classifier for TBM-processed MR images. Weights are relative, and have no applicable units. Blue indicates negative weights, associated with AD, while green indicates zero or neutral weight, while red indicates positively weighted regions associated with healthy status. Green bars in the axial and sagittal views correspond to coronal slices.

FIGURE 6

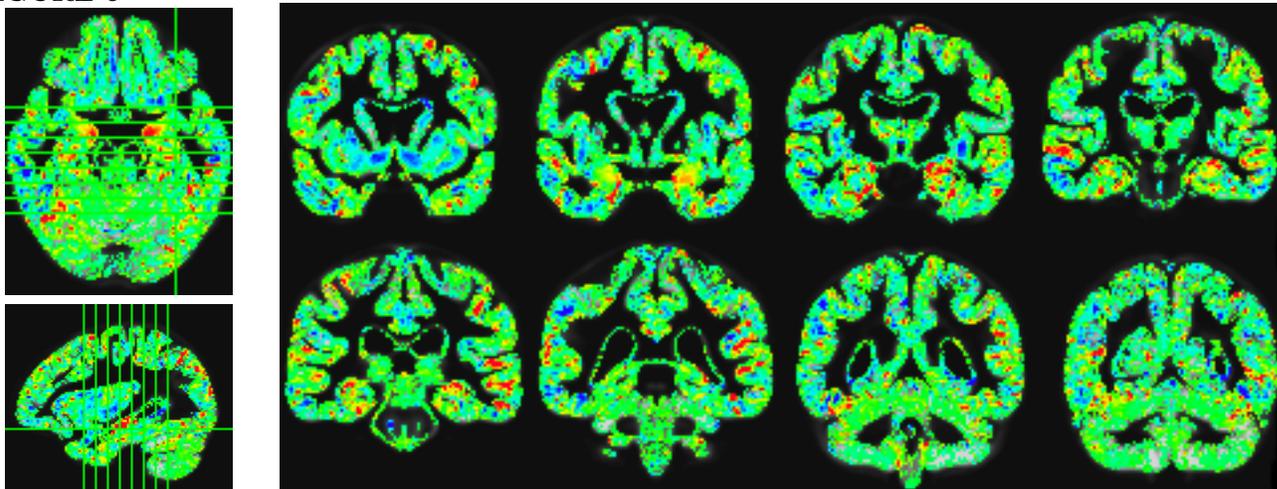


Figure 6: Voxels used in the classifier for VBM-processed (GM density) MR images. Weights are relative, and have no applicable units. Blue indicates negative weights, associated with AD, while green indicates zero or neutral weight, while red indicates positively weighted regions associated with healthy status. Green bars in the axial and sagittal views correspond to coronal slices.

504 based on CSF measures and APOE genetic markers, which did not show any ability to differentiate the 2  
 505 groups. An encouraging sign is that none of the reverting subjects were given negative scores.

506 In order to quantify these differences, we evaluated the degree of group-wise separation between progress-  
 507 ing, reverting, and stable MCI subjects, under each of the available modalities, using a  $t$ -test. As shown  
 508 in Table 6, the resulting  $p$ -values of the imaging-based MMDM (in separating progressing subjects from  
 509 non-progressing) are several orders of magnitude lower than those based on NPSEs at 24 months, and two  
 510 orders lower at baseline, suggesting that *imaging modalities offer a better view of future disease progression*  
 511 *than current cognitive status*. We believe this is an interesting result of our analysis.

512 Area under ROC curve results are shown in Table 7; the corresponding ROC curves are shown in Figure  
 513 9. For ROCs showing separation between progressing and reverting subjects, the AUCs are very high, as

FIGURE 7

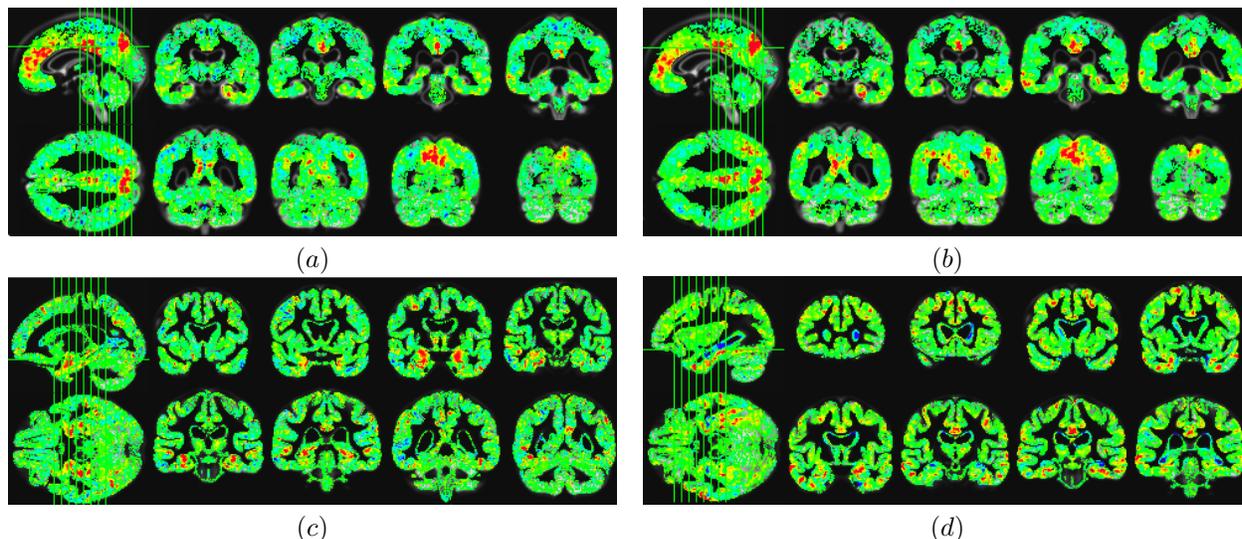


Figure 7: Voxel weights assigned by the MKL classifier when the outlier subjects were removed. (a) FDG-PET baseline images; (b) FDG-PET images at 24 months; (c) VBM-processed baseline MR images; (d) TBM-processed longitudinal MR scans.

514 we would expect. These curves are shown on the left in Figure 9. For comparison, we also computed  
 515 ROC curves for single modalities, which are also shown in the figure. Of special relevance is the fact that *the*  
 516 *MMDM based on imaging data alone outperformed all others*, both at baseline and at 24 months. **The second**  
 517 **comparison we made via ROC curves was between progressing subjects and all others. We accomplish this**  
 518 **by using a different ground truth for computing the ROC curves.** In this case, the task is to understand  
 519 which of the MCI subjects will progress to AD in the near term (2-3 years), and which will remain stable or  
 520 revert. These curves are shown on the right in Figure 9. In this case, the imaging-based MMDM, (shown  
 521 in green) outperformed all others, most significantly at 24 months. The AUC for the image-based MMDM  
 522 was 0.79, while that of the NPSE-based MMDM was 0.74. **The highest leave-one-out accuracy achieved**  
 523 **by the image-based MMDM was 0.723. For the NPSE the highest accuracy was 0.681** For the Biological  
 524 measure-based MMDMs, it was not possible to achieve an accuracy greater than chance.

525

**TABLE 6** t-statistic p-values for comparisons between MMDMs of stable MCI subjects, progressing subjects, and reverting subjects.

Modalities used	Reverting vs. rest	Progressing vs. rest
Biological measures (baseline)	0.65	0.58
Imaging Data (baseline)	$1.31 \times 10^{-3}$	$1.78 \times 10^{-6}$
Imaging Data (longitudinal)	$5.69 \times 10^{-4}$	$3.29 \times 10^{-7}$
NPSEs (baseline)	$2.63 \times 10^{-3}$	$5.51 \times 10^{-4}$
NPSEs (longitudinal)	$2.44 \times 10^{-4}$	$2.19 \times 10^{-6}$

Table 6: Significance of group-level differences in MMDM scores assigned to MCI subjects. There are 3 groups of MCI subjects - those who reverted to normal status, those who remained stable for 3 years, and those who progressed to full AD in 3 years.

526

FIGURE 8

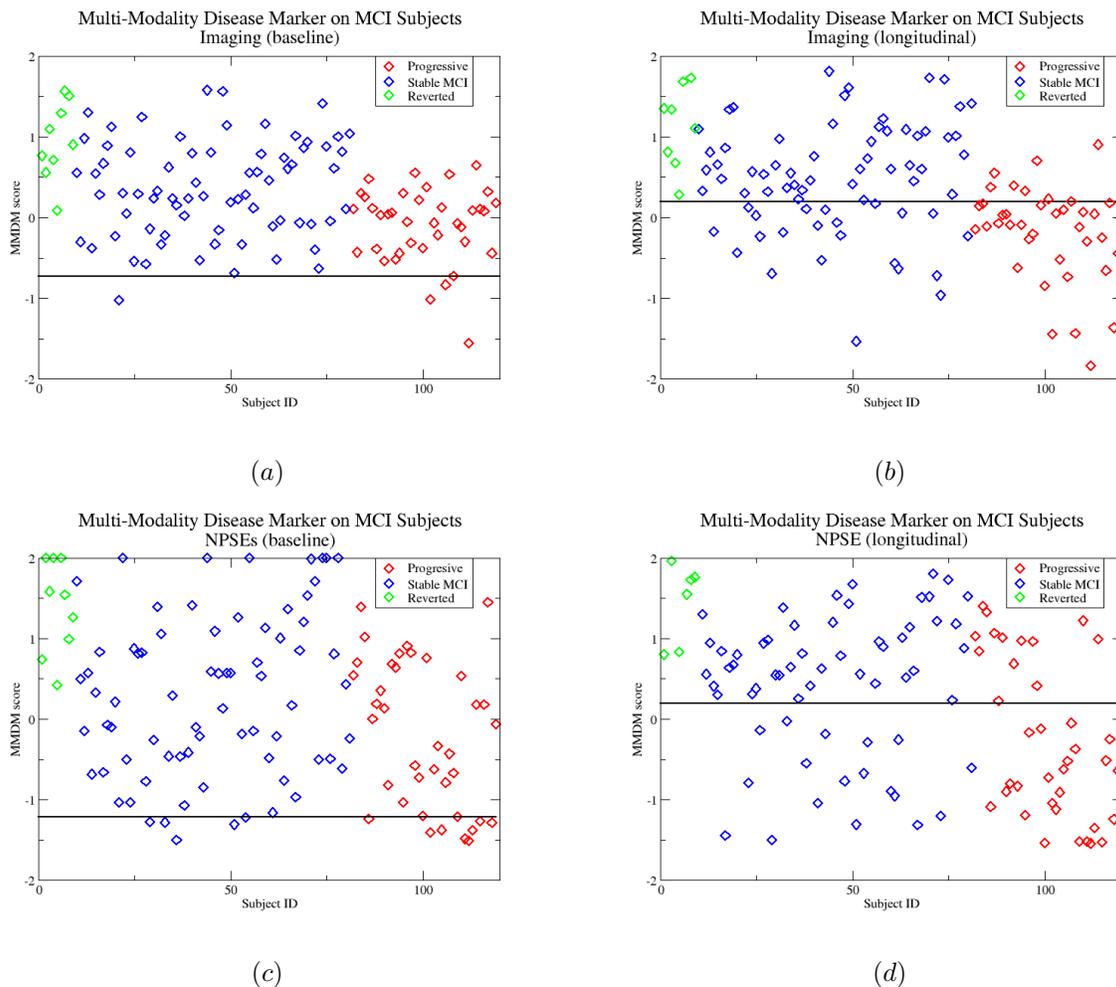


Figure 8: MMDMs applied to the MCI population. Subjects which remained stable are shown in blue; subjects which progressed to AD are shown in red; subjects which reverted to normal cognitive status are shown in green. In each figure, a line giving maximal post-hoc accuracy is shown. Note that in some cases, the best accuracy can be achieved by simply labeling all subjects as the majority class. In some cases, MMDM scores were truncated to  $\pm 2$  so as to preserve the relative scales. On the left (a,c) are shown MMDMs based on information available at baseline. Note the homogeneity of the groups, leading to poor separability. Imaging-based MMDMs are shown at the top (a), while MMDMs based on NPSEs are shown below (c). On the right (b,d) are shown MMDMs based on all modalities available at 24 months. Note the improved separability between the progressing (red) and stable (blue) MCI subjects. Note that the imaging-based marker above (b) shows slightly greater separation of the 2 groups.

## 5 Discussion

527

528 We have shown in our experiments that our approach can offer a flexible means of integrating multiple sources  
 529 of data into a single automated classification framework. As more types of information about subjects become  
 530 available, either through new scanning modalities or new processing methods, they can simply be added to  
 531 this framework as additional kernel matrices in a seamless manner. For instance, rather than choose whether  
 532 to use TBM or VBM in our experiments, we used *both* by delegating the task of choosing the better (*i.e.*,  
 533 more discriminative) view of the data to our model.

**TABLE 7** Area Under ROC results for different classes of MMDMs in predicting MCI progression to AD.

Modalities used	Progressing vs. Reverting	Progressing vs. Rest
Biological measures (baseline)	0.4368	0.5292
Imaging Data (baseline)	0.9532	0.7378
Imaging Data (longitudinal)	<b>0.9737</b>	<b>0.7911</b>
NPSEs (baseline)	0.9298	0.6693
NPSEs (longitudinal)	0.9415	0.7385
All Modalities	0.9708	0.7667

Table 7: Area under ROC curves for predicting whether MCI subjects will progress to AD or not. In the left column are AU ROCs for the task of separating only progressing subjects from reverting subjects, while ignoring stable MCI subjects. On the right are AU ROCs for separating progressing subjects from all other subjects.

534 The principal novelty of this work is to introduce a new machine learning algorithm, Multi-Kernel Learn-  
535 ing, to the application of discriminating different stages of AD using neuroimaging and other biological  
536 measures. Many existing works (Davatzikos et al., 2008a,b; Fan et al., 2008b,a; Vemuri et al., 2008; Duch-  
537 esne et al., 2008; Davatzikos et al., 2009; Querbes et al., 2009; Klöppel et al., 2008; Ramírez et al., 2009;  
538 Kohannim et al., 2010; Walhovd et al., 2010), use either general linear models based on summary statistics,  
539 or machine learning algorithms such as SVMs, logistic regression, or AdaBoost, with extensive pre- and  
540 post-processing of imaging data which adapts these methods to the particular application. Of the machine  
541 learning methods mentioned here, all three are discriminative max-margin learning algorithms. Logistic re-  
542 gression uses a sigmoid function to approximate the hinge-loss function, and must be optimized via iterative  
543 methods. AdaBoost implicitly finds a margin by iteratively increasing the importance of examples which  
544 are misclassified, much the same way that examples inside the margin become support vectors in the SVM  
545 framework. Our method shares some commonalities in the sense that pre-processing of brain scans is also  
546 required before a classifier can be trained. However, by incorporating MKL, we can extend this framework  
547 to allow seamless integration of multiple sources of data while controlling the complexity of the resulting  
548 classifier without the need for creating summary statistics, (which discard a large amount of information).

549 We note that several studies have reported better *raw* performance at classifying AD and control subjects.  
550 There are several factors which can affect such results. First, there is the issue of the severity of the disease,  
551 and of the availability of gold-standard diagnosis. For instance, the authors of (Klöppel et al., 2008) reported  
552 that their accuracy suffered when autopsy data were not available due to the difficulty of diagnosing AD *in*  
553 *vivo*. The ADNI data set, on which our experiments were based, consists entirely of living subjects, having  
554 relatively mild AD. (See Table 1). Other studies have used ADNI subject data (Davatzikos et al., 2009;  
555 Querbes et al., 2009; Fan et al., 2008a), and while some have reported better performance than we have,  
556 issues such as image registration and warping, subject inclusion criteria (*e.g.*, image quality), or choice of  
557 feature extraction / representation might have a greater effect on final outcomes. A recent study, Cuingnet  
558 et al. (2010), addressed exactly these issues, finding that when these issues are controlled, the accuracy  
559 results are closer to those reported in this study. (See Table 4.) For example, if a pre-processing method is  
560 found to be particularly useful for discriminative purposes, that method can be swapped with our current  
561 pre-processing methods, or incorporated as additional kernels. The more important comparison is between  
562 single modality and multi-modality methods, *using the same data and pre-processing pipeline*. In addition,  
563 our experiments comparing MKL with a concatenated-features SVM show that MKL has advantages in the  
564 presence of non-informative kernels.

565 **Single-modality results** Our experiments in single-modality AD classification give an indication of the  
566 relative merits of various scanning modalities. We note first that in FDG-PET scans, the top performing ker-  
567 nels are those which make use of at least 65,000 voxels, indicating that a performance gain of five percentage  
568 points or more can be made from using the *entire* brain volume, rather than using smaller selected regions.  
569 <sup>5</sup> That is, while most subjects can be identified by examining smaller regions, some subjects can only be  
570 identified by examination of whole-brain atrophy. This suggests that there is a small group of subjects having

<sup>5</sup>The authors of (Fan et al., 2008b) found similar results in FDG-PET images.

FIGURE 9

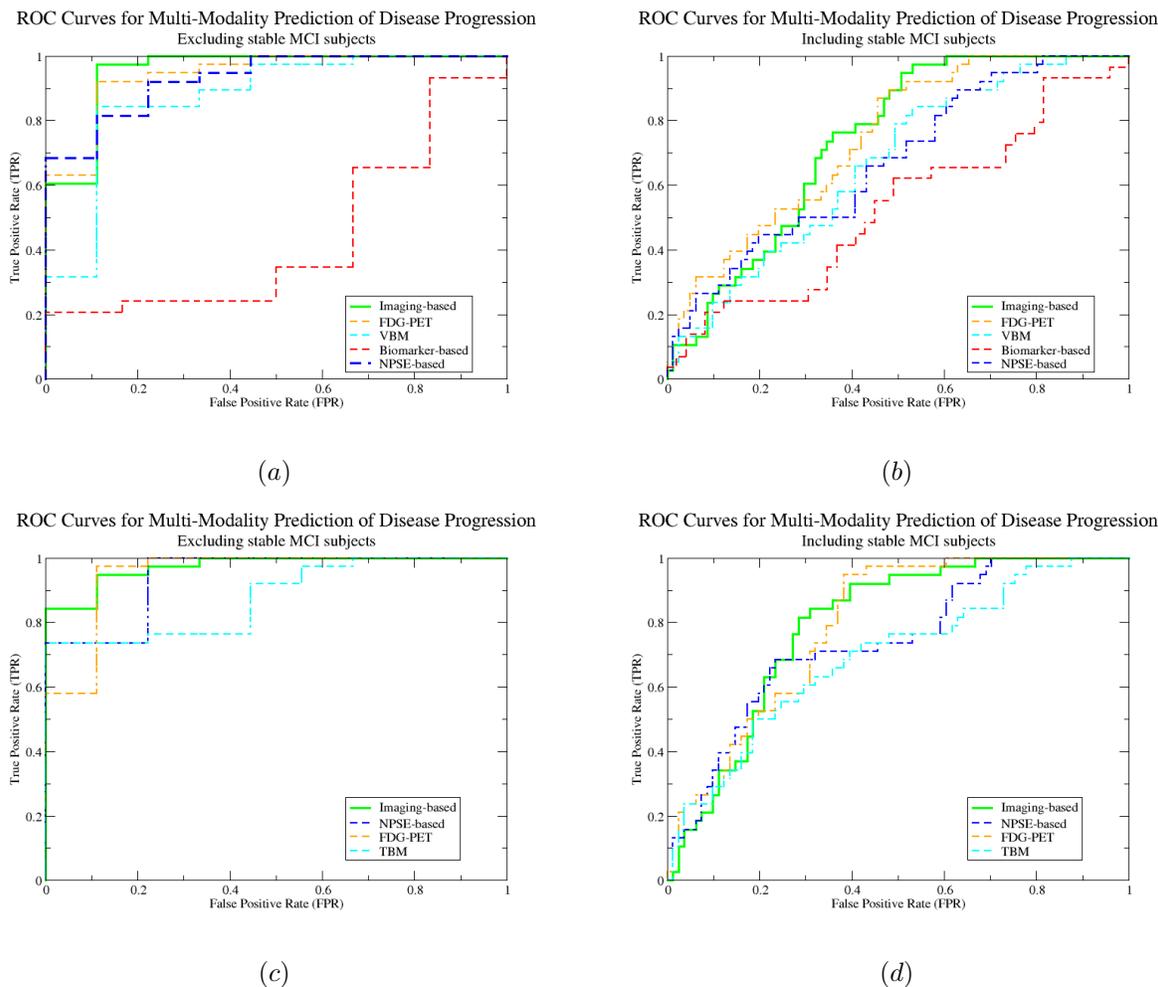


Figure 9: ROC curves for multi-modality learning on disease progression of MCI subjects using various disease markers. The ROC curves for separating progressing and reverting MCI subjects on the left (a,c). The ROC curves for separating progressing MCI subjects from all others are shown on the right, (b,d). The top row (a,b) shows the curves derived from information available at baseline, while those on the bottom (c,d) were derived from scans and markers taken at both baseline and 24-months.

571 atypical disease progression (in the case of AD subjects) or that some control subjects may show early signs  
 572 of disease. A somewhat surprising result is that longitudinal analysis of FDG-PET images did not have  
 573 much discriminative power. Neither of the two methods we considered (voxel-wise temporal difference, and  
 574 voxel-wise temporal ratio) had accuracy higher than about 65%. This is perhaps an indication that signs  
 575 of atrophy in FDG-PET images accumulate slowly enough that changes over a 2-year period alone are not  
 576 enough to distinguish AD with high accuracy.

577 In the MR-based modalities, we can see that in baseline VBM images, the highest performing kernels are  
 578 those that focus on small brain regions of a few thousand voxels, while in TBM images, the best performance  
 579 is obtained from larger regions of about 25,000 voxels. We interpret this to mean that (in classifying AD  
 580 and control subjects,) the most indicative signs of atrophy already present at baseline can be found in  
 581 hippocampal and para-hippocampal regions (not shown), but the atrophy occurring at the stage of full AD  
 582 (*i.e.*, that which occurs in the two years following diagnosis), is more diffuse. This suggests that early signs  
 583 of AD are more likely to be concentrated in smaller regions, such as the hippocampus, and other structures

584 known to be affected by AD.

585 Secondly, we note that linear kernels performed as well as, or better than quadratic and polynomial  
586 kernels in all modalities examined, indicating that there are few quadratic or exponential effects which can  
587 be used for discriminative purposes. This can be interpreted that indications of pathology in each voxel  
588 contribute independently and cumulatively to the final diagnosis.

589 **Multi-modality results** An interesting comparison which arose in our experiments was between the  
590 various imaging-based kernels *individually*, (see Figure 1), and the MKL experiments combining groups of  
591 modalities (see Table 4). MKL produces *linear* combinations of kernels, and therefore does not examine the  
592 interactions between them when evaluating new subjects. This means that the ideal situation is where the  
593 errors present in each kernel matrix are drawn randomly and independently. When combining modalities  
594 with strong similarities, it is therefore expected that some errors will cancel out, to the extent that those  
595 errors do not themselves arise from shared properties of both modalities. The rationale for combining  
596 modalities into groups for comparison is that while imaging modalities *are* expected to contain distinct (and  
597 useful) information about each subject, we expect that they will have some information in common. For  
598 instance, properties such as total inter-cranial volume or particular anatomical artifacts will be present in  
599 different scanning modalities, but not in other biological measures. Thus, we first examine MKL’s ability  
600 to integrate groups of similar measures and modalities, before examining its ability to combine dissimilar  
601 sources of information.

602 First, we note that none of the individual kernels derived from imaging modalities achieved an accuracy  
603 greater than MKL when given the combination of imaging modalities. Moreover, when MKL was given the  
604 *entire* set of kernels from all available sources of information, it outperformed any of the groups of modalities,  
605 except for the NPSEs, where the differences were not significant. This is expected, because clinical diagnosis  
606 is already known, meaning that the disease has already reached a stage where cognitive status effects are  
607 measurable, in contrast to earlier stages, in which anatomical and physiological changes have begun to occur,  
608 but outward signs have not. Indeed, in the analysis of MCI progression (Tables 6 and 7), it is the imaging-  
609 based modalities which have the strongest performance. Finally, it is interesting that for the biological  
610 measures, such as CSF assays and APOE genotypes, while there is certainly some information contained in  
611 the kernels generated from these measures, by themselves they do not have nearly the discriminative power  
612 of either the imaging modalities, or the NPSEs. This may be due in part to the fact that these measures are  
613 not available for all subjects.

614 In Table 7 it may be surprising that the MMDM trained on all available modalities underperformed the  
615 one trained only on longitudinal imaging modalities. This is likely due to the fact that the training task and  
616 evaluation task were closely related, but slightly different. Thus, the subkernel weights estimated to give the  
617 optimal performance on the training task (AD vs. controls), may have been slightly less than optimal on the  
618 related task, (MCI progression). Despite this, the disparity in performance is small, and the MMDM using  
619 all combined modalities still outperformed all other MMDMs. It is also interesting to note that while the  
620 NPSEs dominated in the AD vs. control task of Section 4.1, in this task, the longitudinal NPSEs are roughly  
621 at parity with the baseline imaging modalities. (See Tables 6 and 7.) This suggests that signs of impending  
622 progression from MCI to AD are present in the imaging modalities approximately *two years* ahead of clinical  
623 psychological measures.

624 **MKL-norm results** In our experiments with varying MKL norm, we found that norms which encouraged  
625 sparsity performed slightly worse than those which do not, suggesting that information is being needlessly  
626 discarded. The results in Table 5 show that above about 1.5, sparsity makes less of a difference, but at 1  
627 or 1.25, sparsity is encouraged enough to affect MKL’s performance. In contrast, the concatenated-features  
628 SVM’s performance was significantly lower overall, as it has no mechanism for discarding non-informative  
629 kernels, especially when there are more kernels from many different sources. When given only kernels from  
630 a single modality, the SVM’s performance was closer to parity with MKL, however, this is expected, due to  
631 the relative ease of combining kernels from similar sources of information. Rather, it is when there is greater  
632 variety in the information content of the various kernels that MKL incrementally shows an advantage over  
633 the concatenated-features SVM. This demonstrates that regardless of the norm chosen, MKL has the ability  
634 to automatically detect and discard sets of features which do not contribute significantly to the optimal  
635 classifier. One could, in theory, manually select which features to include, and how to weight them, but this  
636 would essentially emulate the MKL process by hand using a regular SVM. With the proper construction of  
637 kernels, it is even conceivable that MKL could be used to automatically select ROIs.

638 **Brain regions selected** The classifier chosen by MKL consists of a set of kernel combination weights  
639  $\beta$ , as well as a set of example combination weights  $\alpha$ . These weights can be combined to give a single linear  
640 classifier based on voxel-wise features. The distribution of these voxel-weights chosen by the MKL algorithm  
641 therefore gives some insight into the relative importance of various brain regions, and we expect that a good  
642 classifier will place greater weight on regions known to be involved in AD.

643 It is well known that the Posterior Cingulate Cortex is involved in memory retrieval and related self  
644 referential processes (Northoff and Bermpohl, 2004; Piefke et al., 2003; Shannon and Buckner, 2004). As part  
645 of the limbic system, it has reciprocal connections with other memory areas including the dorsomedial and  
646 dorsolateral prefrontal cortex, the posterior parahippocampal cortex, presubiculum, hippocampus, entorhinal  
647 cortex, and thalamus (Mesulam, 2000). Previous imaging studies suggest the PCC is affected in AD even  
648 before clinical symptoms appear, consistent with the very early memory symptoms in AD (Xu et al., 2009;  
649 Ries et al., 2006). Interestingly, the earliest cerebral hypometabolism finding in AD involves the PCC-  
650 precuneus rather than the hippocampus (Villain et al., 2008). Although the mechanism connecting cortical  
651 atrophy and hypometabolism in neurodegenerative disorders is not fully understood, intuitively, a positive  
652 relationship is expected. Both brain atrophy and cerebral hypometabolism reflect loss of neurons/synapses  
653 (Bobinski et al., 1999) and decrease in synaptic density/activity (Rocher et al., 2003). As mentioned in  
654 section 4.2, the brain regions selected by the MKL algorithm in FDG-PET images, as show in Figures 3 to  
655 4, include the PCC and precuneus, the lateral parietal lobules, hippocampal and medial temporal regions,  
656 and the pre-frontal midline.

657 In MR longitudinal images (TBM, Figure 5), regions well-known to be atrophic in AD, such as the  
658 hippocampus, parahippocampal gyri, fusiform gyri and other middle temporal structures (Braak and Braak,  
659 1991) are well highlighted. Expansion, (or reduced contraction) is associated with healthy status, and thus  
660 these regions are given positive weights, shown in red. Conversely, expansion in ventricles, and in the CSF  
661 surrounding the hippocampus is shown in blue. Expansion in these regions is correlated with AD pathology,  
662 and so these regions are given negative weights. In the baseline gray matter density images, (VBM, Figure  
663 6) similar hippocampal and medial temporal regions are shown.

664 **MCI conversion** The task of predicting conversion from MCI to full AD is known to be difficult,  
665 (Querbes et al., 2009; Davatzikos et al., 2009), and presents challenges *beyond* that of classifying AD and  
666 control subjects, or even that of classifying AD/control and MCI subjects. This difficulty arises largely from  
667 the “lag” between brain atrophy and cognitive decline. There are several interesting aspects of the MMDMs  
668 we have examined. First, we note that at baseline, neither NPSEs nor imaging modalities have a strong  
669 ability to detect which subjects will convert to AD. This may be a result of the ADNI selection criteria for  
670 MCI subjects – that is, MCI subjects are chosen so as to have very homogeneous cognitive characteristics at  
671 baseline, and so we expect that NPSEs will not be able to differentiate between progressing and stable MCI  
672 subjects very well. While the MMDM based on all combined imaging modalities does have a better AUC at  
673 baseline than the NPSEs, the improvement shown by the MMDM based on longitudinal imaging modalities  
674 suggests that a significant portion of the neurodegeneration responsible for the subjects’ conversion to AD  
675 takes place after MCI diagnosis. In addition, between baseline and 24 months, the imaging-based MMDM  
676 outperforms the NPSE-based MMDM by an even wider margin, as shown by the AUCs and  $p$ -values in  
677 Tables 6 and 7. This leads us to believe that while NPSEs can be a better marker for subjects who *already*  
678 are showing AD-related cognitive decline, the imaging modalities have slightly better predictive value for  
679 future decline. We expect that further progress can be made in adapting multi-kernel methods to work  
680 specifically with imaging data, allowing greater accuracy in identifying future patterns. Finally, we find it  
681 interesting that combining all imaging markers into a single MMDM offers a slight improvement over the  
682 best single imaging modality, which tends to be FDG-PET. This improvement is relatively stable over time,  
683 between baseline and 24 months.

## 684 6 Conclusion

685 In this paper we have presented a new application of recent developments from the machine learning literature  
686 to early detection of AD-related pathology. Using this measure of AD pathology, we constructed a predictive  
687 marker for MCI progression to AD. This method is fully *multi-modal* – that is, it incorporates all available  
688 sources of input relating to subjects, yielding a unified Multi-Modal Disease Marker (MMDM). Our results  
689 on the ADNI population indicate that this method has the potential to detect subtle changes in MCI subjects

690 which may provide clues as to whether a subject will convert to AD, or remain stable. In particular, we have  
691 shown that imaging modalities have better ability to predict such outcomes than baseline neuropsychological  
692 scores, which is consistent with the view that neurological changes detected in neuroimages can *precede*  
693 clinically detectable declines in cognitive status. Our ongoing work focuses on further developing this method  
694 – which will permit even higher accuracy and sensitivity, and allow predictions at the level of individual  
695 subjects to be made with high confidence.

## 696 Acknowledgments

697 This research was supported in part by NIH grants R21-AG034315 (Singh) and R01-AG021155 (Johnson).  
698 Hinrichs is funded via a University of Wisconsin–Madison CIBM (Computation and Informatics in Biology  
699 and Medicine) fellowship (National Library of Medicine Award 5T15LM007359). Partial support for this  
700 research was also provided by the University of Wisconsin-Madison UW ICTR through an NIH Clinical  
701 and Translational Science Award (CTSA) 1UL1RR025011, a Merit Review Grant from the Department of  
702 Veterans Affairs, the Wisconsin Comprehensive Memory Program, and the Society for Imaging Informatices  
703 in Medicine (SIIM). The authors also acknowledge the facilities and resources at the William S. Middleton  
704 Memorial Veterans Hospital, and the Geriatric Research, Education, and Clinical Center (GRECC).

705 Data collection and sharing for this project was funded by the Alzheimer’s Disease Neuroimaging Initia-  
706 tive (ADNI) (National Institutes of Health Grant U01 AG024904). ADNI is funded by the National Institute  
707 on Aging, the National Institute of Biomedical Imaging and Bioengineering, and through generous contribu-  
708 tions from the following: Abbott, AstraZeneca AB, Bayer Schering Pharma AG, Bristol-Myers Squibb, Eisai  
709 Global Clinical Development, Elan Corporation, Genentech, GE Healthcare, GlaxoSmithKline, Innogenetics,  
710 Johnson and Johnson, Eli Lilly and Co., Medpace, Inc., Merck and Co., Inc., Novartis AG, Pfizer Inc, F.  
711 Hoffman-La Roche, Schering-Plough, Synarc, Inc., as well as non-profit partners the Alzheimer’s Association  
712 and Alzheimer’s Drug Discovery Foundation, with participation from the U.S. Food and Drug Administra-  
713 tion. Private sector contributions to ADNI are facilitated by the Foundation for the National Institutes  
714 of Health (www.fnih.org). The grantee organization is the Northern California Institute for Research and  
715 Education, and the study is coordinated by the Alzheimer’s Disease Cooperative Study at the University of  
716 California, San Diego. ADNI data are disseminated by the Laboratory for Neuro Imaging at the University  
717 of California, Los Angeles. This research was also supported by NIH grants P30 AG010129, K01 AG030514,  
718 and the Dana Foundation.

719 The authors are grateful to Donald McLaren, Moo K. Chung and Sanjay Asthana for many suggestions  
720 and ideas.

## 721 References

- 722 M. S. Albert, M. B. Moss, R. Tanzi, and K. Jones. Preclinical prediction of AD using neuropsychological tests.  
723 *Journal of the International Neuropsychological Society*, 7(05):631–639, 2001.
- 724 H. Arimura, T. Yoshiura, S. Kumazawa, K. Tanaka, H. Koga, F. Mihara, H. Honda, S. Sakai, F. Toyofuku, and  
725 Y. Higashida. Automated method for identification of patients with Alzheimer’s disease based on three-dimensional  
726 MR images. *Academic Radiology*, 15(3):274–284, 2008.
- 727 J. Ashburner. A fast diffeomorphic image registration algorithm. *Neuroimage*, 38(1):95, 113 2007.
- 728 J. Ashburner and K. J. Friston. Voxel-Based Morphometry - The Methods . *Neuroimage*, 11(6):805–821, 2000.
- 729 G. Bakir, T. Hofmann, and B. Schölkopf. *Predicting structured data*. The MIT Press, 2007.
- 730 C. Bishop. *Pattern Recognition and Machine Learning*. Springer New York, 2006.
- 731 M. Bobinski, M. J. De Leon, J. Wegiel, S. Desanti, A. Convit, L. A. Saint Louis, H. Rusinek, and H. M. Wis-  
732 niewski. The histological validation of post mortem magnetic resonance imaging-determined hippocampal volume  
733 in Alzheimer’s disease. *Neuroscience*, 95(3):721–725, 1999.
- 734 E. Braak, K. Griffin, K. Arai, J. Bohl, H. Bratzke, and H. Braak. Neuropathology of Alzheimer’s disease: what is  
735 new since A. Alzheimer? *European Archives of Psychiatry and Clinical Neuroscience*, 249(9):14–22, 1999.

- 736 H. Braak and E. Braak. Neuropathological staging of Alzheimer-related changes. *Acta neuropathologica*, 82(4):  
737 239–259, 1991.
- 738 E. Canu, D. G. McLaren, M. E. Fitzgerald, B. B. Bendlin, G. Zoccatelli, F. Alessandrini, F. B. Pizzini, G. K. Ricciardi,  
739 A. Beltramello, S. C. Johnson, et al. Microstructural Diffusion Changes are Independent of Macrostructural Volume  
740 Loss in Moderate to Severe Alzheimer’s Disease. *Journal of Alzheimer’s Disease*, 2010.
- 741 C. Cortes and V. Vapnik. Support-vector networks. *Machine learning*, 20(3):273–297, 1995.
- 742 R. Cuingnet, E. Gérardin, J. Tessieras, G. Auzias, S. Lehéricy, and M. O. Habert. Automatic classification of  
743 patients with Alzheimer’s disease from structural MRI: A comparison of ten methods using the ADNI database.  
744 *NeuroImage*, 2010.
- 745 C. Davatzikos, Y. Fan, X. Wu, D. Shen, and S.M. Resnick. Detection of prodromal Alzheimer’s disease via pattern  
746 classification of magnetic resonance imaging. *Neurobiology of Aging*, 29(4):514–523, 2008a.
- 747 C. Davatzikos, S.M. Resnick, X. Wu, P. Parmpi, and C.M. Clark. Individual patient diagnosis of AD and FTD via  
748 high-dimensional pattern classification of MRI. *Neuroimage*, 41(4):1220–1227, 2008b.
- 749 C. Davatzikos, F. Xu, Y. An, Y. Fan, and S. M. Resnick. Longitudinal progression of Alzheimer’s-like patterns of  
750 atrophy in normal older adults: the SPARE-AD index. *Brain*, 132(8):2026–2035, 2009.
- 751 O. Demirci, V. P. Clark, and V. D. Calhoun. A projection pursuit algorithm to classify individuals using fMRI data:  
752 Application to schizophrenia. *Neuroimage*, 39(4):1774–1782, 2008.
- 753 L. deToledo-Morrell, T. R. Stoub, M. Bulgakova, RS Wilson, DA Bennett, S. Leurgans, J. Wu, and DA Turner.  
754 MRI-derived entorhinal volume is a good predictor of conversion from MCI to AD. *Neurobiology of Aging*, 25(9):  
755 1197–1203, 2004.
- 756 B. C. Dickerson, I. Goncharova, M. P. Sullivan, C. Forchetti, R. S<sub>z</sub> Wilson, D. A. Bennett, L. A. Beckett, and  
757 L. deToledo-Morrell. MRI-derived entorhinal and hippocampal atrophy in incipient and very mild Alzheimer’s  
758 disease. *Neurobiology of aging*, 22(5):747–754, 2001.
- 759 S. Duchesne, A. Caroli, C. Geroldi, C. Barillot, G. B. Frisoni, and D. L. Collins. MRI-Based Automated Computer  
760 Classification of Probable AD Versus Normal Controls. *IEEE Transactions on Medical Imaging*, 27(4):509–520,  
761 2008.
- 762 Y. Fan, N. Batmanghelich, C.M. Clark, and C. Davatzikos. Spatial patterns of brain atrophy in MCI patients,  
763 identified via high-dimensional pattern classification, predict subsequent cognitive decline. *Neuroimage*, 39(4):  
764 1731–1743, 2008a.
- 765 Y. Fan, S. M. Resnick, X. Wu, and C. Davatzikos. Structural and functional biomarkers of prodromal Alzheimer’s  
766 disease: a high-dimensional pattern classification study. *Neuroimage*, 41(2):277–285, 2008b.
- 767 P. V. Gehler and S. Nowozin. Let the kernel figure it out; principled learning of pre-processing for kernel classifiers.  
768 *Computer Vision and Pattern Recognition*, pages 2836–2843, 2009.
- 769 C. Hinrichs, V. Singh, L. Mukherjee, G. Xu, M. K. Chung, and S. C. Johnson. Spatially augmented LPBoosting for  
770 AD classification with evaluations on the ADNI dataset. *NeuroImage*, 48(1):138–149, 2009a.
- 771 C. Hinrichs, V. Singh, G. Xu, and S. C. Johnson. MKL for Robust Multi-modality AD Classification . *Medical Image*  
772 *Computing and Computer-Assisted Intervention*, 5762:786–794, 2009b.
- 773 J. M. Hoffman, K. A. Welsh-Bohmer, M. Hanson, B. Crain, C. Hulette, N. Earl, and R.E. Coleman. FDG PET  
774 imaging in patients with pathologically verified dementia. *Journal of Nuclear Medicine*, 41(11):1920–1928, 2000.
- 775 X. Hua, A. D. Leow, N. Parikshak, S. Lee, M. C. Chiang, A. W. Toga, C. R. Jack Jr., M. W. Weiner, and P. M.  
776 Thompson. Tensor-based morphometry as a neuroimaging biomarker for Alzheimer’s disease: an MRI study of  
777 676 AD, MCI, and normal subjects. *Neuroimage*, 43(3):458–469, 2008.
- 778 X. Hua, S. Lee, I. Yanovsky, A.D. Leow, Y.Y. Chou, A.J. Ho, B. Gutman, A.W. Toga, C.R. Jack Jr, M.A. Bernstein,  
779 et al. Optimizing power to track brain degeneration in Alzheimer’s disease and mild cognitive impairment with  
780 tensor-based morphometry: An ADNI study of 515 subjects. *NeuroImage*, 48(4):668–681, 2009.

- 781 K. Ishii, H. Sasaki, A. K. Kono, N. Miyamoto, T. Fukuda, and E. Mori. Comparison of gray matter and metabolic  
782 reduction in mild Alzheimers disease using FDG-PET and voxel-based morphometric MR studies. *European Journal*  
783 *of Nuclear Medicine and Molecular Imaging*, 32(8):959–963, 2005.
- 784 C. R. Jack Jr., M. M. Shiung, S. D. Weigand, P. C. O’Brien, J. L. Gunter, B. F. Boeve, D. S. Knopman, G. E. Smith,  
785 R. J. Ivnik, E. G. Tangalos, et al. Brain atrophy rates predict subsequent clinical conversion in normal elderly and  
786 amnesic MCI. *Neurology*, 65(8):1227–1231, 2005.
- 787 S. C. Johnson, T. W. Schmitz, M. A. Trivedi, M. L. Ries, B. M. Torgerson, C. M. Carlsson, S. Asthana, B. P.  
788 Hermann, and M. A. Sager. The influence of Alzheimer disease family history and apolipoprotein E varepsilon4  
789 on mesial temporal lobe activation. *Journal of Neuroscience*, 26(22):6069–6076, 2006.
- 790 M. Kloft, U. Brefeld, S. Sonnenburg, and A. Zien. Non-sparse regularization and efficient training with multiple  
791 kernels. 2010.
- 792 S. Klöppel, C.M. Stonnington, C. Chu, B. Draganski, R.I. Scahill, J.D. Rohrer, N.C. Fox, C.R. Jack, J. Ashburner,  
793 and R.S. Frackowiak. Automatic classification of MR scans in Alzheimer’s disease. *Brain*, 131(3):681–689, 2008.
- 794 W. E. Klunk, H. Engler, A. Nordberg, Y. Wang, G. Blomqvist, D. P. Holt, M. Bergström, I. Savitcheva, G. F.  
795 Huang, S. Estrada, et al. Imaging brain amyloid in Alzheimer’s disease with Pittsburgh Compound-B. *Annals of*  
796 *neurology*, 55(3):306–319, 2004.
- 797 O. Kohannim, X. Hua, D.P. Hibar, S. Lee, Y.Y. Chou, A.W. Toga, C.R. Jack, M.W. Weiner, and P.M. Thompson.  
798 Boosting power for clinical trials using classifiers based on multiple biomarkers. *Neurobiology of Aging*, 2010.
- 799 G. R. G. Lanckriet, N. Cristianini, P. Bartlett, L. E. Ghaoui, and M. I. Jordan. Learning the kernel matrix with  
800 semidefinite programming. *Journal of Machine Learning Research*, 5:27–72, 2004.
- 801 H. Matsuda. Cerebral blood flow and metabolic abnormalities in Alzheimer’s disease. *Annals of Nuclear Medicine*,  
802 15(2):85–92, 2001.
- 803 M. M. Mesulam. *Principles of behavioral and cognitive neurology*. Oxford University Press, USA, 2000.
- 804 S. Minoshima, N. L. Foster, and D. E. Kuhl. Posterior cingulate cortex in Alzheimer’s disease. *Lancet*, 344(8926):  
805 895, 1994.
- 806 C. Misra, Y. Fan, and C. Davatzikos. Baseline and longitudinal patterns of brain atrophy in MCI patients, and their  
807 use in prediction of short-term conversion to AD: Results from ADNI. *Neuroimage*, 44(4):1415–1422, 2008.
- 808 S. G. Mueller, M. W. Weiner, L.J. Thal, R. C. Petersen, C. R. Jack, W. Jagust, J. Q. Trojanowski, A. W. Toga,  
809 and L. Beckett. Ways toward an early diagnosis in Alzheimers disease: The Alzheimers Disease Neuroimaging  
810 Initiative (ADNI). *Journal of the Alzheimer’s Association*, 1(1):55–66, 2005.
- 811 L. Mukherjee, V. Singh, J. Peng, and C. Hinrichs. Learning Kernels for variants of Normalized Cuts: Convex  
812 Relaxations and Applications. *Computer Vision and Pattern Recognition*, 2010.
- 813 G. Northoff and F. Bermpohl. Cortical midline structures and the self. *Trends in Cognitive Sciences*, 8(3):102–107,  
814 2004.
- 815 M. Piefke, P. H. Weiss, K. Zilles, H. J. Markowitsch, and G. R. Fink. Differential remoteness and emotional tone  
816 modulate the neural correlates of autobiographical memory. *Brain*, 126(3):650–668, 2003.
- 817 O. Querbes, F. Aubry, J. Pariente, J. A. Lotterie, J. F. Demonet, V. Duret, M. Puel, I. Berry, J. C. Fort, and  
818 P. Celsis. Early diagnosis of Alzheimer’s disease using cortical thickness: impact of cognitive reserve. *Brain*, 132  
819 (8):2036–2047, 2009.
- 820 A. Rakotomamonjy, F. Bach and S. Canu, and Y. Grandvalet. SimpleMKL. *Journal of Machine Learning Research*,  
821 9:2491–2521, 2008.
- 822 J. Ramírez, J. M. Górriz and D. Salas-Gonzalez, A. Romero, M. López, I. Álvarez, and M. Gómez-Río. Computer-  
823 aided diagnosis of Alzheimer’s type dementia combining support vector machines and discriminant set of features.  
824 *Information Sciences*, 2009.
- 825 E. M. Reiman, R. J. Caselli, L. S. Yun, K. Chen, D. Bandy, S. Minoshima, S.N. Thibodeau, and D. Osborne.  
826 Preclinical Evidence of Alzheimer’s Disease in Persons Homozygous for the  $\epsilon 4$  Allele for Apolipoprotein E. *New*  
827 *England Journal of Medicine*, 334(12):752–758, 1996.

- 828 M. L. Ries, T. W. Schmitz, T. N. Kawahara, B. M. Torgerson, M. A. Trivedi, and S. C. Johnson. Task-dependent  
829 posterior cingulate activation in mild cognitive impairment. *Neuroimage*, 29(2):485–492, 2006.
- 830 A. B. Rocher, F. Chapon, X. Blaizot, J. C. Baron, and C. Chavoix. Resting-state brain glucose utilization as measured  
831 by PET is directly related to regional synaptophysin levels: a study in baboons. *Neuroimage*, 20(3):1894–1898,  
832 2003.
- 833 B. Schoelkopf and A. Smola. *Learning from Kernels*. MIT Press, 2002.
- 834 M. L. Schroeter, T. Stein, N. Maslowski, and J. Neumann. Neural correlates of Alzheimer’s disease and mild cognitive  
835 impairment: A systematic and quantitative meta-analysis involving 1351 patients. *NeuroImage*, 47(4):1196–1206,  
836 2009.
- 837 B. J. Shannon and R. L. Buckner. Functional-anatomic correlates of memory retrieval that suggest nontraditional  
838 processing roles for multiple distinct regions within posterior parietal cortex. *Journal of Neuroscience*, 24(45):  
839 10084–10092, 2004.
- 840 L. Shen, J. Ford, F. Makedon, and A. Saykin. Hippocampal shape analysis: surface-based representation and  
841 classification. In *Proceedings of SPIE*, volume 5032, pages 253–264, 2003.
- 842 N. Shock, R. Greulich, and R. Andres et al. Normal human aging: the Baltimore Longitudinal Study of Aging.  
843 Washington, DC: US Government Printing Office, 1984.
- 844 G. Small, L. M. Ercoli, D. H. Silverman, S.C. Huang, S. Komo, S.Y. Bookheimer, H. Lavretsky, K. Miller, P. Siddarth,  
845 N.L. Rasgon, et al. Cerebral metabolic and cognitive decline in persons at genetic risk for Alzheimer’s disease.  
846 *Proceedings of the National Academies of Science USA*, 97(11):6037–6042, 2000.
- 847 S. Sonnenburg, G. Rätsch, C. Schäfer, and B. Schölkopf. Large scale multiple kernel learning. *Journal of Machine  
848 Learning Research*, 7:1531–1565, 2006.
- 849 C. Soriano-Mas, J. Pujol, P. Alonso, N. Cardoner, J. M. Menchn, B. J. Harrison, J. Deus, J. Vallejo, and C. Gaser.  
850 Identifying patients with obsessive-compulsive disorder using whole-brain anatomy. *Neuroimage*, 35(3), 2007.
- 851 P. M. Thompson and L.G. Apostolova. Computational anatomical methods as applied to ageing and dementia.  
852 *British Journal of Radiology*, 80(2):78–91, 2007.
- 853 P. M. Thompson, M. S. Mega, R. P. Woods, C. I. Zoumalan, C. J. Lindshield, R. E. Blanton, J. Moussaïl, C. J.  
854 Holmes, J. L. Cummings, and A. W. Toga. Cortical change in Alzheimer’s disease detected with a disease-specific  
855 population-based brain atlas. *Cerebral Cortex*, 11(1):1–16, 2001.
- 856 P. Vemuri, J.L. Gunter, M. L. Senjem, J. L. Whitwell, K. Kantarci, D. S. Knopman, B. F. Boeve, R. C. Petersen,  
857 and C. R. Jack Jr. Alzheimer’s disease diagnosis in individual subjects using structural MR images: validation  
858 studies. *Neuroimage*, 39(3):1186–1197, 2008.
- 859 N. Villain, B. Desgranges, F. Viader, V. de la Sayette, F. Mezenge, B. Landeau, J. C. Baron, F. Eustache,  
860 and G. Chetelat. Relationships between hippocampal atrophy, white matter disruption, and gray matter hy-  
861 pometabolism in Alzheimer’s disease. *Journal of Neuroscience*, 28(24):6174–6181, 2008.
- 862 KB Walhovd, AM Fjell, J. Brewer, LK McEvoy, C. Fennema-Notestine, DJ Hagler Jr, RG Jennings, D. Karow, and  
863 AM Dale. Combining MR Imaging, Positron-Emission Tomography, and CSF Biomarkers in the Diagnosis and  
864 Prognosis of Alzheimer Disease. *American Journal of Neuroradiology*, 31(2):347, 2010.
- 865 J. L. Whitwell, S. A. Przybelski, S. D. Weigand, D. S. Knopman, B. F. Boeve, R. C. Petersen, and C. R. Jack Jr. 3D  
866 maps from multiple MRI illustrate changing atrophy patterns as subjects progress from mild cognitive impairment  
867 to Alzheimer’s disease. *Brain*, 130(7):1777–1786, 2007.
- 868 G. Xu, D. G. McLaren, M. L. Ries, M. E. Fitzgerald, B. B. Bendlin, H. A. Rowley, M. A. Sager, C. Atwood, S. Asthana,  
869 and S. C. Johnson. The influence of parental history of Alzheimer’s disease and apolipoprotein E  $\{\epsilon\}$  4 on the  
870 BOLD signal during recognition memory. *Brain*, 132(2):383, 2009.